

ECONOMETRIC MODELS OF TREATMENT
WITH APPLICATION TO LABOUR MARKET
OUTCOMES

University College London

Jeffrey Rowley

Department of Economics

Submitted towards the degree of Doctor of Philosophy (PhD) in Economics

I, Jeffrey Rowley, confirm that the work presented in this thesis is my own. Where information has been derived from other sources or where work has been completed in collaboration with others, I confirm that this has been indicated in the thesis and that the contribution of each author has been recognised.

DETERMINING ASSIGNMENT AND IDENTIFYING EFFECT

ABSTRACT

This thesis is concerned with the related problems of whether individuals benefit from treatment and determining who should receive treatment—both from an empirical and a theoretical viewpoint—and is broadly divided into two parts (with Chapter A outlining the notation that is used throughout).

The first part is mainly theoretical and examines the problem of how to allocate individuals to treatment in order to maximise a welfare criterion. Chapter B collects several results pertaining to modified Bessel functions of the first kind that are then used in Chapter C to derive a characterisation of the statistical divergence of a von Mises-Fisher distribution and its asymptotic behaviour. Chapter D develops a framework for estimating stochastic—rather than deterministic—assignment rules that are drawn from von Mises-Fisher distributions. This method for estimating assignment rules is motivated by a variational Bayes approximation of the infeasible optimal posterior distribution (i.e., the optimal stochastic assignment rule). The optimal posterior distribution is obtained from the maintained prior via an updating procedure that is based upon an empirical welfare criterion and that reflects the objective of a utilitarian social planner with access to the results of a randomised control trial or otherwise suitable data. The use of stochastic assignment rules has strong theoretical justification, and the proposed method is shown to achieve low regret with high probability. Experimental data from the National Job Training Partnership Act Study is used to illustrate the implementation and performance of this framework.

The second part is mainly empirical and examines the problem of identifying how additional children affect maternal labour supply. Chapter E studies a model of maternal behaviour that allows for a rich set of behaviours by mothers and that is applicable to a range of other empirical problems. Census (and related) data is used to estimate whether maternal labour supply increases or decreases in the presence of additional children.

IMPACT STATEMENT

Several novel results are introduced that can be built upon by or directly used in future academic research, both inside and outside of the field of economics. To give a specific example, von Mises-Fisher distributions have been studied extensively in the field of directional statistics, and statistical divergences appear frequently in the field of machine learning as a means of quantifying information; characterisation of the statistical divergence of a von Mises-Fisher distribution and its asymptotic behaviour is arguably of interest to students in both fields. Indeed, these and many of the other concepts that are discussed

are likely more familiar to students of other fields than to those of economics despite their applicability to many types of economic research. Aside from the intrinsic value of these results, their introduction to the field of economics and the suggestion of alternative approaches to empirical and theoretical research can be of benefit to students of economics—especially to those studying cyclical data, or settings that are homomorphic or isomorphic to spheres.

The tools that are developed are with specific economic environments in mind, and are motivated by particular questions that have a basis in policy design and evaluation. Despite this specificity, however, they are suitable for application to a broad range of other economic environments that have a similar composition. The evidence that these tools provide about labour market outcomes is intended to inform public policy on this subject. How should educational and skills training be directed at the unemployed? How do additional children affect maternal employment at the extensive margin? Such questions are of national and international importance, and the evidence that these tools provide is particularly credible due to the parsimony that both frameworks exhibit. In the case of the first framework, this manifests as a lack of specificity about how education and previous earnings influence wages, and is accompanied by the interpretability of the approximating class of assignment rules and their suitability for application to a broader population. In the case of the second framework, this manifests as a lack of restriction upon how mothers choose their employment status following the birth of additional children. It is important to emphasise that parsimony does not simply lend a tool credibility; it is likewise useful for understanding what information further restrictions can yield and for reconciling any existing evidence that is contradictory.

UCL RESEARCH PAPER DECLARATION FORM: REFERENCING THE DOCTORAL CANDIDATE'S OWN PUBLISHED WORKS

This declaration is supported by a further statement (see the Table of Contents).

1. For a research manuscript that has already been published (if not yet published, please skip to section 2):

- (a) What is the title of the manuscript?
- (b) Please include a link to or doi for the work:
- (c) Where was the work published?
- (d) Who published the work?
- (e) When was the work published?
- (f) List the manuscript's authors in the order they appear on the publication:
- (g) Was the work peer reviewed?
- (h) Have you retained the copyright?
- (i) Was an earlier form of the manuscript uploaded to a preprint server (e.g. medRxiv)? If 'Yes', please give a link or doi
If 'No', please seek permission from the relevant publisher and check the box next to the below statement:

I acknowledge permission of the publisher named under 1d to include in this thesis portions of the publication named as included in 1c.

2. For a research manuscript prepared for publication but that has not yet been published (if already published, please skip to section 3):

- (a) What is the current title of the manuscript?

von Mises-Fisher distributions and their statistical divergence

Stochastic Treatment Choice with Empirical Welfare Updating

(b) Has the manuscript been uploaded to a preprint server 'e.g. medRxiv'?

If 'Yes', please give a link or doi:

arXiv:2202.05192

arXiv:2211.01537

(c) Where is the work intended to be published?

As a technical note in a lower tier statistics journal

As an article in a higher tier economics journal

(d) List the manuscript's authors in the intended authorship order:

Toru Kitagawa and Jeff Rowley

Toru Kitagawa, Hugo Lopez and Jeff Rowley

(e) Stage of publication:

In submission

In submission

3. For multi-authored work, please give a statement of contribution covering all authors (if single-author, please skip to section 4):

Kitagawa (senior author and supervisor)—Contributed to the introduction, and to the main theoretical results (suggested the use of the moment-generating function in place of the general definition of a moment to simplify mathematics). Reviewed.

Rowley—All other work.

Kitagawa (senior author and supervisor)—Contributed to the introduction, to the framework, to the main theoretical results (first theorem and its proof, second theorem and its proof, third theorem and its proof, and first lemma and its proof), and to the implementation. Reviewed.

Lopez—Contributed to the introduction, to the framework, to the main theoretical results (first theorem and its proof, second theorem and its proof, third theorem and its proof, and first lemma and its proof), and to the implementation. Prepared initial draft and edited.

Rowley—Contributed to the main theoretical results (proof of the second theorem, third theorem

and its proof, second lemma and its proof, third lemma and its proof), and to the empirical application. Prepared subsequent draft and edited.

4. In which chapter(s) of your thesis can this material be found?

Chapter C

Chapter D

e-Signatures confirming that the information above is accurate (this form should be co-signed by the supervisor/ senior author unless this is not appropriate, e.g. if the paper was a single-author work):

Candidate:

Date:

Supervisor/Senior Author signature (where appropriate):

Date:

ACKNOWLEDGEMENTS

As I look back now upon my time at University College London, it is with a sense of detachment—and the perspective that affords. Throughout my journey from young undergraduate to the man that I am now, I am fortunate to have learnt from some truly exceptional individuals.

These notably include—in what is certainly a non-exhaustive list: Professors Martin Cripps and Vasiliki Skreta—for the challenge of teaching a changing and stimulating syllabus of microeconomics; Professor Roger Koenker—for insights into computation and programming; Professor Guy Laroque—for discussion of a broad range of topics; Mr Hugo Lopez—for an enjoyable and rewarding collaboration; Professor Thomas Sargent—for curing a natural aversion to macroeconomics; and Dr Michela Tincani.

Professors Andrew Chesher and Toru Kitagawa, in particular, have been central to my academic development, and I would like to extend my most heartfelt thanks to both for the patient and supportive guidance that they have offered me. Their influence on the way that I think about economics cannot be understated—one quote stands out above all others, for the many meanings it holds for me.

If mean independence then why not median independence; and if median independence why not full statistical independence?

I recognise that so many of the opportunities that I—and others—have enjoyed is because of them both: Andrew through his patronage of the discipline via The Centre for Microdata Methods and Practice that he established and continues to advance; and Toru through his patronage of annual intercollegiate and interdisciplinary reading groups. I know that so many hold Andrew in high regard not only for his excellent scholarship but also for his generosity—I count myself amongst them, the warmth with which he greeted me when we ran into each other at the excellent Bristol Econometric Study Group Annual Conference last year just one example—and the very same can be said of Toru.

Indeed, to Toru, I would like to extend a special word of thanks for employing me as a research assistant over the past several years. That employment—and, preceding that, an Economic and Social Research Council grant (Number 1329842; Application of minimally restrictive econometric models)—has been crucial to my undertaking this doctoral programme for the financial support that it has given me, but also for the diverse tasks that Toru has exposed me to and the training that these tasks have provided. Quite simply, I cannot overstate how much I—and this work—owes to Toru, his faith in me, and what I have learnt from him. It has been a great privilege of mine to work for and with Toru, and I wish him and his family every happiness as they build their life together in Rhode Island.

I acknowledge and thank these individuals—and those too that I have not mentioned by name, but whom I have been taught by or worked with—for their influence on my journey.

Along the way I have made some dear and lifelong friends—schoolfriends—with whom I have shared so many precious experiences, and with whom I hope to make many more memories. Some have even ventured with me on this last part of my journey: Dr Rui Costa—especially for a memorable trip to Lisbon; Dr Gavin Kader; and Dr Jeremy McCauley. More great adventures await each of us, I am sure. For Mr Noel Thomas too—whose comradeship was so reassuring—with his retirement.

My family have, of course, been with me throughout my journey, and before too. To them, I have so much that I want to say, and so much to thank them for—I hope that finally seeing my work in print, to know that my struggle is at an end, is reward enough for now.

As my journey ends, I look forward to the next with my beloved fiancée, Debbie, and little Percival and Peregrine. The greatest opportunity that I have had—and ever will—was meeting you.

Dedicated to the memory of my grandmother

Rita Rowley

CONTENTS

A	NOTATION	15
	↔FIGURE A.1: THE TYPOGRAPHY USED IN THE MAIN TEXT	16
	↔FIGURE A.2: THE TYPOGRAPHY USED TO DENOTE RANDOM VARIABLES, SETS AND VECTORS	17
	↔FIGURE A.3: THE TYPOGRAPHY USED TO DENOTE THE SUPPORTS OF STOCHASTIC VARIABLES	17
	↔FIGURE A.4: THE TYPOGRAPHY USED TO DENOTE SOME COMMON SPACES . .	19
	↔FIGURE A.5: THE TYPOGRAPHY USED TO DENOTE FUNCTIONS	20
	↔FIGURE A.6: THE TYPOGRAPHY USED TO DENOTE PARAMETERS	21
	↔FIGURE A.7: THE TYPOGRAPHY USED TO DENOTE INTERMEDIATE OBJECTS .	22
B	MODIFIED BESSEL FUNCTIONS OF THE FIRST KIND	24
1	DEFINITION	24
	↔FIGURE B.1: MODIFIED BESSEL FUNCTIONS OF VARYING ORDER	25
2	RECCURENCE RELATIONS	25
3	LIMITING BEHAVIOUR	26
C	THE f -DIVERGENCE OF A VON MISES-FISHER DISTRIBUTION FROM SOME REFERENCE DIS- TRIBUTIONS	28
	↔FIGURE C.1: VON MISES-FISHER DISTRIBUTIONS ON THE CIRCLE	30
	↔FIGURE C.2: VON MISES-FISHER DISTRIBUTIONS ON THE SPHERE	31
1	THE PROBABILITY DENSITY FUNCTION AND MOMENTS OF THE VON MISES-FISHER FAMILY	32
2	MEASURING THE f -DIVERGENCE OF AN OBTAINED DISTRIBUTION FROM A REFERENCE DISTRIBUTION	34
	2.A RÉNYI DIVERGENCE	36
	2.B χ -SQUARE DISTANCE	37
	2.C SQUARED-HELLINGER DISTANCE	38
	2.D KULLBACK-LEIBLER DIVERGENCE	39
	2.E TOTAL VARIATION DISTANCE	40
	APPENDIX C.1: PROOFS	41
	APPENDIX C.2: HANKEL EXPANSION OF THE CIRCULAR VARIANCE	50
D	STOCHASTIC TREATMENT CHOICE WITH EMPIRICAL WELFARE UPDATING	51
1	FRAMEWORK	56

2	POSTERIOR DISTRIBUTIONS AND STOCHASTIC ASSIGNMENT RULES	60
3	OPTIMAL STOCHASTIC ASSIGNMENT AND CONVERGENCE OF WELFARE	61
3.A	BOUNDING EXPECTED WELFARE RISK	61
3.B	OPTIMAL UPDATING RULE	63
3.C	VARIATIONAL APPROXIMATION OF THE OPTIMAL STOCHASTIC ASSIGNMENT RULE	64
4	IMPLEMENTATION	67
5	EMPIRICAL ILLUSTRATION	68
	⇨FIGURE D.1: VARIATION IN TREATMENT PROPENSITY ACROSS INDIVIDUALS IN	
	THE JTPA STUDY SAMPLE	69
	⇨FIGURE D.2: BEHAVIOUR OF THE OBJECTIVE FUNCTION AT μ^* GIVEN VARIA-	
	TION IN κ	70
	⇨FIGURE D.3: DETERMINISTIC ASSIGNMENT RULES AND EMPIRICAL WELFARE	
	RISK	72
	⇨FIGURE D.4: BEHAVIOUR OF THE OBJECTIVE FUNCTION AT κ^* GIVEN VARIA-	
	TION IN THE MEAN DIRECTION μ	72
	APPENDIX D.1: PROOFS	73
	APPENDIX D.2: ACCOUNTING FOR THE COST OF TREATMENT IN THE JTPA STUDY SAMPLE	80
	⇨FIGURE D.5: VARIATION IN TREATMENT PROPENSITY ACROSS INDIVIDUALS IN	
	THE COST-ADJUSTED JTPA STUDY SAMPLE	81
	⇨FIGURE D.6: BEHAVIOUR OF THE OBJECTIVE FUNCTION AT μ^* GIVEN VARIA-	
	TION IN κ	81
	⇨FIGURE D.7: DETERMINISTIC ASSIGNMENT RULES AND EMPIRICAL WELFARE	
	RISK	82
	⇨FIGURE D.8: BEHAVIOUR OF THE OBJECTIVE FUNCTION AT κ^a GIVEN VARIA-	
	TION IN μ	82
	⇨FIGURE D.9: DISTRIBUTION OF TREATMENT PROPENSITY ACROSS INDIVIDU-	
	ALS (ADJUSTED) AT $\{\kappa^a, \mu^a\}$	83
	⇨FIGURE D.10: DISTRIBUTION OF TREATMENT PROPENSITY ACROSS INDIVID-	
	UALS (RAW) AT $\{\kappa^*, \mu^*\}$	83
	APPENDIX D.3: NUMERICAL SIMULATIONS	84
	⇨FIGURE D.11: BEHAVIOUR OF THE OBJECTIVE FUNCTION AT μ^* GIVEN VARIA-	
	TION IN κ	87
	⇨FIGURE D.12: (EXPERIMENT 1) VARIATION IN TREATMENT PROPENSITY	
	ACROSS INDIVIDUALS	88

↔	FIGURE D.13: (EXPERIMENTS 1–3) BEHAVIOUR OF THE OBJECTIVE FUNCTION AT μ^* GIVEN VARIATION IN κ	89
↔	FIGURE D.14: (EXPERIMENT 4) VARIATION IN TREATMENT PROPENSITY ACROSS INDIVIDUALS	91
↔	FIGURE D.15: (EXPERIMENTS 4–6) BEHAVIOUR OF THE OBJECTIVE FUNCTION AT μ^* GIVEN VARIATION IN κ	92
↔	FIGURE D.16: (EXPERIMENTS 7–8) VARIATION IN TREATMENT PROPENSITY ACROSS INDIVIDUALS	95
↔	FIGURE D.17: (EXPERIMENTS 7–8) BEHAVIOUR OF THE OBJECTIVE FUNCTION AT μ^* GIVEN VARIATION IN κ	96
↔	FIGURE D.18: (EXPERIMENTS 9–10) VARIATION IN TREATMENT PROPENSITY ACROSS INDIVIDUALS	97
↔	FIGURE D.19: (EXPERIMENTS 9–10) BEHAVIOUR OF THE OBJECTIVE FUNC- TION AT μ^* GIVEN VARIATION IN κ	98
E	THE EFFECT OF ADDITIONAL CHILDREN ON MATERNAL LABOUR SUPPLY	100
1	FRAMEWORK	104
1.A	THE MEANING AND CREDIBILITY OF THE RESTRICTIONS	106
2	EQUIVALENCE CLASSES AND LATENT TYPES	108
3	IDENTIFICATION	110
4	DATA AND SAMPLE COMPOSITION	112
4.A	VARIABLE DEFINITION	113
↔	FIGURE E.1: FERTILITY AND LABOUR SUPPLY MEASURES	114
↔	FIGURE E.2: FERTILITY MEASURES, CHILD GENDER, AND LABOUR MARKET OUTCOMES	115
↔	FIGURE E.3: PARENTAL CHOICE OVER FAMILY SIZE GIVEN SEX OF CHILDREN	118
5	LOCAL ESTIMATION	118
↔	FIGURE E.4: DIFFERENCES IN MEANS FOR DEMOGRAPHIC VARIABLES BY IN- STRUMENTAL VARIABLE DEFINITION	119
↔	FIGURE E.5: ESTIMATES OF LOCAL TREATMENT EFFECTS	120
6	ESTIMATION	120
	APPENDIX E.1: THE SUITABILITY OF AN AUXILIARY MODEL	121
↔	FIGURE E.6: ESTIMATES OF GLOBAL TREATMENT EFFECTS WITHOUT COVARI- ATES	122
↔	FIGURE E.7: ESTIMATES OF GLOBAL TREATMENT EFFECTS WITH COVARIATES	123

APPENDIX E.2: DISCUSSION OF THE GENERALISED BALKE-PEARL BOUNDS	125
↔FIGURE E.8: LATENT TYPES AND THE IDENTIFICATION PROBLEM	126
APPENDIX E.3: DERIVING THE GENERALISED BALKE-PEARL BOUNDS	127
F BIBLIOGRAPHY	135
G DIVISION OF LABOUR	150

—| NOTATION |—

CONVENTIONS AND TYPOGRAPHY

This chapter constitutes a record of the various mathematical and typographical conventions that I have adopted—a sort of standards manual if you will. I include some additional information that is otherwise taken to be implicit and underline when I feel that additional emphasis is warranted (something that I continue to do throughout this thesis). Some of these works are full or partial reproductions of collaborations with with co-authors. Where this is the case, I include only that material that I have contributed towards or else that is necessary for the reader’s comprehension—omitting, for instance, any related proofs—and use the third-person singular.

I adopt one and a half spacing throughout. All margins are set to $3/4$ inch except the binding (side) margin, which is set to $3/2$ inch.

The main text is typeset in Computer Modern regular serif font (Figure A.1), or in Computer Modern italicised serif font (Figure A.1) whenever this is appropriate. Figures and tables are centred, framed and titled, with the title appearing above each object and typeset in Computer Modern capitalised serif font (Figure A.1). A caption is included below a figure or a table when I feel that further description is necessary, and typeset in Computer Modern italicised serif font. All external references are typeset in Computer Modern capitalised serif typeface, whilst internal references are regularly typeset. Paragraphs are marked using vertical line spacing.

Mathematical expressions are typeset in one of two fonts and several typographical effects depending upon what they are intended to represent, and are either presented inline (if they comprise a simple equality or inequality relation, or are the product of no more than two objects) or as equations. Equations are left-aligned and labelled, and are preceded by the conditions under which they hold. For instance, an equation might hold for all or for some values of an argument. Where these conditions are omitted then equations should be taken to hold for all values of their arguments. Equations are not punctuated as a matter of style.

FIGURE A.1
The typography used in the main text

Computer Modern regular serif font

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z
a b c d e f g h i j k l m n o p q r s t u v w x y z

Computer Modern italicised serif font

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z
a b c d e f g h i j k l m n o p q r s t u v w x y z

Computer Modern capitalised serif font

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z
A B C D E F G H I J K *l* M N O P Q R S T U V W X Y Z

Underpinning my analysis and associated exposition is The Probability Approach in Econometrics (HAAVELMO, 1944), and I assume throughout that there exists a (standard) probability space equipped with the Borel algebra, and upon which are defined a multitude of stochastic variables that stand-in for real economic variates.

I denote stochastic variables using the Latin alphabet, distinguishing single-valued stochastic variables from their multiple-valued counterparts (Figure A.2). I refer to single-valued stochastic variables as random variables, which I denote using Computer Modern italicised serif font; and I refer to multiple-valued stochastic variables as random sets or random vectors, which I denote using Computer Modern bold italicised serif font. I reserve the upper-case alphabet for stochastic variables, whilst I reserve the lower-case alphabet for their realisations (actual or hypothetical). Where I refer to more than one realisation of the same stochastic variable, I distinguish these two realisations using a combination of accents and subscripted characters from the Latin alphabet. So as to minimise overlap or possible confusion, I avoid using those characters from the alphabet that are difficult to distinguish from Arabic numerals or that I use extensively for other purposes,¹ such as to represent mathematical objects and operators that I rely upon.

I denote the collection of values that a stochastic variable can take using the Latin alphabet combined with Fraktur font (Figure A.3). By support, I mean one of two things, which I do not distinguish between: the collection of values that a stochastic variable can take; and the collection of values that a stochastic variable does take (i.e., the collection of values that occur with positive probability). I reserve the upper-case alphabet for the supports of stochastic variables, whilst I reserve the lower-case alphabet for their cardinalities. Where I am concerned with the support of a stochastic variable conditional on

¹ I present the full alphabet in Figure A.2 in any case, for completeness.

but also include some special spaces such as the spherical manifold, the unit continuum and the unit simplex. These special spaces warrant some additional description. The spherical manifold is defined for higher dimensional Euclidean spaces as the collection of unit vectors; the unit continuum is defined for higher dimensional Euclidean spaces as the collection of vectors whose elements are all between zero and one; and the unit simplex is defined for higher dimensional Euclidean spaces as the collection of vectors of non-negative elements whose elements sum to one. I otherwise refer to the spherical manifold as the hypersphere, to the unit continuum as the hypercube and to the unit simplex as the simplex. I subscript spaces that can contain zero by an asterisk to render the space exclusive of zero. I superscript the space by infinity to render the power set of (or Borel sets over) the space; I superscript the space by a finite integer (or by i , as a placeholder for an unspecified finite integer) to render the product space; and I use a colon separator (plus symbol preceding infinity) to render the power set over the product space. To summarise, I use the syntax

Space^∞ = the power set of the space

Space^i = the product set of the space A.2

$\text{Space}_*^{i:\infty}$ = the power set of the product set of the space exclusive of zero

which is a convention that I also adopt when considering the support of a stochastic variable (in place of a space). I am not specific as to whether these spaces or their subsets are closed or open or neither unless providing this detail is useful. The notation that I use for the space that a random variable inhabits is consistent with the notation that I use for the space that a function inhabits.

To distinguish cases where a random set or random vector is such that several of their elements can take the same value from cases where each of their elements is distinct, I use \in (the in symbol) and \subset (the subset symbol), respectively. To avoid any confusion, particularly with respect to closedness and openness, I use inequality relations to define intervals rather than brackets and parentheses; I use brackets to group objects (in advance of applying an operation to them say) and parentheses to enclose function arguments.

Where I am concerned with the particular elements of a set or vector, I use subscripted Arabic numerals to index and reference individual elements. I denote the concatenation of elements of a set or of a vector using braces, separating individual elements by a comma or a semi-colon depending upon which is clearer. Both sets and vectors are, in essence, lists of varying dimension or length, and the notation that I use does not distinguish between them in view of this. Often a set or vector comprises elements

FIGURE A.4
The typography used to denote some common spaces

General spaces	Special spaces
\mathbb{C} \mathbb{N} \mathbb{R} \mathbb{Z}	\mathbb{D} \mathbb{S} \mathbb{U}
<i>Complex numbers</i> <i>Natural numbers</i> <i>Real numbers</i> <i>Integer numbers</i>	<i>Unit simplex</i> <i>Spherical manifold</i> <i>Unit continuum</i>

that have a common definition or form, for which I use the syntax

$$\text{List} = \{\text{List}_i : \text{Definition of the } i^{\text{th}} \text{ element (and all other elements) of the list}\} \tag{A.3}$$

additionally specifying the cardinality of these objects or other information when I feel that this is warranted. To transpose a list, I use \top (the transpose symbol), which I write in superscript.

Some of the characters of the Latin alphabet already have a well-established meaning in the body of existing economic work. For instance, the characters Y and T are typically interpreted as observable endogenous economic variates (response and treatment, respectively); the characters Z and X as observable exogenous economic variates (instrument and covariates, respectively); and the character U as a latent exogenous economic variate (heterogeneity). I maintain this convention for the most part, whether as random variables as is written above or as random sets and random vectors.

What distinguishes endogenous economic variates from their exogenous counterparts is how they are determined. I say that an economic variate is endogenous if it is determined by other economic variates in the overarching system, and is exogenous if it is given (i.e., it is determined outside of the overarching system, or is determined by an economic agent other than the one whose choice is of principal focus). The relationships between the endogenous and exogenous components of the economic environment are described by structural functions, for which I adopt the same conventions as I do for functions.

I denote functions (and operators) using characters from the Greek or Latin alphabets combined with either Computer Modern italicised serif font or Computer Modern regular serif font, omitting those characters from the Greek alphabet that overlap with or appear similar to characters from the Latin alphabet (Figure A.5). The reader is undoubtedly familiar with many of the functions that are associated with these characters, but some may be unfamiliar. I enclose the arguments of functions inside parentheses. In cases where several functions share a purpose, I denote these functions using the same character; to

FIGURE A.5
The typography used to denote functions

Greek alphabet										
Γ	Δ	Θ	Λ	Ξ	Π	Σ	Υ	Φ	Ψ	Ω
<i>Gamma</i>	<i>Delta</i>	<i>Theta</i>	<i>Lambda</i>	<i>Xi</i>	<i>Pi</i>	<i>Sigma</i>	<i>Upsilon</i>	<i>Phi</i>	<i>Psi</i>	<i>Omega</i>

distinguish between each function, I employ a combination of accents and characters from the Latin alphabet. For example, in the classical model of endogenous choice described by the equations

$$\begin{aligned}
 y &= h_y(t, x, u) \\
 t &= h_t(z, x, u)
 \end{aligned}
 \tag{A.4}$$

I write h to refer to h_y and h_t together. This example also serves to illustrate how I refer to functions in the main text. Whilst I typeset non-Latin characters exclusively in regular font, I typeset Latin characters in a mixture of italicised and regular font; as a general rule, I use regular font to denote common (perhaps more accurately described as standard, since some of the objects that I work with are uncommon) constants, functions and operators, otherwise using italicised font. Particular examples include the characters D (the differential operator), E (the elementary vector, the expectation operator, or the exponential constant), I (the identity matrix or vector, or the modified Bessel function of the first kind), M (the annihilation matrix), O (the order operator), and P (the projection matrix), to list but a few examples.

Often, the functions that I define are characterised or multiplied by one or several parameters. Where these parameters are of finite dimension, I term the associated function parametric; where these parameters are of infinite dimension, I term the associated function non-parametric; and where some of these parameters are of finite dimension and others of infinite dimension, I term the associated function semi-parametric. I denote parameters using the Greek alphabet combined with Computer Modern italicised serif font (Figure A.6), omitting those characters that overlap with or appear similar to characters in the Latin alphabet or that are mathematical constants that I rely upon. To distinguish the parameters of a function from its other inputs, I use a semi-colon separator.

In some limited cases, the functions that I define rely upon other functions. That is, they are functions of functions—otherwise known as functionals. To distinguish the functions of a functional from its other

FIGURE A.6
The typography used to denote parameters

α	β	γ	δ	ε	ζ	η	θ	κ	λ	μ	ξ	σ	ϕ	χ	ψ	ω
Alpha	Beta	Gamma	Delta	Epsilon	Zeta	Eta	Theta	Kappa	Lambda	Mu	Xi	Sigma	Phi	Chi	Psi	Omega

inputs and parameters, I use a semi-colon separator. This syntax can be summarised as

$$\text{Output} = \text{Function}(\text{Input}; \text{Parameter}) \tag{A.5}$$

$$\text{Output} = \text{Functional}(\text{Function}; \underline{\text{Functional input}}; \underline{\text{Functional parameter}})$$

with non-applicable arguments simply omitted as and when this is required. I otherwise do not treat functionals differently from other types of functions.

I emphasise that functionals are not the same thing as composite functions. Where I sequentially apply more than two functions and no detail is omitted by doing so (i.e., a later function does not require additional arguments to be specified beyond the output of the preceding function), I use \circ (the composition symbol). This syntax can be summarised as

$$\text{Output} = \text{Second function}(\text{First function}(\text{Input})) \tag{A.6}$$

$$\text{Output} = \text{Third function} \circ \text{Second function}(\text{First function}(\text{Input}))$$

with the composition symbol replaced by more standard parentheses if the third function is reliant on another input or on a parameter. The composition symbol is not to be confused with \cdot (the dot symbol) that I use for multiplication. I also use the composition symbol to make notation more concise when several functions with the same input are added, divided, multiplied or subtracted.

The functions that I define are, for the most part, regular, in that they map from one or several single-valued inputs (whether deterministic or random) to a single-valued output. I do, however, allow for a function's generalisation to set- or vector-valued inputs that, accordingly, map to set- or vector-valued outputs. For example, returning to the classical model of endogenous choice described by Equation A.4, I require that h_y satisfies such equality relations as

$$h_y(t, x, \mathbf{u}) = \{\mathbf{y}_i : \mathbf{y}_i = h_y(t, x, \mathbf{u}_i)\} \tag{A.7}$$

say, so that h_y is well-defined irrespective of the dimension of its inputs. I require that functions have well-defined inverses; with the exception of the inverse operator, I write operations applied to functions

FIGURE A.7
The typography used to denote intermediate objects

<u>Cyrillic alphabet</u>									
Ђ	Ж	Љ	Њ	Ю	Я	Æ	Ђ	Ѓ	Ќ
<i>Đ_e</i>	<i>Ž_e</i>	<i>Ĺ_e</i>	<i>Ń_e</i>	<i>Yu</i>	<i>Ya</i>	<i>Æ_e</i>	<i>Đ_e</i>	<i>Ǻ_e</i>	<i>Yus</i>

in superscript following both the function and the parentheses enclosing its arguments.

Differentiation and integration are two operations that I apply to functions but do not use superscripted notation for; rather, I write the differential and integral operators in full, making explicit the variable of differentiation or integration.

To characterise the behaviour of functions as the principal argument takes extreme values I use either the big-O operator (from big-O notation, otherwise referred to as Bachmann-Landau notation in LATTIMORE and SZEPESVÁRI, 2020) or \simeq (the approximately equal symbol). Specifically, I use the big-O operator to characterise the behaviour of functions as their principal argument becomes large (i.e., approaches infinity) and the approximately equal symbol to characterise the behaviour of functions as their principal argument becomes small (i.e., approaches zero). In each case, I express the function in terms of another function. This syntax can be summarised as

$$\begin{aligned} \text{Function (Principal argument)} &= O(\text{Known function (Principal argument)}) \\ \text{Function (Principal argument)} &\simeq \text{Known function (Principal argument)} \end{aligned} \tag{A.8}$$

which should be taken to mean that the function on the left-hand side grows or decays at a rate that is proportional to the known function on the right-hand side, or is else is approximately equal to the known function on the right-hand side; the growth or decay rate of the known function on the right-hand side is then what I refer to as the rate of the function on the left-hand side.

To characterise the neighbourhood that a function occupies for particular values of its inputs, it is necessary to specify constants. This is also the case when talking about how likely something is to occur. I follow convention in reserving the characters Delta and Epsilon from the Greek alphabet (Figure A.6) to denote positive and infinitesimal or specified constants, as is typical elsewhere. I capture the probability with which an event occurs, which is typically the probability with which a stochastic variable takes one or several possible values, using the probability operator. Where this probability is

conditional, I use the conditional probability operator. I summarise this syntax as

$$\Pr(\text{Stochastic variable} = \text{Realisation of stochastic variable})$$

A.9

$$\Pr(\text{Stochastic variable} = \text{Realisation of stochastic variable} | \text{Conditioning event})$$

depending upon whether it is an unconditional or conditional probability that is of interest, and where I write the conditioning event simply as the realisation of a stochastic variable (rather than as the full event that the stochastic variable takes that realisation) if this is sufficient to make the full event clear.

Since each of the following chapters are self-contained and standalone (with the exception of some of the appendices, that precede or follow certain chapters), I recycle notation between them. The reader should not, therefore, carry-over notation and terminology from one chapter to another. In several places, it is necessary for me to define functions that play a very limited role in the overall analysis; I denote such intermediate functions using characters from the Slavic or non-Slavic Cyrillic alphabets combined with Computer Modern regular serif font, omitting those characters from the Cyrillic alphabet that overlap with or appear similar to characters from the Latin alphabet (Figure A.7).

—| CHAPTER B |—

MODIFIED BESSEL FUNCTIONS OF THE FIRST KIND

Hereafter, where I refer to the modified Bessel function, I intend this to mean the modified Bessel function of the first kind.

—| SECTION 1 |—

DEFINITION

The modified Bessel function is defined in NIST (2021, §10.25 and §10.32) and, its precursor, ABRAMOWITZ and STEGUN (1964, §9.6), and I refer the reader to those references for further (detailed) information about this function. The modified Bessel function is defined, for all $a > 0$, as

$$I_m(a) \doteq (a/2)^m \cdot \sum_{k=0}^{\infty} \frac{(a/2)^{2k}}{k! \cdot \Gamma(m+k+1)} = \frac{(a/2)^m}{\sqrt{\pi} \cdot \Gamma(m+1/2)} \cdot \int_0^\pi \exp(\pm a \cdot \cos(b)) \cdot \sin(b)^{2m} \cdot db \quad \text{B.1}$$

or as

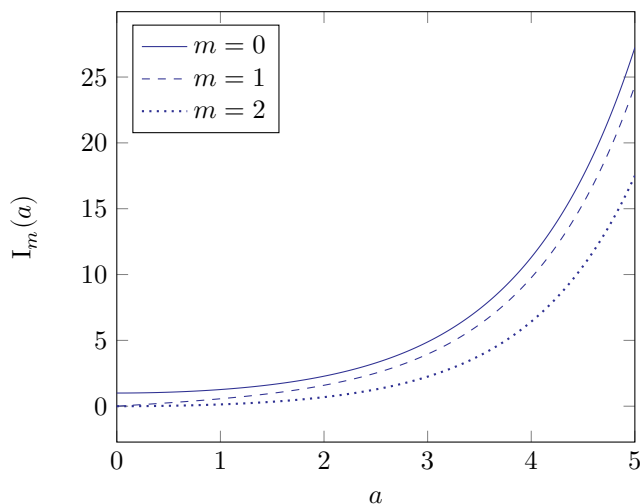
$$I_m(a) \doteq \frac{1}{\pi} \cdot \int_0^\pi \exp(a \cdot \cos(b)) \cdot \cos(m \cdot b) \cdot db - \frac{\sin(m \cdot \pi)}{\pi} \cdot \int_0^\infty \exp(-m \cdot c - a \cdot \cosh(c)) \cdot dc \quad \text{B.2}$$

where m is said to be the order and a is said to be the argument. Various other definitions of the modified Bessel function exist, including as the purely imaginary solution to

$$a^2 \cdot \frac{d^2 f(a)}{da^2} + a \cdot \frac{df(a)}{da} - (a^2 + m^2) \cdot f(a) = 0 \quad \text{B.3}$$

which is a modification of Bessel's equation. This modification relates to the final term on the left-hand side of Equation B.3; Bessel's equation is generally written with the final term on the left-hand side entering additively. This is what distinguishes a Bessel function (that admits real solutions, and is a generalisation of a sine wave) from a modified Bessel function (that admits imaginary solutions, and is an increasing function). Bessel's equation and the modified Bessel function relate to Laplace's equation and harmonic functions (that describe the propagation of a wave along a taut string). The modified

FIGURE B.1
Modified Bessel functions of varying order



Modified Bessel functions of increasing order are plotted for several values of their argument.

Bessel function is a type of generalised hypergeometric function.

The modified Bessel function is plotted in Figure B.1 for various integer orders, although the function is equally well-defined for non-integer orders too, including negative ones.

—| SECTION 2 |—

RECCURENCE RELATIONS

The modified Bessel function satisfies the recurrence relation (NIST, 2021, §10.29)

$$I_m(a) = \frac{a}{2m} \cdot (I_{m-1}(a) - I_{m+1}(a)) \tag{B.4}$$

that arises from

$$I_{m-1}(a) - \frac{m}{a} \cdot I_m(a) = \frac{d}{da} I_m(a) = I_{m+1}(a) + \frac{m}{a} \cdot I_m(a) \tag{B.5}$$

such that the modified Bessel function of a particular order can be computed recursively from the modified Bessel function of a lower order (indeed, this is how many computer programmes calculate the modified Bessel function).

AMOS (1974) introduces several further results that are not only useful for the characterisation of the modified Bessel function, but also for its computation. Adapting the notation of that paper to suit,

AMOS (1974) shows that, for all $\kappa_q \geq \kappa_r > 0$,

$$\underline{\mathbb{B}}(\kappa_q, \kappa_r, m) \leq \ln(I_m(\kappa_q)) \leq \overline{\mathbb{B}}(\kappa_q, \kappa_r, m) \quad \text{B.6}$$

where

$$\underline{\mathbb{B}}(\kappa_q, \kappa_r, m) \doteq \ln(I_m(\kappa_r)) + m \cdot \ln\left(\frac{\kappa_q}{\kappa_r}\right) + \underline{a}_m \cdot \ln\left(\frac{\underline{a}_m + \sqrt{\kappa_r^2 + \underline{a}_m^2}}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}}\right) + \frac{\kappa_q^2 - \kappa_r^2}{\sqrt{\kappa_q^2 + \underline{a}_m^2} + \sqrt{\kappa_r^2 + \underline{a}_m^2}} \quad \text{B.7}$$

and

$$\overline{\mathbb{B}}(\kappa_q, \kappa_r, m) \doteq \ln(I_m(\kappa_r)) + m \cdot \ln\left(\frac{\kappa_q}{\kappa_r}\right) + \underline{a}_m \cdot \ln\left(\frac{\underline{a}_m + \sqrt{\kappa_r^2 + \underline{a}_m^2}}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}}\right) + \frac{\kappa_q^2 - \kappa_r^2}{\sqrt{\kappa_q^2 + \underline{a}_m^2} + \sqrt{\kappa_r^2 + \underline{a}_m^2}} \quad \text{B.8}$$

given $\underline{a}_m \doteq m + 1/2$ and $\bar{a}_m \doteq m + 3/2$. I note that Equation B.6 is reversed when the arguments instead satisfy $\kappa_r > \kappa_q$. In the special case where $\kappa_r \rightarrow 0$, AMOS (1974) shows that Equations B.7 and B.8 reduce to

$$\underline{\mathbb{B}}(\kappa_q, 0, m) = \frac{1}{2} \cdot \ln\left(\frac{2}{\kappa_q}\right) - \ln(\Gamma(m+1)) + \underline{a}_m \cdot \ln\left(\frac{\kappa_q \cdot (\underline{a}_m + \bar{a}_m)/2}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}}\right) + \frac{\kappa_q^2}{\bar{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}} \quad \text{B.9}$$

and

$$\overline{\mathbb{B}}(\kappa_q, 0, m) = \frac{1}{2} \cdot \ln\left(\frac{2}{\kappa_q}\right) - \ln(\Gamma(m+1)) + \underline{a}_m \cdot \ln\left(\frac{\kappa_q \cdot (\underline{a}_m + \underline{a}_m)/2}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}}\right) + \frac{\kappa_q^2}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}} \quad \text{B.10}$$

respectively. It is self-evident that Equation B.10 exceeds Equation B.9 as is required for Equation B.6 to be non-empty. An implication of Equation B.6 is that

$$\frac{\kappa_q}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}} \leq \frac{I_{m+1}(\kappa)}{I_m(\kappa)} \leq \frac{\kappa_q}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}} \quad \text{B.11}$$

which defines a subset of the unit interval as is evident from Equation B.4.

—| SECTION 3 |—

LIMITING BEHAVIOUR

The modified Bessel function satisfies the limiting behaviour (NIST, 2021, §10.30)

$$\frac{I_m(a)}{a^m} \simeq \frac{1}{2^m \cdot \Gamma(v+1)} \quad \text{B.12}$$

which is proportional to the surface area of the hypersphere. Conversely, the modified Bessel function admits, for all $n \in \mathbb{N}_*$, the Poincaré asymptotic expansion (NIST, 2021, §2.1)

$$I_m(a) = \frac{\exp(a)}{\sqrt{2\pi \cdot a}} \cdot \left(\sum_{j=0}^{n-1} (-1)^j \cdot \frac{\text{Pochhammer}_j(m)}{a^j} + O\left(\frac{1}{a^n}\right) \right) \tag{B.13}$$

which is derived from Hankel's expansion (NIST, 2021, §10.17 and §10.40) and which is defined in terms of a factorial sequence (NIST, 2021, §5.2—Pochhammer's symbol) that is polynomial in its argument. It is of no surprise that the modified Bessel function grows at an exponential rate given the appearance of Figure B.1, but this is easily verified from Equation B.13.

THE f -DIVERGENCE OF A VON MISES-FISHER DISTRIBUTION FROM SOME REFERENCE
DISTRIBUTIONS

The von Mises-Fisher family of distributions is well-known in the field of directional statistics¹ but is foreign to economics and, as such, warrants some introduction. Also known as the Langevin family (WATSON, 1984), the von Mises-Fisher family recognises those two titans of statistics, Sir Ronald Fisher and Richard von Mises, for their seminal contributions in considering Gaussianity on the circle (VON MISES, 1918) and on the sphere (R. A. FISHER, 1953). Subsequent work has generalised the von Mises-Fisher family to \mathbb{S}^i , and has led to the definition of other related parametric distributions such as the Bingham family (BINGHAM, 1974) and the Fisher-Bingham or Kent family (KENT, 1982).

A von Mises-Fisher distribution assigns probability mass to the surface of the unit ball—the hypersphere. As such, the von Mises-Fisher family is relevant to situations where the researcher is interested in either the sampling of directional vectors—i.e., vectors of unit length—or in the clustering of some phenomenon on a circular object, such as occurs if data is periodic. Applications range from the study of sea turtle navigation (HILLEN et al., 2017), to the study of perihelia of long-tailed comets (MARDIA, 1975) and near-earth objects (SEI et al., 2013), as well as to the study of patient arrival data (MARDIA, 1975). SABELFELD (2018) even links the von Mises-Fisher family to the solving of high-dimensional diffusion-advection-reaction equations. The von Mises-Fisher family is a two-parameter family, summarised by a concentration parameter (or, simply, concentration), which we denote by $\kappa > 0$, and a mean direction, which we denote by $\mu \in \mathbb{S}^i$.

The main contribution of this chapter is to provide analytical expressions for the f -divergence of a von Mises-Fisher distribution from two relevant reference distributions given several common choices of function. We study the broad class of Rényi divergence of simple order as well as several other measures of (statistical) divergence that relate to the Rényi class—the χ -square distance, the squared-Hellinger distance and the Kullback-Leibler divergence. Each is, of course, a measure of the difference between

¹ See MARDIA and JUPP (2009) for a summary of important results in the field of directional statistics.

two probability distributions. Several well-known inequalities relate these measures to the total variation distance (BRETAGNOLLE and HUBER, 1978; PINSKER, 1964), which is often of interest. The reference distributions that we specify are another, distinct, von Mises-Fisher distribution on \mathbb{S}^i and the uniform distribution on the hypersphere. We are unaware of such expressions being available elsewhere. Alongside these expressions, we characterise how the various measures of divergence that we consider change as a von Mises-Fisher distribution becomes increasingly concentrated and degenerate. In particular, we clarify the leading order terms of asymptotic expansions,² which relies on results in AMOS (1974). These asymptotic expansions offer analytically tractable polynomial approximations of each of the measures of divergence that we consider in terms of the concentration parameter, with these approximations accurate when this parameter takes large values.

Obtaining analytical expressions of statistical divergence is useful for implementing minimum distance-type or penalised estimation methods, and also for characterising the statistical performance of these procedures. See, for instance, KITAGAWA et al. (2022b), which builds upon the analytical expression of the Kullback-Leibler divergence that is derived in this chapter to estimate the randomised treatment assignment rule that minimises a penalised empirical welfare criterion. Polynomial approximation of statistical divergences for parameter values corresponding to high concentration are useful for characterising the convergence behaviour of numerous objects including the concentration rate of the posterior distribution and of the PAC-Bayes regret bounds that are formulated in KITAGAWA et al. (2022b), and that are exploited in that paper.

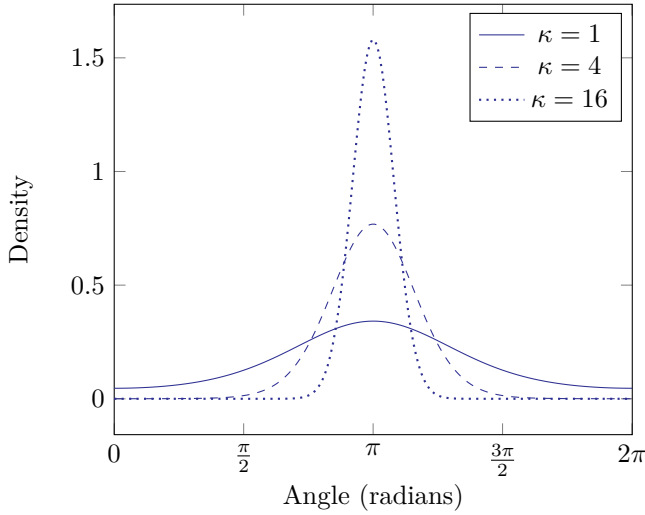
Analytical expressions for the moments and other known distributional features of the von Mises-Fisher family are, to varying extents, available elsewhere. These include statement of the first two moments of a von Mises-Fisher distribution (available in or adaptable from MARDIA and JUPP, 2009, §9.3 and §9.6, respectively) and its associated Fisher information matrix (HORNIK and GRÜN, 2013). DHILLON and SRA (2003), HILLEN et al. (2017), and HORNIK and GRÜN (2013) demonstrate three distinct approaches to obtaining expressions for these moments. These are integration by substitution following transformation to spherical coordinates, application of the divergence theorem,³ and differentiation of the moment-generating function, respectively. We use several of the results in these papers directly.⁴ In particular, knowledge of the first moment is essential to characterising the divergence of a von Mises-Fisher distribution from our chosen reference distributions.

2 In the appendices, we demonstrate how Hankel expansions can be performed to complement results given in KITAGAWA et al. (2022b).

3 The divergence theorem relates the area of a surface integral to a volume integral.

4 We otherwise establish these results in Section 2.e for the benefit of the reader, should they be unfamiliar with spherical distributions and their associated concepts. Our use of the moment-generating function is very much a case of horses for courses, to use the British idiom: integration by substitution and the divergence theorem can also be used, but their application is more involved in this setting.

FIGURE C.1
von Mises-Fisher distributions on the circle



The density function of a von Mises-Fisher distribution on the circle with mean direction $(-1, 0)$ —i.e., a polar orientation with reference angle equal to π radians—for several values of the concentration parameter. A spherical coordinate system is used.

Directional objects are common in economics, and the von Mises-Fisher family of distributions is relevant to many environments and several methods. For instance, consider the canonical binary choice model with latent random utility. The rational choice of the individual, which we denote by $t \in \{0, 1\}$, is determined according to the linear index equation

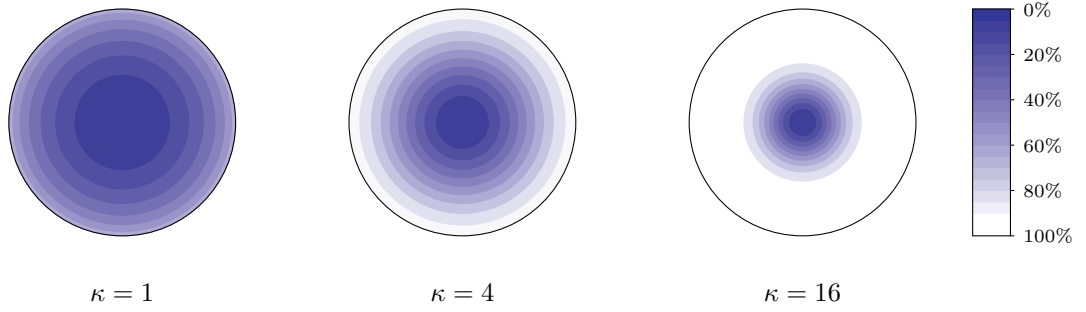
$$t = 1(\mathbf{x}^\top \boldsymbol{\gamma} - u \geq 0) \tag{C.1}$$

where $\mathbf{x} \in \mathbb{R}^i$ and $u \in \mathbb{R}^i$ denote the individual’s observable characteristics (including an intercept) and latent heterogeneity, respectively. The conditional zero-median restriction inherent in the semiparametric maximum score approach of MANSKI (1975, 1985) does not identify the scale of the utility coefficients $\boldsymbol{\gamma} \in \mathbb{R}^i$. It is common to normalise the parameter space of $\boldsymbol{\gamma}$ to the collection of vectors satisfying $\|\boldsymbol{\gamma}\|_2 = 1$ —i.e., to the hypersphere. Similarly, in the context of statistical treatment choice (MANSKI, 2004a), KITAGAWA and TETENOV (2018a) considers individualised treatment assignment rules based upon a linear index,

$$t = 1(\mathbf{x}^\top \boldsymbol{\gamma} \geq 0) \tag{C.2}$$

where $\mathbf{x}^\top \boldsymbol{\gamma} \in \mathbb{R}$ is an eligibility score that aggregates the individual’s observable characteristics and determines whether she should be assigned to treatment—i.e., if $\mathbf{x} \in \mathbb{R}^i$ maps to $t = 1$ —or to non-treatment—i.e., if $\mathbf{x} \in \mathbb{R}^i$ maps to $t = 0$. Such assignment rules are invariant to multiplication of the eligibility score by a positive constant and can be uniquely indexed by a parameter vector on the

FIGURE C.2
von Mises-Fisher distributions on the sphere



Orthogonal projection of the sphere, oriented to the mean direction, for several values of the concentration parameter. Each contour describes a region in which the von Mises-Fisher distribution assigns 10% mass, with contours distinguished by their shading. As the value of the concentration parameter is increased, more mass is assigned to the vicinity of the pole.

hypersphere.

Optimising maximum score or an empirical welfare criterion is difficult, however, and complicates estimation of and inference on γ . This motivates a quasi-Bayesian approach as considered in CHERNOZHUKOV and HONG (2003a). The von Mises-Fisher family offers a parsimonious and convenient prior specification over γ within the quasi-Bayesian framework, with prior elicitation facilitated by knowledge of the moments of the distribution. In a related but different context, PAC-Bayesian analysis, which is widely studied in machine learning (ALQUIER et al., 2016; CATONI, 2007; GERMAIN et al., 2009; MCALLESTER, 2003, to name but a few relevant papers), considers exponentiated negative empirical risk as a quasi-likelihood, and forms a posterior distribution over prediction rules. The von Mises-Fisher family is then not only useful as a specified prior over directional parameters, but can also be used to approximate a posterior distribution over γ in the linear classification rule or linear index treatment assignment rule settings (KITAGAWA et al., 2022b).

Another context where the parameter space is isomorphic to the hypersphere is the class of underidentifying linear simultaneous equation models in which the imposed model restrictions identify structural parameters up to sets of orthonormal transformations. See, for instance, ARIAS et al. (2018), GIACOMINI and KITAGAWA (2021a), and UHLIG (2005) for a class of set-identified structural vector autoregressions in which the identified set of an impulse-response is spanned by the class of orthonormal matrices. The isomorphism of the hypersphere and the orthogonal group suggests that a spherical distribution—and a von Mises-Fisher distribution in particular—can be used as a prior distribution for the non-identified orthonormal matrices. When combined with prior elicitation of the reduced-form parameters, the moments of a von Mises-Fisher distribution can be used to translate a belief about the structural parameters into a prior distribution over the non-identified orthonormal matrices. Like the Gaussian family of distribu-

tions on the hyperplane from which it can be derived, the von Mises-Fisher family is highly restrictive⁵ but, nonetheless, forms an interesting baseline case to study.

We are unaware of any paper that characterises the f -divergence of a von Mises-Fisher distribution as we do. For instance, DIETHE (2015) similarly studies the Kullback-Leibler divergence of von Mises-Fisher distributions. Whereas we provide exact analytical expressions, DIETHE (2015) either provides upper bounds on the Kullback-Leibler divergence, or else provides analytical expressions that rely on an approximation that is valid only when the von Mises-Fisher distribution is close to the uniform distribution over the hypersphere, something that we do not rely on. Where updating of the von Mises-Fisher distribution has been considered, this appears to have mainly centred on the likelihood function and its characterisation rather than on measures of divergence *per se*. We refer to LIN et al. (2017) and MARDIA and EL-ATOUM (1976) as pertinent examples.

A von Mises-Fisher distribution constitutes a conjugate prior (MARDIA and EL-ATOUM, 1976). Our choice of reference distributions—another, distinct, von Mises-Fisher distribution on \mathbb{S}^i and the uniform distribution on the hypersphere—reflects both this and the prevalence of the uniform distribution in practice. Another commonly invoked choice that we do not consider is the Jeffreys prior, which is proportional to the square root of the determinant of the Fisher information matrix relative to the parametrisation employed.⁶ HORNIK and GRÜN (2013) derive the Fisher information matrix and its determinant and show that the Jeffreys prior is improper in this setting.

—| SECTION 1 |—

THE PROBABILITY DENSITY FUNCTION AND MOMENTS OF THE VON MISES-FISHER FAMILY

Throughout this chapter, we exploit the fact that the von Mises-Fisher family is an exponential family and, accordingly, we adopt the terminology that is used in conjunction with that well-known class. Defining $m \doteq p/2 - 1$ for convenience and maintaining $\kappa \geq 0$ and $\|\mu\|_2 = 1$, we write the probability density function of a von Mises-Fisher random vector, which we denote by $\mathbf{V} \in \mathbb{S}^{p-1}$, for integer $p > 1$, as

$$f(\mathbf{v}; \kappa, \mu) \doteq \frac{\exp(\kappa \cdot \mu^\top \mathbf{v})}{\mathfrak{K}(\kappa, m)} \tag{C.3}$$

5 The von Mises-Fisher family is akin to the class of Gaussian distributions that feature a diagonal variance matrix—i.e., statistically independent Gaussian random variables—with the same entry on each element of the diagonal.

6 Our parametrisation of the von Mises-Fisher family specifies a concentration and a mean direction—what HORNIK and GRÜN (2013) call a polar coordinate parametrisation. Equally, one could combine the concentration and mean direction into a single parameter vector, η say, such that the concentration corresponds to $\|\eta\|_2$ and the mean direction corresponds to $\eta/\|\eta\|_2$. This is the parametrisation that HORNIK and GRÜN (2013) specifies.

where we reiterate that κ and μ are the concentration and mean direction, respectively, and where, for all $\kappa \geq 0$,

$$\mathbb{K}(\kappa, m) \doteq \int_{\mathbb{S}^{p-1}} \exp(\kappa \cdot \mu^\top \mathbf{v}) \cdot d\mathbf{v} = \begin{cases} [2\pi]^{m+1} \cdot \mathbf{I}_m(\kappa) / \kappa^m & \text{if } \kappa > 0 \\ 2\pi^{m+1} / \Gamma(m+1) & \text{if } \kappa = 0 \end{cases} \quad \text{C.4}$$

thereby guaranteeing that the density function integrates to one.⁷ We recall the definition of the modified Bessel function in Chapter B. We refer to the exponentiation—i.e., the numerator in Equation C.3—as the kernel of the density function, and to the normalising constant—i.e., the denominator in Equation C.3—as the partition function. We emphasise that the integral in Equation C.4 is over the hypersphere and it is this fact that makes derivation of statistical features of the von Mises-Fisher random vector difficult. Moreover, we note that our choice of parametrisation is but one way that a von Mises-Fisher distribution can be parametrised. Another parametrisation that is better suited to certain analyses of von Mises-Fisher distributions (in particular, derivation of their moments) is presented in HORNIK and GRÜN (2013).

The von Mises-Fisher family is the hyperspherical analogue of the Gaussian family, which is informative as to its shape: the von Mises-Fisher family is unimodal and symmetric about its mean direction, with the concentration parameter determining the degeneracy (when $\kappa \rightarrow \infty$) and uniformity (when $\kappa \rightarrow 0$) of the distribution, and it assigns positive density to the entirety of the hypersphere. This relationship is shown by appropriately normalising the probability density function of statistically independent Gaussian random variables with variance $1/\kappa$ that are distributed on the hypersphere. Importantly, the von Mises-Fisher family is rotationally equivariant, which is the hyperspherical analogue of the translation equivariance property that is exhibited by the Gaussian family.

The (centred) moments of the von Mises-Fisher family can be obtained by differentiating its moment-generating function. Since the von Mises-Fisher family is an exponential family, the moment-generating function is equal to the log-partition function. Whilst the partition function of the von Mises-Fisher family has a closed-form expression, this is not necessarily true for other named directional families. For instance, we are unaware of any closed-form expressions for the partition functions of both the Kent and von Mises families of distributions beyond the bivariate case.⁸ Differentiating the log-partition function,

⁷ See DHILLON and SRA (2003, §B.2) for proof of this statement, and also for derivation of the determinant of the Jacobian matrix in the generalised spherical coordinate transform of this integral.

⁸ Whilst numerical methods can be used to compute the moments of the Kent and von Mises families (see BEST and N. I. FISHER, 1979; KENT et al., 2013 for suitable rejection sampling routines; also MARDIA et al., 2008, which establishes that the conditional von Mises density is itself von Mises), the lack of closed-form expressions for these other named directional families certainly limits their tractability and usefulness in applications featuring a high-dimensional covariate or outcome vector.

we obtain

$$\mathbb{E}(\mathbf{V}) = \frac{\mathbf{I}_{m+1}(\kappa)}{\mathbf{I}_m(\kappa)} \cdot \boldsymbol{\mu} \tag{C.5}$$

which exploits results in NIST (2021, §10.29). The higher-order moments of the von Mises-Fisher family can be derived via recursive differentiation.⁹

Equation C.5 is adapted from results that are available in MARDIA and JUPP (2009), which also discusses several other important results, including the asymptotic and high-concentration behaviour of von Mises-Fisher random vectors and their tangent normal vectors. We emphasise that the expression for the variance that we present here is distinct from the circular variance, which is simply the distance of the mean resultant length from the surface—i.e., one minus the mean resultant length.

—| SECTION 2 |—

MEASURING THE f -DIVERGENCE OF AN OBTAINED DISTRIBUTION FROM A REFERENCE DISTRIBUTION

To facilitate consideration of the f -divergence of an obtained distribution from a reference distribution, we suppose that, for all $\mathbf{a} \subset \mathbb{S}^{p-1}$,

$$\begin{aligned} \Pr(\mathbf{Q} \in \mathbf{a}) &= \int_{\mathbf{a}} f(\mathbf{v}; \kappa_q, \mu_q) \cdot d\mathbf{v} \\ \Pr(\mathbf{R} \in \mathbf{a}) &= \int_{\mathbf{a}} f(\mathbf{v}; \kappa_r, \mu_r) \cdot d\mathbf{v} \end{aligned} \tag{C.8}$$

such that $\mathbf{Q} \in \mathbb{S}^{p-1}$ and $\mathbf{R} \in \mathbb{S}^{p-1}$ are two von Mises-Fisher random vectors with the corresponding mean direction and concentration parameters satisfying all of the usual properties that we maintain.¹⁰ We use the notation above to distinguish the two von Mises-Fisher random vectors from the κ -concentrated $\boldsymbol{\mu}$ -oriented \mathbf{V} that is previously discussed.

We adopt the convention of referring to the probability distribution of the random vector \mathbf{Q} as the ob-

⁹ For reference, the second moment of the von Mises-Fisher family is

$$\text{Variance}(\mathbf{V}) = \frac{\mathbf{I}_{m+1}(\kappa)}{\mathbf{I}_m(\kappa)} \cdot \left[\frac{1}{\kappa} \cdot \mathbf{I}_p + \left[\frac{\mathbf{I}_{m+2}(\kappa)}{\mathbf{I}_{m+1}(\kappa)} - \frac{\mathbf{I}_{m+1}(\kappa)}{\mathbf{I}_m(\kappa)} \right] \cdot \boldsymbol{\mu} \boldsymbol{\mu}^\top \right] \tag{C.6}$$

where \mathbf{I}_p is the identity matrix of subscripted dimension; this expression can otherwise be written as

$$\text{Variance}(\mathbf{V}) = \frac{1}{\kappa} \cdot \frac{\mathbf{I}_{m+1}(\kappa)}{\mathbf{I}_m(\kappa)} \cdot \mathbf{I}_p + \left[1 - \frac{\mathbf{I}_{m+1}(\kappa)}{\mathbf{I}_m(\kappa)} \cdot \left[\frac{p}{\kappa} - \frac{\mathbf{I}_{m+1}(\kappa)}{\mathbf{I}_m(\kappa)} \right] \right] \cdot \boldsymbol{\mu} \boldsymbol{\mu}^\top \tag{C.7}$$

which exploits recurrence relations in NIST (2021, §10.29). As noted in HORNIK and GRÜN (2013), given that the von Mises-Fisher family is an exponential family, its second moment coincides with its Fisher information matrix under the parametrisation considered in that paper. We derive the first and second moments of the von Mises-Fisher family in the appendices.

¹⁰ A necessary condition that we require is that the probability distribution of the random vector \mathbf{Q} is absolutely continuous with respect to the probability distribution of the random vector \mathbf{R} (LATTIMORE and SZEPESVÁRI, 2020, §2.7 and §14.5), thereby guaranteeing that each of the measures that we consider is well-defined. The von Mises-Fisher family clearly satisfies this requirement, assigning positive density to all points on the hypersphere.

tained distribution and to the probability distribution of the random vector \mathbf{R} as the reference distribution, and of measuring the divergence of the obtained distribution from the reference distribution.

We reiterate that the uniform distribution on the hypersphere corresponds to the limiting case where the concentration of the von Mises-Fisher distribution is zero. As such, we emphasise that our framework is also compatible with the reference distribution being equal to the uniform distribution on the hypersphere.

For all of the measures that we consider, each measure approaches its maximum possible value when the obtained distribution is degenerate. We defer proof of each proposition in this section to the appendices.

With the exception of the total variation distance, each of the measures that we consider has a closed-form expression that depends only upon the parameters of the two distributions whose divergence they capture (and even then, most of the expressions that we obtain simply feature the concentration parameter). These expressions feature ratios of modified Bessel functions, whether of consecutive order or of the same order. Various statistical programming languages have inbuilt functionality to compute modified Bessel functions that perform well for comparatively small values of their argument (i.e., for low concentrations) but that can break for large values. AMOS (1974) derives some inequalities (see the appendices) that are robust to large arguments and that can be used when inbuilt functionality breaks.

Whilst one or more of these measures may be of ultimate interest, it is perhaps more likely that they form only a constituent part of what the researcher is interested in (see, for instance, KITAGAWA et al., 2022b). The rejection sampling algorithm of WOOD (1994; see TSAGRIS et al., 2022, which implements the suggested algorithm, for further details) can be used to sample vectors from a von Mises-Fisher distribution as part of a routine to approximate related objects that do not have a closed-form expression.

Asymptotic calculations can be performed using Hankel series expansion (see NIST, 2021, §10.17 and §10.40, which rely upon §2.1 and the Poincaré expression therein). We include a demonstration in the appendices, applying a Hankel series expansion to the circular variance of the von Mises-Fisher family of distributions,¹¹ which is defined as one minus the ratio of modified Bessel functions—i.e., one minus the mean resultant length. We note that Hankel series expansion is appropriate when $\kappa \rightarrow \infty$, which is the limiting behaviour that we are interested in. We emphasise that the circular variance is strictly contained in the unit interval, which is asymptotically guaranteed by the fact that $p \geq 2$ —i.e., the minimal hypersphere (the circle) is defined on the real plane.

¹¹ See MARDIA and JUPP (2009) for further, general, details about the circular variance.

RÉNYI DIVERGENCE

Following VAN ERVEN and HARREMOËS (2014), we define the Rényi divergence of simple order $0 < \alpha < 1$ or $1 < \alpha < \infty$ as

$$d_\alpha(\mathbf{Q}, \mathbf{R}) \doteq \frac{1}{\alpha - 1} \cdot \ln \left(\int_{\mathbb{S}^{p-1}} f(\mathbf{v}; \kappa_q, \mu_q)^\alpha \cdot f(\mathbf{v}; \kappa_r, \mu_r)^{1-\alpha} \cdot d\mathbf{v} \right) \quad \text{C.9}$$

The Rényi divergence is the most general measure of f -divergence that we consider and describes a broad class that relates to the other measures that we study—namely, the χ -squared distance, the squared-Hellinger distance and the Kullback-Leibler divergence.

Proposition C.1. *Given $\mathbf{Q} \in \mathbb{S}^{p-1}$ and $\mathbf{R} \in \mathbb{S}^{p-1}$ are distributed as p -variate von Mises-Fisher random vectors with concentrations $\kappa_q > 0$ and $\kappa_r > 0$, respectively, and mean directions $\mu_q \in \mathbb{S}^{p-1}$ and $\mu_r \in \mathbb{S}^{p-1}$, respectively, and recalling the definition of m , the Rényi divergence of simple order $0 < \alpha < 1$ or $1 < \alpha < \infty$ is*

$$d_\alpha(\mathbf{Q}, \mathbf{R}) = \frac{m}{\alpha - 1} \cdot \ln \left(\frac{\kappa_q^\alpha \cdot \kappa_r^{1-\alpha}}{\kappa_\alpha} \right) + \frac{\alpha}{\alpha - 1} \cdot \ln \left(\frac{I_m(\kappa_\alpha)}{I_m(\kappa_q)} \right) - \ln \left(\frac{I_m(\kappa_\alpha)}{I_m(\kappa_r)} \right) \quad \text{C.10}$$

where

$$\kappa_\alpha \doteq \|\alpha \cdot \kappa_q \cdot \mu_q + (1 - \alpha) \cdot \kappa_r \cdot \mu_r\|_2 \quad \text{C.11}$$

In the special case where $\kappa_\alpha = 0$,

$$d_\alpha(\mathbf{Q}, \mathbf{R}) = \frac{m}{\alpha - 1} \cdot \ln \left(\frac{\kappa_q^\alpha \cdot \kappa_r^{1-\alpha}}{2} \right) + \frac{\alpha}{\alpha - 1} \cdot \ln \left(\frac{1/\Gamma(m+1)}{I_m(\kappa_q)} \right) - \ln \left(\frac{1/\Gamma(m+1)}{I_m(\kappa_r)} \right) \quad \text{C.12}$$

with

$$\kappa_r = \frac{\alpha}{|1 - \alpha|} \cdot \kappa_q \text{ and } \mu_r + \text{Sign}(1 - \alpha) \cdot \mu_q = 0 \quad \text{C.13}$$

In the special case where \mathbf{R} is, instead, uniformly distributed on the hypersphere,

$$d_\alpha(\mathbf{Q}, \mathbf{R}) = \frac{m}{\alpha - 1} \cdot \ln \left(\frac{2^{1-\alpha}}{\alpha \cdot \kappa_q^{1-\alpha}} \right) + \frac{\alpha}{\alpha - 1} \cdot \ln \left(\frac{I_m(\alpha \cdot \kappa_q)}{I_m(\kappa_q)} \right) - \ln \left(\frac{I_m(\alpha \cdot \kappa_q)}{1/\Gamma(m+1)} \right) \quad \text{C.14}$$

The Rényi divergence is an $O(\ln(\kappa_q))$ function.

We define the χ -square distance as

$$d_\chi(\mathbf{Q}, \mathbf{R}) \doteq \int_{\mathbb{S}^{p-1}} \frac{[f(\mathbf{v}; \kappa_q, \mu_q) - f(\mathbf{v}; \kappa_r, \mu_r)]^2}{f(\mathbf{v}; \kappa_r, \mu_r)} \cdot d\mathbf{v} \quad \text{C.15}$$

$$= \int_{\mathbb{S}^{p-1}} \frac{f(\mathbf{v}; \kappa_q, \mu_q)^2}{f(\mathbf{v}; \kappa_r, \mu_r)} \cdot d\mathbf{v} - 1 \quad \text{C.16}$$

with Equation C.16 being the more convenient definition to work with. We observe that the χ -square distance relates to the Rényi class of measures via the equality relation

$$d_\alpha(\mathbf{Q}, \mathbf{R})|_{\alpha=2} = \ln(1 + d_\chi(\mathbf{Q}, \mathbf{R})) \quad \text{C.17}$$

as per VAN ERVEN and HARREMOËS (2014), which serves to illustrate the breadth of the Rényi class.

Proposition C.2. *Given $\mathbf{Q} \in \mathbb{S}^{p-1}$ and $\mathbf{R} \in \mathbb{S}^{p-1}$ are distributed as p -variate von Mises-Fisher random vectors with concentrations $\kappa_q > 0$ and $\kappa_r > 0$, respectively, and mean directions $\mu_q \in \mathbb{S}^{p-1}$ and $\mu_r \in \mathbb{S}^{p-1}$, respectively, and recalling the definition of m , the χ -square distance is*

$$d_\chi(\mathbf{Q}, \mathbf{R}) = \frac{\kappa_q^{2m} \cdot \mathbb{I}_m(\kappa_\chi) \cdot \mathbb{I}_m(\kappa_r)}{\kappa_\chi^m \cdot \kappa_r^m \cdot \mathbb{I}_m(\kappa_q)^2} - 1 \quad \text{C.18}$$

where

$$\kappa_\chi \doteq \|2\kappa_q \cdot \mu_q - \kappa_r \cdot \mu_r\|_2 \quad \text{C.19}$$

In the special case where $\kappa_\chi = 0$,

$$d_\chi(\mathbf{Q}, \mathbf{R}) = \frac{\kappa_q^m \cdot \mathbb{I}_m(2\kappa_q)}{4^m \cdot \Gamma(m+1) \cdot \mathbb{I}_m(\kappa_q)^2} - 1 \quad \text{C.20}$$

with

$$\kappa_r = 2\kappa_q \text{ and } \mu_r - \mu_q = 0 \quad \text{C.21}$$

In the special case where \mathbf{R} is, instead, uniformly distributed on the hypersphere, the χ -square distance coincides with Equation C.20. The χ -square distance is an $\mathcal{O}(\kappa_q)$ function.

SQUARED-HELLINGER DISTANCE

We define the squared-Hellinger distance as

$$d_h(\mathbf{Q}, \mathbf{R})^2 \doteq \int_{\mathbb{S}^{p-1}} \left[\sqrt{f(\mathbf{v}; \kappa_q, \mu_q)} - \sqrt{f(\mathbf{v}; \kappa_r, \mu_r)} \right]^2 \cdot d\mathbf{v} \quad \text{C.22}$$

$$= 2 \left[1 - \int_{\mathbb{S}^{p-1}} \sqrt{f(\mathbf{v}; \kappa_q, \mu_q) \cdot f(\mathbf{v}; \kappa_r, \mu_r)} \cdot d\mathbf{v} \right] \quad \text{C.23}$$

with Equation C.23 being the more convenient definition to work with. We observe that the squared-Hellinger distance relates to the Rényi class of measures via the equality relation

$$d_\alpha(\mathbf{Q}, \mathbf{R})|_{\alpha=1/2} = -2 \ln \left(1 - \frac{d_h(\mathbf{Q}, \mathbf{R})^2}{2} \right) \quad \text{C.24}$$

as per VAN ERVEN and HARREMOËS (2014), which serves to illustrate the breadth of the Rényi class and is also interpretable as twice the negative logarithm of the Bhattacharyya coefficient (BHATTACHARYYA, 1943). The Bhattacharyya coefficient is an approximate measure of the amount of overlap between two probability distributions, such that when the obtained and reference distributions are close (i.e., they have a similar concentration and mean direction) then the squared-Hellinger distance is small in absolute value, as is to be expected.

Proposition C.3. *Given $\mathbf{Q} \in \mathbb{S}^{p-1}$ and $\mathbf{R} \in \mathbb{S}^{p-1}$ are distributed as p -variate von Mises-Fisher random vectors with concentrations $\kappa_q > 0$ and $\kappa_r > 0$, respectively, and mean directions $\mu_q \in \mathbb{S}^{p-1}$ and $\mu_r \in \mathbb{S}^{p-1}$, respectively, and recalling the definition of m , the squared-Hellinger distance is*

$$d_h(\mathbf{Q}, \mathbf{R})^2 = 2 \left[1 - \sqrt{\frac{\kappa_q^m \cdot \kappa_r^m \cdot \text{I}_m(\kappa_h)^2}{\kappa_h^{2m} \cdot \text{I}_m(\kappa_q) \cdot \text{I}_m(\kappa_r)}} \right] \quad \text{C.25}$$

where we define

$$\kappa_h \doteq \frac{1}{2} \cdot \|\kappa_q \cdot \mu_q + \kappa_r \cdot \mu_r\|_2 \quad \text{C.26}$$

In the special case where $\kappa_h = 0$,

$$d_h(\mathbf{Q}, \mathbf{R})^2 = 2 \left[1 - \frac{\kappa_q^m}{2^m \cdot \text{I}_m(\kappa_q) \cdot \Gamma(m+1)} \right] \quad \text{C.27}$$

with

$$\kappa_r = \kappa_q \text{ and } \mu_r + \mu_q = 0 \tag{C.28}$$

In the special case where \mathbf{R} is, instead, uniformly distributed on the hypersphere,

$$d_h(\mathbf{Q}, \mathbf{R})^2 = 2 \left[1 - \sqrt{\frac{2^{3m} \cdot I_m(\kappa_q/2)^2 \cdot \Gamma(m+1)}{\kappa_q^m \cdot I_m(\kappa_q)}} \right] \tag{C.29}$$

The squared-Hellinger distance is an $O(1 - 1/\kappa_q)$ function.

—| SUBSECTION 2.D |—

KULLBACK-LEIBLER DIVERGENCE

We define the Kullback-Leibler divergence as

$$d_\ell(\mathbf{Q}, \mathbf{R}) \doteq \int_{\mathbb{S}^p} \ln \left(\frac{f(\mathbf{v}; \kappa_q, \mu_q)}{f(\mathbf{v}; \kappa_r, \mu_r)} \right) \cdot f(\mathbf{v}; \kappa_q, \mu_q) \cdot d\mathbf{v} \tag{C.30}$$

We observe that the Kullback-Leibler divergence is a limiting case of the Rényi divergence. That is,

$$\lim_{\alpha \rightarrow 1} d_\alpha(\mathbf{Q}, \mathbf{R}) = d_\ell(\mathbf{Q}, \mathbf{R}) \tag{C.31}$$

as per VAN ERVEN and HARREMOËS (2014), which serves to illustrate the breadth of the Rényi class.

Proposition C.4. *Given $\mathbf{Q} \in \mathbb{S}^{p-1}$ and $\mathbf{R} \in \mathbb{S}^{p-1}$ are distributed as p -variate von Mises-Fisher random vectors with concentrations $\kappa_q > 0$ and $\kappa_r > 0$, respectively, and mean directions $\mu_q \in \mathbb{S}^{p-1}$ and $\mu_r \in \mathbb{S}^{p-1}$, respectively, and recalling that $m \doteq p/2 - 1$, the Kullback-Leibler divergence is*

$$d_\ell(\mathbf{Q}, \mathbf{R}) = m \cdot \ln \left(\frac{\kappa_q}{\kappa_r} \right) - \ln \left(\frac{I_m(\kappa_q)}{I_m(\kappa_r)} \right) + \frac{I_{m+1}(\kappa_q)}{I_m(\kappa_q)} \cdot [\kappa_q \cdot \mu_q - \kappa_r \cdot \mu_r]^\top \mu_q \tag{C.32}$$

In the special case where \mathbf{R} is, instead, uniformly distributed on the hypersphere,

$$d_\ell(\mathbf{Q}, \mathbf{R}) = m \cdot \ln \left(\frac{\kappa_q}{2} \right) - \ln(I_m(\kappa_q)) - \ln(\Gamma(m+1)) + \frac{I_{m+1}(\kappa_q)}{I_m(\kappa_q)} \cdot \kappa_q \tag{C.33}$$

The Kullback-Leibler divergence is an $O(\ln(\kappa_q))$ function.

We note that the final terms of Equations C.32 and C.33 are proportional to the first moment of a von Mises-Fisher distribution (specifically, the first moment of the obtained distribution). In particular, we

note that the final term of Equation C.32 can otherwise be written as

$$\frac{I_{m+1}(\kappa_q)}{I_m(\kappa_q)} \cdot [\kappa_q \cdot \mu_q - \kappa_r \cdot \mu_r]^\top \mu_q = \frac{I_{m+1}(\kappa_q)}{I_m(\kappa_q)} \cdot [\kappa_q - \kappa_r \cdot \mu_r^\top \mu_q] \quad \text{C.34}$$

which we contrast with the corresponding term in Equation C.33. The difference between Equations C.33 and C.34 arises because the mean direction does not enter the probability density function of the uniform distribution on the hypersphere. The intuition here is that, in the special case where the reference distribution is the uniform distribution on the hypersphere, the mean direction can always be taken to be equal to the mean direction.

—| SUBSECTION 2.E |—

TOTAL VARIATION DISTANCE

We define the total variation distance, for all measurable $\mathbf{a} \subset \mathbb{S}^{p-1}$, as

$$d_t(\mathbf{Q}, \mathbf{R}) \doteq \sup_{\mathbf{a}} \left| \int_{\mathbf{a}} f(\mathbf{v}; \kappa_q, \mu_q) - f(\mathbf{v}; \kappa_r, \mu_r) \cdot d\mathbf{v} \right| \quad \text{C.35}$$

which is well-defined if the underlying probability space is endowed with the Borel σ -algebra—something that we, implicitly, maintain throughout all parts of our analysis.

Corollary C.1. *Given $\mathbf{Q} \in \mathbb{S}^{p-1}$ and $\mathbf{R} \in \mathbb{S}^{p-1}$ are distributed as p -variate von Mises-Fisher random vectors with concentrations $\kappa_q > 0$ and $\kappa_r > 0$, respectively, and mean directions $\mu_q \in \mathbb{S}^{p-1}$ and $\mu_r \in \mathbb{S}^{p-1}$, respectively, and recalling that $m \doteq p/2 - 1$, the total variation distance satisfies the inequality relations*

$$d_t(\mathbf{Q}, \mathbf{R})^2 \leq d_h(\mathbf{Q}, \mathbf{R})^2 \leq d_\ell(\mathbf{Q}, \mathbf{R}) \leq d_\chi(\mathbf{Q}, \mathbf{R}) \quad \text{C.36}$$

$$d_t(\mathbf{Q}, \mathbf{R}) \leq \sqrt{1 - \exp(-d_\ell(\mathbf{Q}, \mathbf{R}))} \leq 1 - \frac{d_\ell(\mathbf{Q}, \mathbf{R})}{2} \quad \text{C.37}$$

and, for all orders $0 < \alpha \leq 1$,

$$d_t(\mathbf{Q}, \mathbf{R})^2 \leq \frac{\alpha \cdot d_\alpha(\mathbf{Q}, \mathbf{R})}{2} \quad \text{C.38}$$

where each measure of divergence is as defined in Section 2.

The inequalities in Corollary C.1 are well-known (see, for instance, LATTIMORE and SZEPESVÁRI, 2020, §14.3 and references therein, and VAN ERVEN and HARREMOËS, 2014). The total variation distance is difficult to characterise in this setting due to the need to integrate over a region of the hypersphere. The inequality relations stated in Corollary C.1 provide a means to bound the total variation distance

from above using measures that are more easily characterised. In particular, Equation C.38 is a generalisation of Pinsker's inequality, yielding the classic statement of that inequality at the limit (see Equation C.31).

—| APPENDIX C.1 |—

PROOFS

\hookrightarrow *Proof of Equation C.5.* We reiterate that the von Mises-Fisher family is an exponential family and, as such, its moments can be obtained via differentiation of the log-partition function, with the first moment equal to the first derivative, the second moment equal to the second derivative, and so forth. To facilitate our analysis of the log-partition function and its derivatives, we rewrite the probability density function of the von Mises-Fisher family in terms of a single parameter vector. We let

$$\begin{aligned}\kappa &\doteq \|\eta\|_2 \\ \mu &\doteq \eta/\|\eta\|_2\end{aligned}\tag{C.39}$$

where $\eta \in \mathbb{R}^p$. The log-partition function is equal to

$$\ln(\mathfrak{K}(\kappa, m)) = \ln\left(\int_{\mathbb{S}^{p-1}} \exp(\kappa \cdot \mu^\top \mathbf{v}) \cdot d\mathbf{v}\right) = -m \cdot \ln(\kappa) + [m+1] \cdot \ln(2\pi) + \ln(I_m(\kappa))\tag{C.40}$$

The moments of the von Mises-Fisher family of distributions are obtained by recursively differentiating this function with respect to the new parameter vector.

We now present some derivatives that are useful for the construction of the first, second and higher-order moments. First,

$$\begin{aligned}\frac{d}{d\eta} \kappa &= \mu \\ \frac{d}{d\eta^\top} \mu &= \frac{1}{\kappa} \cdot [\mathbf{I}_p - \mu\mu^\top]\end{aligned}\tag{C.41}$$

Second, for all $n \in \mathbb{N}_*$,

$$\begin{aligned}\frac{d}{d\eta^\top} \frac{I_{m+n}(\kappa)}{I_m(\kappa)} &= \frac{I_m \cdot (\kappa) \cdot I'_{m+n}(\kappa) - I_{m+n}(\kappa) \cdot I'_m(\kappa)}{I_m(\kappa)^2} \cdot \mu^\top \\ &= \left[\frac{I_{m+n+1}(\kappa)}{I_m(\kappa)} + \frac{n \cdot I_{m+n}(\kappa)}{\kappa \cdot I_m(\kappa)} - \frac{I_{m+n}(\kappa)}{I_m(\kappa)} \cdot \frac{I_{m+1}(\kappa)}{I_m(\kappa)} \right] \cdot \mu^\top \\ &= \left[\prod_{j=0}^{n-1} \frac{I_{m+j+1}(\kappa)}{I_{m+j}(\kappa)} \right] \cdot \left[\frac{I_{m+n+1}(\kappa)}{I_{m+n}(\kappa)} + \frac{n}{\kappa} - \frac{I_{m+1}(\kappa)}{I_m(\kappa)} \right] \cdot \mu^\top\end{aligned}\tag{C.42}$$

which exploits Equation B.4 and the telescoping property of the ratios. We now differentiate Equa-

tion C.40 (i.e., we take the first derivative of the log-partition function), which yields

$$\frac{d}{d\eta} \ln(\mathfrak{K}(\kappa, m)) = \frac{I_{m+1}(\kappa)}{I_m(\kappa)} \cdot \mu \quad \text{C.43}$$

where we use Equation C.41 and Equation B.4. This proves the first result of the corollary. We then differentiate Equation C.43 (i.e., we take the second derivative of the log-partition function), which yields

$$\frac{d}{d\eta^\top} \frac{I_{m+1}(\kappa)}{I_m(\kappa)} \cdot \mu = \frac{I_{m+1}(\kappa)}{I_m(\kappa)} \cdot \left[\frac{1}{\kappa} \cdot \mathbf{I}_p + \left[\frac{I_{m+2}(\kappa)}{I_{m+1}(\kappa)} - \frac{I_{m+1}(\kappa)}{I_m(\kappa)} \right] \cdot \mu \mu^\top \right] \quad \text{C.44}$$

where we use Equations B.4, C.41 and C.42. Substituting

$$\frac{I_{m+1}(\kappa)}{I_m(\kappa)} \cdot \frac{I_{m+2}(\kappa)}{I_{m+1}(\kappa)} = 1 - \frac{p}{\kappa} \cdot \frac{I_{m+1}(\kappa)}{I_m(\kappa)} \quad \text{C.45}$$

which is valid by Equation B.4, we obtain an alternative expression for the variance. This proves the second result of the corollary. The higher-order moments of the von Mises-Fisher family can be obtained via recursive differentiation of these expressions, with application of the product rule of differentiation and the derivatives that are presented in Equations C.41 and C.42 then sufficient to construct these moments. ■

To facilitate characterisation of the limiting behaviour of the various measures of divergence that we consider, we apply Equation B.6.

Corollary C.2. *Given that $\kappa > 1$, then for $\underline{\mathfrak{H}}(\kappa, 0, m)$ and $\overline{\mathfrak{B}}(\kappa, 0, m)$ defined as in Equations B.9 and B.10, with $m \geq 0$*

$$\underline{\mathfrak{H}}(\kappa, 0, m) \geq \kappa - \frac{1}{2} \cdot \ln(\kappa) - m \cdot \ln(2) - \ln(\Gamma(m+1)) - \bar{a}_m \quad \text{C.46}$$

and

$$\overline{\mathfrak{B}}(\kappa, 0, m) \leq \kappa - \frac{1}{2} \cdot \ln(\kappa) - m \cdot \ln(2) - \ln(\Gamma(m+1)) + \underline{a}_m \cdot \ln(2\underline{a}_m) \quad \text{C.47}$$

such that $\ln(I_m(\kappa))$ is an $O(\kappa - \ln(\kappa)/2)$ function.

Corollary C.2 follows from Equations B.9 and B.10 and the simple exploitation of the properties of increasing concave functions. Moreover, that we choose to state that the logarithm of the modified Bessel function is an $O(\kappa - \ln(\kappa)/2)$ function rather than a linear function is for practical reasons. Specifically, we rely on Corollary C.2 to prove many of the statements in Section 1, for which we find

that the linear component typically cancels. We observe that Corollary C.2 aligns with results elsewhere NIST (2021, §10.30).

↪ *Proof of Corollary C.2.* We rewrite Equations B.9 and B.10 using the formula for the difference of two squares as

$$\underline{\mathbb{B}}(\kappa, 0, m) = \frac{1}{2} \cdot \ln\left(\frac{2}{\kappa}\right) - \ln(\Gamma(m+1)) + \underline{a}_m \cdot \ln\left(\frac{\kappa \cdot [\underline{a}_m + \bar{a}_m]/2}{\underline{a}_m + \sqrt{\kappa^2 + \bar{a}_m^2}}\right) + \sqrt{\kappa^2 + \bar{a}_m^2} - \bar{a}_m \quad \text{C.48}$$

$$\geq \frac{1}{2} \cdot \ln\left(\frac{2}{\kappa}\right) - \ln(\Gamma(m+1)) + \underline{a}_m \cdot \ln\left(\frac{\kappa \cdot [\underline{a}_m + \bar{a}_m]/2}{\underline{a}_m + \kappa + \bar{a}_m}\right) + \kappa - \bar{a}_m \quad \text{C.49}$$

$$= \frac{1}{2} \cdot \ln\left(\frac{2}{\kappa}\right) - \ln(\Gamma(m+1)) + \underline{a}_m \cdot \ln(\kappa) + \underline{a}_m \cdot \ln\left(\frac{[\underline{a}_m + \bar{a}_m]/2}{\underline{a}_m + \bar{a}_m + \kappa}\right) + \kappa - \bar{a}_m \quad \text{C.50}$$

$$\geq \frac{1}{2} \cdot \ln\left(\frac{2}{\kappa}\right) - \ln(\Gamma(m+1)) + (\underline{a}_m - \bar{a}_m) \cdot \ln(\kappa) + \underline{a}_m \cdot \ln\left(\frac{\underline{a}_m + \bar{a}_m}{2[\underline{a}_m + \bar{a}_m]}\right) + \kappa - \bar{a}_m \quad \text{C.51}$$

$$= \kappa - \frac{1}{2} \cdot \ln(\kappa) - m \cdot \ln(2) - \ln(\Gamma(m+1)) - \bar{a}_m \quad \text{C.52}$$

and

$$\bar{\mathbb{B}}(\kappa, 0, m) = \frac{1}{2} \cdot \ln\left(\frac{2}{\kappa}\right) - \ln(\Gamma(m+1)) + \underline{a}_m \cdot \ln\left(\frac{\kappa \cdot [\underline{a}_m + \bar{a}_m]/2}{\underline{a}_m + \sqrt{\kappa^2 + \bar{a}_m^2}}\right) + \sqrt{\kappa^2 + \bar{a}_m^2} - \underline{a}_m \quad \text{C.53}$$

$$\leq \frac{1}{2} \cdot \ln\left(\frac{2}{\kappa}\right) - \ln(\Gamma(m+1)) + \underline{a}_m \cdot \ln\left(\frac{2\underline{a}_m \cdot \kappa}{2\kappa}\right) + \kappa + \underline{a}_m - \bar{a}_m \quad \text{C.54}$$

$$= \kappa - \frac{1}{2} \cdot \ln(\kappa) - m \cdot \ln(2) - \ln(\Gamma(m+1)) + \underline{a}_m \cdot \ln(2\underline{a}_m) \quad \text{C.55}$$

respectively. Here, we exploit the concavity of the logarithmic and square root functions, and set some terms equal to zero where it is useful to do so. This establishes the veracity of the first part of Corollary C.2. We now need to show that the logarithm of the modified Bessel function, which we know lies between these two bounds for all $\kappa > 0$, is an $O(\kappa - \ln(\kappa)/2)$ function. To do so, we show that the maximum of the absolute value of Equations C.52 and C.55 is bounded from above by a function that has the required order. Specifically, for all $\kappa > 0$,

$$|\ln(\mathbb{I}_m(\kappa))| \leq \kappa - \frac{1}{2} \cdot \ln(\kappa) + m \cdot \ln(2) + \ln(\Gamma(m+1)) + \max(\bar{a}_m, \underline{a}_m \cdot \ln(2\underline{a}_m)) \quad \text{C.56}$$

$$\leq \left[\kappa - \frac{1}{2} \cdot \ln(\kappa) \right] \cdot \left[1 + \frac{m \cdot \ln(2) + \ln(\Gamma(m+1)) + \max(\bar{a}_m, \underline{a}_m \cdot \ln(2\underline{a}_m))}{(1 + \ln(2))/2} \right] \quad \text{C.57}$$

where the denominator in the second term of Equation C.57 is the minimum value that the first term of Equation C.57 can attain, which establishes the result. ■

Each of the various measures of divergence that we consider involves integration of the kernel for some

value of the concentration parameter. The implication of this property is that each measure can be expressed in terms of logarithmic ratios of the form

$$\ln \left(\frac{\mathfrak{K}(\kappa_d, m)}{\mathfrak{K}(\kappa, m)} \right) = \begin{cases} \ln(\mathbb{I}_m(\kappa_d)) - \ln(\mathbb{I}_m(\kappa)) - m \cdot [\ln(\kappa_d) - \ln(\kappa)] & \text{if } \kappa > 0 \\ \ln(\mathbb{I}_m(\kappa_d)) + \ln(\Gamma(m+1)) - m \cdot [\ln(\kappa_d) - \ln(2)] & \text{if } \kappa = 0 \end{cases} \quad \text{C.58}$$

where $\kappa_d > 0$ and $\kappa \geq 0$ are placeholders for either the obtained concentration, the reference concentration, or their weighted average.

\hookrightarrow *Proof of Proposition C.1.* We focus on the integrand in the definition of the Rényi divergence, and note that

$$f(\mathbf{v}; \kappa_q, \mu_q)^\alpha \cdot f(\mathbf{v}; \kappa_r, \mu_r)^{1-\alpha} = \frac{\exp\left(\left(\alpha \cdot \kappa_q \cdot \mu_q + [1-\alpha] \cdot \kappa_r \cdot \mu_r\right)^\top \mathbf{v}\right)}{\mathfrak{K}(\kappa_q, m)^\alpha \cdot \mathfrak{K}(\kappa_r, m)^{1-\alpha}} \quad \text{C.59}$$

We then define

$$\begin{aligned} \kappa_\alpha &\doteq \|\alpha \cdot \kappa_q \cdot \mu_q + [1-\alpha] \cdot \kappa_r \cdot \mu_r\|_2 \\ \mu_\alpha &\doteq [\alpha \cdot \kappa_q \cdot \mu_q + [1-\alpha] \cdot \kappa_r \cdot \mu_r] / \|\alpha \cdot \kappa_q \cdot \mu_q + [1-\alpha] \cdot \kappa_r \cdot \mu_r\|_2 \end{aligned} \quad \text{C.60}$$

such that the resulting expression respects the conventions that we have adopted for the concentration and mean direction—i.e., that the concentration is non-negative and that the mean direction is a unit vector. We substitute Equation C.60 into Equation C.59 and integrate over the hypersphere, finding that

$$\frac{1}{\alpha-1} \cdot \ln \left(\int_{\mathbb{S}^{p-1}} \frac{\exp(\kappa_\alpha \cdot \mu_\alpha^\top \mathbf{v})}{\mathfrak{K}(\kappa_q, m)^\alpha \cdot \mathfrak{K}(\kappa_r, m)^{1-\alpha}} \cdot d\mathbf{v} \right) = \frac{\alpha}{\alpha-1} \cdot \ln \left(\frac{\mathfrak{K}(\kappa_\alpha, m)}{\mathfrak{K}(\kappa_q, m)} \right) - \ln \left(\frac{\mathfrak{K}(\kappa_\alpha, m)}{\mathfrak{K}(\kappa_r, m)} \right) \quad \text{C.61}$$

Hence, applying Equation C.58 to Equation C.61 implies that

$$d_\alpha(\mathbf{Q}, \mathbf{R}) = \frac{\alpha}{\alpha-1} \cdot \ln \left(\frac{\kappa_q^m \cdot \mathbb{I}_m(\kappa_\alpha)}{\kappa_\alpha^m \cdot \mathbb{I}_m(\kappa_q)} \right) - \ln \left(\frac{\kappa_r^m \cdot \mathbb{I}_m(\kappa_\alpha)}{\kappa_\alpha^m \cdot \mathbb{I}_m(\kappa_r)} \right) \quad \text{C.62}$$

which yields Equation C.10 upon rearrangement.

In the special case where $\kappa_\alpha = 0$, the numerator on the right-hand side of Equation C.59 is one, with the corresponding integral over the hypersphere equal to the surface area of the unit ball. The formula

for the surface area of the unit ball is well-known. Hence, Equation C.61 reduces to

$$d_\alpha(\mathbf{Q}, \mathbf{R}) = \frac{m}{\alpha - 1} \cdot \ln \left(\frac{\kappa_q^\alpha \cdot \kappa_r^{1-\alpha}}{2} \right) + \frac{\alpha}{\alpha - 1} \cdot \ln \left(\frac{1/\Gamma(m+1)}{I_m(\kappa_q)} \right) - \ln \left(\frac{1/\Gamma(m+1)}{I_m(\kappa_r)} \right) \quad \text{C.63}$$

We reiterate that this special case occurs only when Equation C.13 holds, which facilitates the restatement of the Rényi divergence in terms of one of the two concentration parameters. We choose to state the Rényi divergence in terms of both concentration parameters because the expression that we obtain is relatively simple.

In the special case where the reference distribution is the uniform distribution on the hypersphere, Equation C.61 reduces to

$$d_\alpha(\mathbf{Q}, \mathbf{R}) = \frac{\alpha}{\alpha - 1} \cdot \ln \left(\frac{\kappa_q^m \cdot I_m(\kappa_\alpha)}{\kappa_\alpha^m \cdot I_m(\kappa_q)} \right) - \ln \left(\frac{2^m \cdot I_m(\kappa_\alpha)}{\kappa_\alpha^m \cdot 1/\Gamma(m+1)} \right) \quad \text{C.64}$$

We then recall that $\kappa_\alpha = \alpha \cdot \kappa_q$ whenever $\kappa_r = 0$, such that

$$d_\alpha(\mathbf{Q}, \mathbf{R}) = \frac{\alpha}{\alpha - 1} \cdot \ln \left(\frac{I_m(\alpha \cdot \kappa_q)}{I_m(\kappa_q)} \right) - \ln \left(\frac{2^m \cdot I_m(\alpha \cdot \kappa_q)}{\kappa_q^m \cdot 1/\Gamma(m+1)} \right) - \frac{m}{\alpha - 1} \cdot \ln(\alpha) \quad \text{C.65}$$

from which we obtain Equation C.14 upon rearrangement.

We now turn our attention to the limiting behaviour of the Rényi divergence with respect to an increase in κ_q given that $\kappa_q > \kappa_r$ and κ_q is sufficiently large such that $\min(\kappa_\alpha, \kappa_q) > 1$, holding $\kappa_r > 0$ fixed. We recall that

$$\underline{\mathbb{H}}(\kappa_q, \kappa_r, m) \leq \ln(I_m(\kappa_q)) \leq \overline{\mathbb{H}}(\kappa_q, \kappa_r, m) \quad \text{C.66}$$

where

$$\underline{\mathbb{H}}(\kappa_q, \kappa_r, m) \doteq \ln(I_m(\kappa_r)) + m \cdot \ln \left(\frac{\kappa_q}{\kappa_r} \right) + \underline{\mathbb{J}}(\kappa_q, \kappa_r, m) \quad \text{C.67}$$

and

$$\overline{\mathbb{H}}(\kappa_q, \kappa_r, m) \doteq \ln(I_m(\kappa_r)) + m \cdot \ln \left(\frac{\kappa_q}{\kappa_r} \right) + \overline{\mathbb{J}}(\kappa_q, \kappa_r, m) \quad \text{C.68}$$

where

$$\underline{\mathbb{J}}\mathbb{B}(\kappa_q, \kappa_r, m) \doteq \underline{a}_m \cdot \ln \left(\frac{\underline{a}_m + \sqrt{\kappa_r^2 + \underline{a}_m^2}}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}} \right) + \sqrt{\kappa_q^2 + \underline{a}_m^2} - \sqrt{\kappa_r^2 + \underline{a}_m^2} \quad \text{C.69}$$

and

$$\overline{\mathbb{J}}\mathbb{B}(\kappa_q, \kappa_r, m) \doteq \underline{a}_m \cdot \ln \left(\frac{\underline{a}_m + \sqrt{\kappa_r^2 + \underline{a}_m^2}}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}} \right) + \sqrt{\kappa_q^2 + \underline{a}_m^2} - \sqrt{\kappa_r^2 + \underline{a}_m^2} \quad \text{C.70}$$

given $\underline{a}_m \doteq m + 1/2$ and $\bar{a}_m \doteq m + 3/2$. We further recall that the Rényi divergence is defined as

$$d_\alpha(\mathbf{Q}, \mathbf{R}) = \frac{m}{\alpha - 1} \cdot \ln \left(\frac{\kappa_q^\alpha \cdot \kappa_r^{1-\alpha}}{\kappa_\alpha} \right) + \frac{\alpha}{\alpha - 1} \cdot \ln \left(\frac{\mathbb{I}_m(\kappa_\alpha)}{\mathbb{I}_m(\kappa_q)} \right) - \ln \left(\frac{\mathbb{I}_m(\kappa_\alpha)}{\mathbb{I}_m(\kappa_r)} \right) \quad \text{C.71}$$

and so satisfies

$$\frac{1}{\alpha - 1} \cdot \textcircled{\text{A}} + \frac{\alpha}{\alpha - 1} \cdot \textcircled{\text{C}} + \frac{1 - \alpha}{\alpha - 1} \cdot \textcircled{\text{D}} \leq d_\alpha(\mathbf{Q}, \mathbf{R}) \leq \frac{1}{\alpha - 1} \cdot \textcircled{\text{B}} - \frac{\alpha}{\alpha - 1} \cdot \textcircled{\text{C}} - \frac{1 - \alpha}{\alpha - 1} \cdot \textcircled{\text{D}} \quad \text{C.72}$$

where

$$\begin{aligned} \textcircled{\text{A}} &\doteq \alpha \cdot \underline{\mathbb{J}}\mathbb{B}(\kappa_\alpha, \kappa_q, m) + [1 - \alpha] \cdot \underline{\mathbb{J}}\mathbb{B}(\kappa_\alpha, \kappa_r, m) \\ \textcircled{\text{B}} &\doteq \alpha \cdot \overline{\mathbb{J}}\mathbb{B}(\kappa_\alpha, \kappa_q, m) + [1 - \alpha] \cdot \overline{\mathbb{J}}\mathbb{B}(\kappa_\alpha, \kappa_r, m) \\ \textcircled{\text{C}} &\doteq 1 \cdot (\text{sign}(\alpha - 1) \neq \text{sign}(\kappa_\alpha - \kappa_q)) \cdot (\overline{\mathbb{J}}\mathbb{B} - \underline{\mathbb{J}}\mathbb{B}) \circ (\kappa_\alpha, \kappa_q, m) \\ \textcircled{\text{D}} &\doteq 1 \cdot (\kappa_\alpha > \kappa_r) \cdot (\overline{\mathbb{J}}\mathbb{B} - \underline{\mathbb{J}}\mathbb{B}) \circ (\kappa_\alpha, \kappa_r, m) \end{aligned} \quad \text{C.73}$$

as per Equation C.66. It immediately follows (by considering each possible ordering of κ and the magnitude of α relative to one) that

$$|d_\alpha(\mathbf{Q}, \mathbf{R})| \leq \frac{1}{|\alpha - 1|} \cdot \left[\underline{a}_m \cdot \left| \ln \left(\frac{\kappa_q^\alpha \cdot \kappa_r^{1-\alpha}}{\kappa_\alpha} \right) \right| + |\kappa_\alpha - \alpha \cdot \kappa_q - [1 - \alpha] \cdot \kappa_r| + \bar{a}_m \right] \quad \text{C.74}$$

which then reduces to

$$|d_\alpha(\mathbf{Q}, \mathbf{R})| \leq \frac{1}{|\alpha - 1|} \cdot \left[\underline{a}_m \cdot \left| \ln \left(\frac{\kappa_q^\alpha \cdot \kappa_r^{1-\alpha}}{\kappa_\alpha} \right) \right| + \bar{a}_m \right] + 2\kappa_r \quad \text{C.75}$$

by virtue of the fact that

$$\alpha \cdot \kappa_q - |1 - \alpha| \cdot \kappa_r \leq \kappa_\alpha \leq \alpha \cdot \kappa_q + |1 - \alpha| \cdot \kappa_r \quad \text{C.76}$$

as per its definition. Moreover, given that

$$\ln \left(\frac{\kappa_q^\alpha \cdot \kappa_r^{1-\alpha}}{\kappa_\alpha} \right) = -\ln \left(\left\| \alpha \cdot \left(\frac{\kappa_q}{\kappa_r} \right)^{1-\alpha} \cdot \mu_q + [1-\alpha] \cdot \left(\frac{\kappa_r}{\kappa_q} \right)^\alpha \cdot \mu_r \right\|_2 \right) \quad \text{C.77}$$

it is then self-evident that the limiting behaviour of the Rényi divergence is then wholly determined by the first term on the right-hand side of Equation C.75, which is $O(\ln(\kappa_q))$ function—and so too the Rényi divergence.¹²

It is straightforward to determine that this result extends to the special case where the reference distribution is the uniform distribution on the hypersphere, since the final two logarithmic terms in Equation C.14 can easily be shown to be of the same order as the first logarithmic term using Corollary C.2. ■

↪ *Proof of Proposition C.2.* Proposition C.2 is implied by Proposition C.1 and Equation C.17. ■

↪ *Proof of Proposition C.3.* Proposition C.3 is implied by Proposition C.1 and Equation C.24. ■

↪ *Proof of Proposition C.4.* We focus on the integrand in the definition of the Kullback-Leibler divergence and note that

$$\ln \left(\frac{f(\mathbf{v}; \kappa_q, \mu_q)}{f(\mathbf{v}; \kappa_r, \mu_r)} \right) \cdot f(\mathbf{v}; \kappa_q, \mu_q) = \left((\kappa_q \cdot \mu_q - \kappa_r \cdot \mu_r)^\top \mathbf{v} + \ln \left(\frac{\mathbb{X}(\kappa_r, m)}{\mathbb{X}(\kappa_q, m)} \right) \right) \cdot f(\mathbf{v}; \kappa_q, \mu_q) \quad \text{C.78}$$

We integrate each term on the right-hand side of Equation C.78, finding that

$$\int_{\mathbb{S}^{p-1}} (\kappa_q \cdot \mu_q - \kappa_r \cdot \mu_r)^\top \mathbf{v} \cdot f(\mathbf{v}; \kappa_q, \mu_q) \cdot d\mathbf{v} = (\kappa_q \cdot \mu_q - \kappa_r \cdot \mu_r)^\top \mathbf{E}(\mathbf{Q}) \quad \text{C.79}$$

and

$$\int_{\mathbb{S}^{p-1}} \ln \left(\frac{\mathbb{X}(\kappa_r, m)}{\mathbb{X}(\kappa_q, m)} \right) \cdot f(\mathbf{v}; \kappa_q, \mu_q) \cdot d\mathbf{v} = \ln \left(\frac{\kappa_q^m \cdot \mathbb{I}_m(\kappa_r)}{\kappa_r^m \cdot \mathbb{I}_m(\kappa_q)} \right) \quad \text{C.80}$$

respectively, which relies on the fact that the density function integrates to one and which we obtain from Equation C.58. Adding Equations C.79 and C.80 and substituting the results of Equation C.5—specifically, the result pertaining to the first moment—we obtain Equation C.32.

In the special case where the reference distribution is the uniform distribution on the hypersphere,

¹² Equation C.77 can be negative and decreasing for some values of κ_q , with the overall value of the function then resembling a check function. Over the range of values for which $\kappa_\alpha > 1$, however, this function is purely increasing in κ_q . From Equation C.76, it is clear that $\kappa_\alpha > 1$ is satisfied if $\alpha \cdot \kappa_q > 1 + |1-\alpha| \cdot \kappa_r$, regardless of the respective mean directions of the obtained and reference distributions.

Equation C.79 simplifies to

$$\int_{\mathbb{S}^{p-1}} \kappa_q \cdot \mu_q^\top \mathbf{v} \cdot f(\mathbf{v}; \kappa_q, \mu_q) \cdot d\mathbf{v} = \kappa_q \cdot \mu_q^\top \mathbf{E}(\mathbf{Q}) \quad \text{C.81}$$

and Equation C.80 simplifies to

$$\int_{\mathbb{S}^{p-1}} \ln \left(\frac{2\pi^{m+1}}{\mathfrak{K}(\kappa_q, m) \cdot \Gamma(m+1)} \right) \cdot f(\mathbf{v}; \kappa_q, \mu_q) \cdot d\mathbf{v} = \ln \left(\frac{(\kappa_q/2)^m}{\mathbf{I}_m(\kappa_q) \cdot \Gamma(m+1)} \right) \quad \text{C.82}$$

which we obtain from Equation C.58. Adding Equations C.81 and C.82 and substituting the results of Equation C.5—specifically, the result pertaining to the first moment—we obtain Equation C.33. In particular, we rely on the fact that the inner product of a single mean direction is one.

We now turn our attention to the limiting behaviour of the Kullback-Leibler divergence with respect to an increase in κ_q given that $\kappa_q \geq \min(1, \kappa_r)$, holding $\kappa_r > 0$ fixed. Given Equation C.31 and Proposition C.1, it is perhaps unsurprising that we find that the Kullback-Leibler divergence is also an $O(\ln(\kappa_q))$ function. To show this, however, requires a slight modification of our approach as compared to the proof of Proposition C.1. Whilst the main idea remains the same (replace the logarithms of modified Bessel functions with their lower or upper bounds and characterise the limiting behaviour), the Kullback-Leibler divergence is also a function of ratios of modified Bessel functions. We must take this into account when deriving bounds. We rely extensively on Corollary C.2, and recall that the Kullback-Leibler divergence is defined as

$$d_\ell(\mathbf{Q}, \mathbf{R}) = m \cdot \ln \left(\frac{\kappa_q}{\kappa_r} \right) - \ln \left(\frac{\mathbf{I}_m(\kappa_q)}{\mathbf{I}_m(\kappa_r)} \right) + \frac{\mathbf{I}_{m+1}(\kappa_q)}{\mathbf{I}_m(\kappa_q)} \cdot (\kappa_q - \kappa_r \cdot \mu_r^\top \mu_q) \quad \text{C.83}$$

We bound the Kullback-Leibler divergence from below, as

$$d_\ell(\mathbf{Q}, \mathbf{R}) \geq m \cdot \ln \left(\frac{\kappa_q}{\kappa_r} \right) - \overline{\mathfrak{B}}(\kappa_q, \kappa_r, m) + \ln(\mathbf{I}_m(\kappa_r)) + \frac{\kappa_q^2}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}} - \frac{\kappa_q \cdot \kappa_r}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}} \quad \text{C.84}$$

$$= \underline{a}_m \cdot \ln \left(\frac{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}}{\underline{a}_m + \sqrt{\kappa_r^2 + \underline{a}_m^2}} \right) - \underline{\mathbf{H}}_\ell(\boldsymbol{\kappa}, m) \quad \text{C.85}$$

where

$$\underline{\mathbf{H}}_\ell(\boldsymbol{\kappa}, m) \doteq \sqrt{\kappa_q^2 + \underline{a}_m^2} - \sqrt{\kappa_r^2 + \underline{a}_m^2} - \frac{\kappa_q^2}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}} + \frac{\kappa_q \cdot \kappa_r}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}} \quad \text{C.86}$$

$$\leq \sqrt{\kappa_q^2 + \underline{a}_m^2} - \sqrt{\kappa_r^2 + \underline{a}_m^2} - \sqrt{\kappa_q^2 + \underline{a}_m^2} + \underline{a}_m + \kappa_r \quad \text{C.87}$$

To move from Equation C.86 to Equation C.87, we replace \underline{a}_m with \bar{a}_m where it is appropriate to do so, apply the formula for the difference of two squares (the other square being zero) to the penultimate term, and set the dimensional constants in the denominator of the final term equal to zero. It then follows that

$$d_\ell(\mathbf{Q}, \mathbf{R}) \geq \underline{a}_m \cdot \ln \left(\frac{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}}{\underline{a}_m + \sqrt{\kappa_r^2 + \underline{a}_m^2}} \right) - \kappa_r - \bar{a}_m + \sqrt{\kappa_r^2 + \underline{a}_m^2} \quad \text{C.88}$$

We bound the Kullback-Leibler divergence from above, as

$$d_\ell(\mathbf{Q}, \mathbf{R}) \leq m \cdot \ln \left(\frac{\kappa_q}{\kappa_r} \right) - \underline{\mathbb{H}}(\kappa_q, \kappa_r, m) + \ln(\mathbb{I}_m(\kappa_r)) + \frac{\kappa_q^2}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}} \quad \text{C.89}$$

$$= \underline{a}_m \cdot \ln \left(\frac{\underline{a}_m + \sqrt{\kappa_q^2 + \bar{a}_m^2}}{\underline{a}_m + \sqrt{\kappa_r^2 + \bar{a}_m^2}} \right) - \bar{\mathbb{H}}_\ell(\boldsymbol{\kappa}, m) \quad \text{C.90}$$

where

$$\bar{\mathbb{H}}_\ell(\boldsymbol{\kappa}, m) \doteq \sqrt{\kappa_q^2 + \bar{a}_m^2} - \sqrt{\kappa_r^2 + \bar{a}_m^2} - \frac{\kappa_q^2}{\underline{a}_m + \sqrt{\kappa_q^2 + \underline{a}_m^2}} \quad \text{C.91}$$

$$\geq \sqrt{\kappa_q^2 + \underline{a}_m^2} - \sqrt{\kappa_r^2 + \bar{a}_m^2} - \sqrt{\kappa_q^2 + \underline{a}_m^2} + \underline{a}_m \quad \text{C.92}$$

To move from Equation C.91 to Equation C.92, we replace \bar{a}_m with \underline{a}_m where it is appropriate to do so and apply the formula for the difference of two squares (the other square being zero) to the final term. It then follows that

$$d_\ell(\mathbf{Q}, \mathbf{R}) \leq \underline{a}_m \cdot \ln \left(\frac{\underline{a}_m + \sqrt{\kappa_q^2 + \bar{a}_m^2}}{\underline{a}_m + \sqrt{\kappa_r^2 + \bar{a}_m^2}} \right) - \underline{a}_m + \sqrt{\kappa_r^2 + \bar{a}_m^2} \quad \text{C.93}$$

Together, Equations C.88 and C.93 imply that

$$|d_\ell(\kappa_q, \kappa_r)| \leq \underline{a}_m \cdot \ln \left(\frac{\underline{a}_m + \sqrt{\kappa_q^2 + \bar{a}_m^2}}{\underline{a}_m + \sqrt{\kappa_r^2 + \bar{a}_m^2}} \right) + \left| \sqrt{\kappa_r^2 + \bar{a}_m^2} - \bar{a}_m - \kappa_r \right| \quad \text{C.94}$$

and so it is trivial to show that the Kullback-Leibler divergence is an $O(\ln(\kappa_q))$ function. \blacksquare

HANKEL EXPANSION OF THE CIRCULAR VARIANCE

One quantity that is often of interest is the circular variance of the von Mises-Fisher family of distributions, which is defined as one minus the ratio of modified Bessel functions—i.e., one minus the mean resultant length. KITAGAWA et al. (2022b) demonstrates that the circular variance is an $O(1/\kappa)$ function. Here, we show that the same result can be obtained via Hankel series expansion (see Equation B.13 for a definition). We note that Hankel series expansion is appropriate when $\kappa \rightarrow \infty$, which is the limiting behaviour that we are interested in. We begin by writing the circular variance in the alternative form

$$1 - \frac{I_{m+1}(\kappa)}{I_m(\kappa)} = 1 - \frac{I_{m+1}(\kappa)}{I_m(\kappa)} = \frac{I_m(\kappa) - I_{m+1}(\kappa)}{I_m(\kappa)} \tag{C.95}$$

to which we apply the expansion. Recalling Equation B.13 and setting $n = 2$ for simplicity,

$$1 - \frac{I_{m+1}(\kappa)}{I_m(\kappa)} = \frac{[\text{Pochhammer}_1(m+1) - \text{Pochhammer}_1(m)]/\kappa + O(1/\kappa^2)}{1 - \text{Pochhammer}_1(m)/\kappa + O(1/\kappa^2)} \tag{C.96}$$

$$= \frac{\text{Pochhammer}_1(m+1) - \text{Pochhammer}_1(m) + O(1/\kappa)}{\kappa - \text{Pochhammer}_1(m) + O(1/\kappa)} \tag{C.97}$$

which is, trivially, an $O(1/\kappa)$ function. More accurate approximations can be attained by increasing n , but a simple approximation suffices to establish the rate of the circular variance.

STOCHASTIC TREATMENT CHOICE WITH EMPIRICAL WELFARE UPDATING

The principal goal of programme evaluation is to inform the social planner as to which individuals within a target population should receive a given treatment. When treatment effects are heterogeneous in individuals' observable characteristics, the social planner can improve social welfare by implementing an individualised treatment assignment rule based upon these characteristics. The literature on statistical treatment choice initiated by MANSKI (2004b) studies how to estimate assignment rules based upon a finite sample and how to assess their welfare performance. Given an experimental or observational sample, existing approaches—including those proposed in ATHEY and WAGER (2021), HIRANO and PORTER (2009), KITAGAWA and TETENOV (2018b), and MANSKI (2004b)—yield deterministic assignment rules, which are functions mapping an individual's observable characteristics to a recommended treatment. That is, individuals who share the same observable characteristics are all assigned the same treatment. Such assignment rules are sharp and address the question of who should be treated? We adopt a broader view of the treatment choice problem by considering stochastic assignment rules that map individual observable characteristics to a probability distribution over the different treatment arms, instead addressing the question of with what probability should an individual be treated?

In static treatment choice problems with outcome distributions that exhibit the monotone likelihood ratio property, deterministic assignment rules form a class of admissible policies (KARLIN and H. RUBIN, 1956; TETENOV, 2012), such that restricting attention to this class is without loss (of welfare). Once we allow the class of outcome distributions to be unconstrained though, there is little theoretical justification for focusing on deterministic rules. In comparison to stochastic assignment rules, deterministic assignment rules have the following three potentially undesirable features. First, deterministic assignment rules cannot incorporate confidence or uncertainty about which treatment is best for each individual, with individuals typically assigned treatment if conventional point estimates suggest that treatment has a positive effect on average. The strength of evidence in support of this conclusion, usually presented in the form of confidence intervals or p-values, is not generally acted upon; what matters is whether empirical evidence supports a positive point estimate, and not whether it is sufficient or insufficient

to reject a non-positive effect. Such a sharp dichotomy of assignment is naturally overconfident in its prescription, and there is no theoretical basis for the incorporation of confidence intervals or p-values into frequentist-based decision-making. Stochastic assignment rules can represent such uncertainty, which arises due to the (finite sample) nature of experimental data or model misspecification, through their probability weighting of treatments. Second, stochastic assignment rules facilitate future evaluation, since implementing a stochastic assignment rule can generate a new experimental sample in which treatment is randomised conditional on individual observable characteristics. Third, unlike deterministic assignment rules for which the probability that a treatment is assigned changes discontinuously at some threshold, stochastic assignment rules feature assignment probabilities that smoothly change with respect to individual characteristics. Such a feature is desirable if a fairness criterion requiring that individuals with similar characteristics have similar probabilities of treatment (DWORK et al., 2012) is enforced.

This chapter proposes novel and general methods for obtaining stochastic individualised assignment rules based on randomised control trial data. Assuming that the social planner assigns individuals to a binary treatment, with her goal being to maximise additive (utilitarian) social welfare as in MANSKI (2004b), we exploit an empirical analogue of the social welfare criterion to generate individualised assignment probabilities. Specifically, we start with a prior distribution over a collection of deterministic assignment rules that we denote by \mathfrak{G} , each g of which partitions the space of individual observable characteristics into a group of characteristics and its complement. An individual is assigned treatment if their characteristics are such that $g(\mathbf{X}) = 1$, and is not assigned treatment if $g(\mathbf{X}) = 0$. We then update the prior distribution based upon an empirical analogue of the social welfare criterion to obtain a posterior distribution over \mathfrak{G} . To generate a stochastic assignment, we draw a $g \in \mathfrak{G}$ according to the posterior distribution over this collection, and implement the policy prescribed by g . In this way, the probability that an individual is treated is equal to the posterior probability that the g that is drawn is such that $g(\mathbf{X}) = 1$ (i.e., the individual is assigned treatment under several possible g , with their overall probability of being assigned treatment equal to the probability of these g being drawn).

One of the main contributions of this chapter is that we derive an optimal updating procedure for obtaining the posterior distribution over \mathfrak{G} . This procedure minimises an upper bound on welfare regret and yields an exponential tilting of the prior over \mathfrak{G} , where the tilting depends upon the empirical welfare criterion. This novel updating formula resembles the quasi-posterior distribution that appears in the Laplace-type estimation that is studied by CHERNOZHUKOV and HONG (2003b), but differs in that the constant factor in the exponential tilting term is determined endogenously by the Lagrange multiplier of the optimisation.

Despite our analytical characterisation of the optimal posterior distribution, computation of this distribution or sampling of g from it is not straightforward. We therefore consider a variational approximation of the optimal posterior distribution by a parametric distribution. In particular, as a specification of \mathfrak{G} , we consider the class of Linear Eligibility Score (LES) rules that assign treatment if $x^\top \gamma$ (the linear score) exceeds some threshold γ_0 (eligibility). Building upon the LES class, we exploit the invariance of the welfare criterion to the ratio of γ to γ_0 (scale invariance) and approximate the optimal posterior distribution using a multivariate von Mises-Fisher distribution.

For the practice of reporting and communicating individualised allocations of treatment, our approach of obtaining a posterior distribution over policies is useful for generating some quantities that existing methods yielding deterministic assignment rules cannot produce. First, the posterior probability that $g(\mathbf{X}) = 1$ offers a personalised probabilistic assessment that individuals with these characteristics favour treatment over no treatment. Reporting such a probability offers a novel alternative to the common practice of using the p-values of hypothesis testing to express confidence in a positive treatment effect, something which does not easily translate to a recommendation about what the social planner should do. Second, viewing the posterior over g as an inferential tool for the welfare-optimality of (deterministic) assignment policies within \mathfrak{G} , we can obtain a credible region for the optimal policy by, for instance, selecting its highest posterior density region. This approach to obtaining confidence sets for the optimal treatment assignment policy is an alternative to the frequentist approach that is studied in RAI (2019). Third, analogous to the practice of using a Bayesian posterior with a noninformative prior as a visual summary of the likelihood function, we can use our variationally approximated posterior over g as a visual summary of the exponentiated empirical welfare criterion function.

To demonstrate how to implement our approach and what it delivers in practice, we apply our methods to the JTPA Study sample that is studied by BLOOM et al. (1997). Given observations of prior earnings and years of education, we ask with what probability should an individual be treated? Restricting attention to linear assignment rules and a variational approximation of the optimal posterior distribution by a multivariate von Mises-Fisher distribution, we estimate a stochastic assignment rule that is more likely to allocate JTPA assistance to individuals with high prior earnings and fewer years of education. KITAGAWA and TETENOV (2018b), which similarly considers the JTPA Study sample, and estimates a deterministic assignment rule, serves as a useful benchmark for comparison. Aside from the obvious difference that we estimate a stochastic rule (i.e., every individual has a non-trivial probability to be allocated JTPA assistance under our rule), our estimated rule allocates JTPA assistance to a smaller fraction of the population than the deterministic rule of KITAGAWA and TETENOV (2018b), which targets individuals with low prior earnings and few years of education for treatment. This difference reflects the shape of

the empirical welfare criterion, which can be captured by our approach but is missed by deterministic policies that are obtained as the mode of the empirical welfare criterion.

This chapter contributes to the growing literature on statistical treatment choice initiated by MANSKI (2004b). Exact minimax regret assignment rules are studied in ISHIHARA and KITAGAWA (2021), SCHLAG (2006), STOYE (2009), STOYE (2012), TETENOV (2012), and YATA (2021). HIRANO and PORTER (2009) analyses asymptotically-optimal assignment rules in limit experiments, and BHATTACHARYA and DUPAS (2012) considers capacity constrained policies. KITAGAWA and TETENOV (2018b) proposes Empirical Welfare Maximisation (EWM) methods for individualised assignment, which maximise a sample analogue of the social welfare function over a constrained class of policies. Similar approaches have been studied in the literature on machine learning and personalised medicine, as in BEYGELZIMER and LANGFORD (2009), SWAMINATHAN and JOACHIMS (2015), ZADROZNY (2003), ZHANG et al. (2012), and ZHAO et al. (2012). Recent advances in learning individualised assignment policies include ADJAHO and CHRISTENSEN (2022), ATHEY and WAGER (2021), D'ADAMO (2021), HAN (2022), KIDO (2022), KITAGAWA and TETENOV (2021), KITAGAWA et al. (2021), LIU (2022), MBAKOP and TABORD-MEEHAN (2021), NIE et al. (2021), SAKAGUCHI (2019), SASAKI and URA (2020), SUN (2021), and VIVIANO (2021), to list but a few works. The assignment rules estimated in these works are all deterministic.

There are some earlier works that investigate the decision-theoretic justification for stochastic (fractional) assignment and the welfare performance of these rules, with MANSKI (2009) providing a detailed review of settings where stochastic rules are preferable. When the welfare criterion is only partially identified, minimax regret-optimal rules are stochastic given knowledge of the identified set (MANSKI, 2000a, 2005, 2007a,b), which remains true even after taking into account uncertainty of estimates of the bounds (MANSKI, 2022; STOYE, 2012; YATA, 2021). As shown by MANSKI and TETENOV (2007) and MANSKI, 2009, stochastic assignment rules can also be justified by a nonlinear welfare criterion in a point-identified setting. KITAGAWA et al. (2022a) shows that, for a wide class of nonlinear welfare regret criteria, admissible assignment rules are stochastic (fractional). In particular, KITAGAWA et al. (2022a) shows that the minimax squared-regret rule is stochastic, with the probability of assignment equal to the posterior probability of a positive treatment effect under the least-favorable null. KITAGAWA et al. (2022a) proposes using this probability as a measure of the strength of evidence for a positive treatment effect, replacing the commonly used p-value of a hypothesis test. In contrast, this chapter obtains the probability of assignment from a posterior probability distribution over assignment rules, rather than over the treatment effect parameters, with a quasi-likelihood built upon the empirical welfare criterion. KOCK et al. (2022) obtains a stochastic assignment rule in a setting where the oracle optimal rule is fractional due to nonlinearity in the social planner's chosen welfare criterion. In contrast to the

static treatment assignment problem, dynamic treatment assignment problems analysed in the multi-arm bandit literature often consider stochastic assignments that balance the exploitation versus exploration trade-off, such as the posterior probability matching algorithm of THOMPSON (1933) does. Thompson sampling algorithms build upon the standard Bayesian posterior distribution for treatment effects such that the allocation algorithm crucially relies on a parametric specification of the data generating process, which our approach does not require.

CHAMBERLAIN (2011) and DEHEJIA (2005) approach the treatment choice problem from a Bayesian perspective. In their framework, the potential outcome distributions are parametric, and it is over the parameters of these distributions that a prior is imposed. For the standard mean welfare criterion, the Bayes optimal allocation rule is deterministic. Our approach differs from these works in that we do not assume a prior distribution over the data generating process. We instead impose few restrictions on the data generating process, and form prior and posterior distributions over the parameters that index assignment rules. Our approach can be advantageous when compared to CHAMBERLAIN (2011) if the social planner is concerned about potential misspecification of the likelihood. If the likelihood is misspecified, the resulting Bayes-optimal assignment rule can be suboptimal even for large samples. In contrast, our approach is guaranteed to yield a distribution over policies that is guaranteed to concentrate on welfare-optimal policies without requiring a specification for the data generating process.

Our approach is also related to that of BISSIRI et al. (2016) and CSABA and SZOKE (2020), where loss function-driven (quasi-Bayes) updating rules are proposed. Rather than follow their approach by adopting exponentiated loss as a quasi-likelihood and solving the quasi-Bayesian decision problem, we obtain an optimal learning rule by minimising a high probability upper bound on welfare regret. This way of establishing optimality is similar to the structural risk minimisation approach of VAPNIK (1998) and the Probably Approximately Correct (henceforth, PAC) analysis proposed by VALIANT (1984), which was extended to the study of randomised predictors in MCALLESTER (1999), and SHAWE-TAYLOR and WILLIAMSON (1997), constituting the development of PAC-Bayes theory. For classification and regression problems, various PAC-Bayes bounds on prediction generalisation errors are obtained in BÉGIN et al. (2014, 2016), CATONI (2007), DERBEKO et al. (2004), GERMAIN et al. (2009), MCALLESTER (2003), PENTINA and LAMPERT (2015), and SEEGER (2002), and can accommodate quasi-Bayesian procedures similar to ours. See GUEDJ (2019) for a recent review of this literature. To our knowledge, the PAC-Bayes bounds that we derive for treatment choice are new to the literature and offer a contribution of independent interest. We also note that PELLATT (2022) makes use of PAC-Bayes theory to analyse the treatment allocation problem with stochastic assignment rules under a budget (or resource) constraint.

Although the treatment choice problem is distinct from prediction problems—as is discussed in KITA-GAWA and TETENOV (2018b)—the EWM approach for treatment choice is closely related to the cost-sensitive binary classification problem, as first pointed out by ZADROZNY (2003). The PAC-Bayes classification analysis with variational posterior approximation that is proposed by ALQUIER et al. (2016) is, therefore, closely related to our analysis. There are, however, important differences with ALQUIER et al. (2016). First, we make use of the approach proposed by BÉGIN et al. (2016), which allows for the construction of a variety of different bounds via a general convex function. This introduces the complication that classification is not standard (i.e., cost is homogeneous) and is instead cost-sensitive. Introducing heterogeneity in the cost of misclassification leads to a non-trivial challenge in deriving the PAC bounds. To address these complications we leverage results in MAURER (2004) for continuous loss functions over the unit interval. Second, ALQUIER et al. (2016) considers approximating the optimal posterior distribution by a Gaussian distribution. We exploit the scale invariance property of the welfare criterion and approximate the optimal posterior distribution by a multivariate von Mises-Fisher distribution over the hypersphere.

—| SECTION 1 |—

FRAMEWORK

We suppose that the social planner (or the econometrician acting on her behalf) observes experimental data that comprises a probability distribution over $\{Y, T, \mathbf{X}\}$ that we denote by P^n and that is constructed from n independent and identically distributed draws. Here, $Y \in \mathbb{R}$ denotes response, which is some measured outcome of interest; $T \in \{0, 1\}$ denotes treatment, which is an indicator for whether an individual is treated or not (i.e., which group—the experimental group or the control group—an individual is a member of in the experimental data); and $\mathbf{X} \in \mathbb{R}^i$ denotes covariates, which are some measured individual characteristics. The population from which the experimental data is drawn comprises a probability distribution over $\{Y_0, Y_1, T, \mathbf{X}\}$ that we denote by P . Here, Y_1 and Y_0 are potential outcomes and relate to Y via the relationship

$$y = y_1 \cdot t + y_0 \cdot [1 - t] \tag{D.1}$$

We refer to P^n as the empirical distribution and to P (together with Equation D.1) as the data-generating process, and assume that \mathbf{X} consists only of those characteristics that the social planner can use to discriminate between individuals in the target population, with budgetary, ethical or legal considerations precluding the use of other characteristics.

Throughout our analysis, we maintain several assumptions. We follow MANSKI (2004b) and the subse-

quent literature in supposing that the social welfare criterion is that of a utilitarian social planner who aims to maximise the average level of individual outcomes. We note that other criteria could also be implemented, such as inequality-averse social welfare and Gini social welfare.¹

Assumption D.1 (External validity). *The population to which policy is to be applied—the target population—has the same distribution over $\{Y_1, Y_0, \mathbf{X}\}$ as the marginal distribution of $\{Y_1, Y_0, \mathbf{X}\}$ that is obtained from the data generating process.*

Assumption D.2 (Unconfoundedness). *The data generating process satisfies $\{Y_1, Y_0\} \perp\!\!\!\perp T | \mathbf{X}$.*

Assumptions D.1 and D.2 are satisfied, for instance, if the experimental data is extracted directly from the target population and the treatment is, conditional on the covariates, randomly assigned,² independently of the potential outcomes (ROSENBAUM and D. B. RUBIN, 1983).

Assumption D.3 (Bounded outcomes). *There exists a constant $0 < c_y < \infty$ such that $0 \leq Y \leq c_y$.*

Assumption D.4 (Strict overlap). *There exists a constant $0 < \psi < 1/2$ such that the propensity score satisfies $\psi \leq e(x) \leq 1 - \psi$, where $e(x) \doteq \mathbb{E}_P(T | \mathbf{X})$.*

As is discussed in KITAGAWA and TETENOV (2018b) and SWAMINATHAN and JOACHIMS (2015), policies that maximise an empirical welfare criterion are not invariant to positive affine transformations of outcomes, which is the case for the empirical welfare criterion that we consider in this chapter. Given Assumptions D.3 and D.4, we can define

$$H \doteq \frac{\psi \cdot Y / c_y}{e(\mathbf{X}) \cdot T + [1 - e(\mathbf{X})] \cdot [1 - T]} \tag{D.2}$$

which is confined to the unit interval, and which we interpret as weights and that are motivated by an unbiased estimator of the (scaled) expected potential outcomes. We solve the planner’s problem using these transformed outcomes. This transformation of outcomes does not affect the welfare ranking over assignment policies, both in the population and according to in-sample welfare criteria, yet facilitates our analysis.

Adopting an additive utilitarian perspective, the average level of social welfare attained by g is propor-

1 For example, KASY (2016) and KITAGAWA and TETENOV (2021) study a setting where the social welfare function is a weighted average of the outcomes with rank-dependent weights, which includes the Gini social welfare function as a special case.

2 KITAGAWA and TETENOV (2018b) considers a setting where the marginal distribution of \mathbf{X} differs between the population of interest and the data generating process. ADJAHO and CHRISTENSEN (2022) and KIDO (2022) study settings that differ also in terms of the distribution of potential outcomes.

tional to

$$W(g) \doteq \mathbb{E}_P(Y_1 \cdot 1(g(\mathbf{X}) = 1) + Y_0 \cdot 1(g(\mathbf{X}) = 0)) \quad \text{D.3}$$

Given Assumptions D.1, D.2 and D.4 and that $Y = Y_0 + T \cdot (Y_1 - Y_0)$, we can re-write Equation D.3 as

$$W(g) = \mathbb{E}_P\left(\frac{Y \cdot T}{e(\mathbf{X})} \cdot 1(g(\mathbf{X}) = 1) + \frac{Y \cdot (1 - T)}{1 - e(\mathbf{X})} \cdot 1(g(\mathbf{X}) = 0)\right) \quad \text{D.4}$$

$$= \mathbb{E}_P\left(\frac{Y \cdot [1 - T]}{1 - e(\mathbf{X})}\right) + \mathbb{E}_P\left(\left[\frac{Y \cdot T}{e(\mathbf{X})} - \frac{Y \cdot [1 - T]}{1 - e(\mathbf{X})}\right] \cdot g(\mathbf{X})\right) \quad \text{D.5}$$

Accordingly, the sample analogue of Equation D.3 can be written as

$$W_n(g) \doteq \frac{1}{n} \sum_{i=1}^n \left[\frac{y_i \cdot t_i}{e(\mathbf{x}_i)} \cdot 1(g(\mathbf{x}_i) = 1) + \frac{y_i \cdot [1 - t_i]}{1 - e(\mathbf{x}_i)} \cdot 1(g(\mathbf{x}_i) = 0) \right] \quad \text{D.6}$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{y_i \cdot [1 - t_i]}{1 - e(\mathbf{x}_i)} + \frac{1}{n} \sum_{i=1}^n \left[\frac{y_i \cdot t_i}{e(\mathbf{x}_i)} - \frac{y_i \cdot [1 - t_i]}{1 - e(\mathbf{x}_i)} \right] \cdot g(\mathbf{x}_i) \quad \text{D.7}$$

where $W_n(g)$ is an unbiased estimator for the true level of welfare that arises from the implementation of a particular g . Given the additive social welfare criterion, the maximal welfare level can be attained by a deterministic policy. Hence, as far as the population welfare maximisation problem is concerned, the social planner wants to select the g that maximises $W(g)$.

Each g can be associated with the set of characteristics that it assigns treatment to. We refer to g as a deterministic assignment rule, or simply as an assignment rule, with \mathfrak{G} constituting the class of assignment rules. We can write Equation D.7 as

$$W_n(g) = \frac{1}{n} \sum_{i=1}^n y_i \cdot \left[\frac{t_i}{e(\mathbf{x}_i)} \cdot g(\mathbf{x}_i) + \frac{1 - t_i}{1 - e(\mathbf{x}_i)} \cdot [1 - g(\mathbf{x}_i)] \right] \quad \text{D.8}$$

$$= \frac{1}{n} \sum_{i=1}^n h_i \cdot \frac{c_y}{\psi} \cdot [t_i \cdot g(\mathbf{x}_i) + [1 - t_i] \cdot [1 - g(\mathbf{x}_i)]] \quad \text{D.9}$$

$$= \frac{1}{n} \sum_{i=1}^n h_i \cdot \frac{c_y}{\psi} \cdot 1(g(\mathbf{x}_i) = t_i) \quad \text{D.10}$$

$$= \frac{1}{n} \sum_{i=1}^n h_i \cdot \frac{c_y}{\psi} - \frac{1}{n} \sum_{i=1}^n h_i \cdot \frac{c_y}{\psi} \cdot 1(g(\mathbf{x}_i) \neq t_i) \quad \text{D.11}$$

where h_i is the realisation of H , as defined in Equation D.2, for a given individual. We observe from Equation D.11 that $W_n(g)$ depends upon g only through its second term, such that

$$\operatorname{argmax}_g W_n(g) = \operatorname{argmin}_g \frac{1}{n} \sum_{i=1}^n h_i \cdot 1(g(\mathbf{x}_i) \neq t_i) \quad \text{D.12}$$

Accordingly, we define

$$R(g) \doteq \mathbb{E}_P(H \cdot \mathbf{1}(g(\mathbf{X}) \neq T)) \tag{D.13}$$

which we term the welfare risk of g , and its empirical analogue

$$R_S(g) \doteq \frac{1}{n} \sum_{i=1}^n h_i \cdot \mathbf{1}(g(\mathbf{x}_i) \neq t_i) \tag{D.14}$$

which we term the empirical welfare risk of g . In view of Equation D.12, the social planner's objective is to minimise Equation D.13 with respect to g , via Equation D.14, following the empirical risk minimisation principle of VAPNIK (1998).

One special set of policies that we draw particular attention to is the LES class that is defined in KITAGAWA and TETENOV (2018b), and that we denote by \mathfrak{L} . Assignment rules in this class are indexed by a finite-dimensional parameter vector γ and a threshold γ_0 , and are associated with a binary function ℓ_β that satisfies

$$\ell_\beta(\mathbf{X}) \doteq \mathbf{1}(\mathbf{X}^\top \gamma \geq \gamma_0) \tag{D.15}$$

where we take β to include both γ and γ_0 (i.e., β is an m -dimensional vector). Each LES rule induces a partitioning of the covariate space into two half-spaces, such that individuals in the upper contour set receive treatment and individuals in the lower contour set do not. By restricting β to the unit hypersphere, we guarantee that each policy is associated with a unique β . In what follows, we exploit the interchangeability of β and the LES rule that it indexes, adopting β as the argument of the loss functions that we consider. For instance, and again with some abuse of notation, whenever we focus on the LES class of decision rules we write

$$R(\beta) = \mathbb{E}_P(H \cdot \mathbf{1}(\ell_\beta(\mathbf{X}) \neq T)) \tag{D.16}$$

and

$$R_S(\beta) = \frac{1}{n} \sum_{i=1}^n h_i \cdot \mathbf{1}(\ell_\beta(\mathbf{x}_i) \neq t_i) \tag{D.17}$$

respectively, in place of Equations D.13 and D.14.

POSTERIOR DISTRIBUTIONS AND STOCHASTIC ASSIGNMENT RULES

We now adapt Equations D.13 and D.14 to handle stochastic assignment rules. We let Π denote a probability distribution over \mathfrak{G} that we interpret as a posterior distribution, assuming that \mathfrak{G} can be embedded in a measurable space.³ We let \mathfrak{P} denote the collection of all posterior distributions.

Definition D.1 (Posterior assignment rule). *Let Π be a probability distribution over \mathfrak{G} that is constructed upon observing the sample. The posterior assignment rule under Π is a stochastic assignment rule that assigns treatment to individuals with probability $Q^\Pi(\mathbf{X}) \doteq \int_{\mathfrak{G}} g(\mathbf{X}) \cdot d\Pi$.*

To implement posterior assignment rules in practice, we randomly draw a g from \mathfrak{G} according to Π for each individual in the target population. In this way, similar individuals, who can have similar assignment probabilities, can be assigned to different treatment arms. Moreover, this approach does not require computation of the probability of treatment.

Definition D.2 (Expected welfare risk under Π). *We define the expected welfare risk under Π as*

$$R^\Pi \doteq \int_{\mathfrak{G}} R(g) \cdot d\Pi(g) = E_P(H \cdot [T \cdot [1 - Q^\Pi(\mathbf{X})] + [1 - T] \cdot Q^\Pi(\mathbf{X})]) \quad \text{D.18}$$

with its empirical analogue taking the form

$$R_S^\Pi \doteq \int_{\mathfrak{G}} R_S(g) \cdot d\Pi(g) = \frac{1}{n} \sum_{i=1}^n h_i \cdot [t_i \cdot [1 - Q^\Pi(\mathbf{x}_i)] + [1 - t_i] \cdot Q^\Pi(\mathbf{x}_i)] \quad \text{D.19}$$

The interpretation of R^Π is the average welfare loss that the social planner expects from stochastic implementation of g in \mathfrak{G} in the target population when g is distributed according to Π .

We reiterate that stochastic assignment is achieved by randomly drawing g according to Π . This way of selecting assignment rules is reminiscent of the Gibbs classifier in statistical learning theory (GERMAIN et al., 2009) and might offer a computational advantage over other methods if drawing g according to Π is easier than maximising empirical welfare or finding the mode of Π , say. An advantage of stochastic assignment is the possibility for sequential treatment evaluation: the induced assignment of treatment and non-treatment to individuals in the target population by Π is random conditional on \mathbf{X} , which allows for estimation of the causal effect of treatment in future studies. In this sense, stochastic assignment is well suited to balancing existing evidence about what constitutes the optimal assignment for each

³ For \mathfrak{G} to be embedded in a measurable space, \mathfrak{G} cannot be too rich. We defer to MOLCHANOV (2005) and GUNSILIUS (2019) for further discussion of this point.

individual against the benefit of further exploration of the treatment effect (MANSKI, 2000b). We do not, however, study this channel and focus on a purely static problem in this chapter.

Our framework allows the social planner to hold some prior as to what constitutes the best policy. We differentiate the subjective beliefs that the social planner holds, which we encode using the prior distribution Π_0 , and their updated beliefs following their observation of sample data, which we encode using Π . In the analysis that follows, we provide finite sample regret guarantees in the form of PAC-Bayes bounds for stochastic assignment. The approach discussed here differs significantly from other Bayesian treatment choice settings, such as those discussed in CHAMBERLAIN (2011), because we do not impose any kind of prior belief on P . We instead choose to model the beliefs that the social planner has regarding the optimal policy. Using a decision procedure that is free from specification of the likelihood comes with a desirable robustness property, as we discuss in due course.

Aside from its more conventional role as a means of expressing existing information about what constitutes the best policy, Π_0 can also play several other roles within our framework. For instance, Π_0 can also embed any constraints that are imposed upon the set of policies through truncation of its support. Such a zero density condition can be imposed in lieu of an explicit restriction on \mathfrak{G} and is easy to implement in practice via rejection sampling, with only those policies that satisfy any budgetary, ethical or legal constraints being retained under the sampling procedure. Moreover, Π_0 can also be used to describe the status quo, with restrictions on the shape of Π_0 governing how much policy can deviate. These interpretations of Π_0 naturally extend to Π .

—| SECTION 3 |—

OPTIMAL STOCHASTIC ASSIGNMENT AND CONVERGENCE OF WELFARE

Having defined a criterion by which to assess stochastic assignment rules, we ask what is the optimal posterior assignment rule? Of course, answering this question requires that we make some link between the performance of an assignment rule in the sample versus in the target population, and that we make some judgement about the feasibility of estimating this rule.

—| SUBSECTION 3.A |—

BOUNDING EXPECTED WELFARE RISK

Seeing as experimental data provides only an insight into the welfare performance of any policy in the target population, a very natural question to ask is how much we can expect R_S^Π to differ from R^Π for any given Π . We provide an answer to this question here.

Theorem D.1. *Suppose that Assumptions D.1 to D.4 are satisfied, and that $n \geq 8$. Then, with at least probability $1 - \varepsilon$ under P^n , for all ε satisfying $0 < \varepsilon < 1$,*

$$R^\Pi \leq R_S^\Pi + \sqrt{\frac{1}{2n} \cdot \left[d_\ell(\Pi, \Pi_0) + \ln \left(\frac{2\sqrt{n}}{\varepsilon} \right) \right]} \tag{D.20}$$

provided that $\Pi \in \mathfrak{P}$ is absolutely continuous with respect to Π_0 .

We present Theorem D.1 without proof (see KITAGAWA et al., 2022b for a proof of Theorem D.1).⁴

We note that Theorem D.1 holds for any rules that update Π_0 and deliver Π , and that we have not committed to any particular updating rule to obtain Equation D.20. A feature of Equation D.20 is that the regularisation term—the square root term containing the Kullback-Leibler divergence of Π from Π_0 —enters additively, which we find is a convenient feature for establishing convergence of our variational approximation (regularisation can otherwise be effected multiplicatively).⁵ The regularisation term, by design, prevents overfitting. To illustrate this point, suppose that Π_0 is uniform over \mathfrak{G} : the best response of the social planner absent regularisation is to concentrate probability mass on the optimal (in-sample) g as suggested by data, such that Π is degenerate. In the presence of regularisation, however, this is no longer a best response, since the Kullback-Leibler divergence infinitely penalises degeneracy vis-à-vis uniformity. Rather, the best response of the social planner is to allocate probability mass on all g in \mathfrak{G} , albeit concentrating more mass on those g that are associated with low empirical welfare risk. Put differently, the regularisation term controls how far away from Π_0 a stochastic assignment rule can be, with this difference governed by the number of observations in the sample.

The Vapnik-Chervonenkis (VC) dimension is the standard measure of complexity in the statistical learning literature. We instead associate complexity with the Kullback-Leibler divergence. If one is willing to impose distributional constraints on a posterior assignment rule, an advantage of the PAC-Bayes approach and of using the Kullback-Leibler divergence is that complexity is then purely in terms of the selected (stochastic) assignment rule Π rather than in terms of the class of possible stochastic assignment rules \mathfrak{P} or the class of underlying deterministic rules \mathfrak{G} . As such, Theorem D.1 does not explicitly require any assumption about the VC dimension of \mathfrak{G} . The influence of VC dimension for \mathfrak{G} is implicit in our setting: the upper bound on the difference between R_S^Π and R^Π implied by Equation D.20 is governed

⁴ The proof builds upon BÉGIN et al. (2016, §Lemma 3), which offers a flexible approach that allows for the recovery of many different PAC-Bayes bounds. The approach centres around a convex function of R^Π and R_S^Π , which we specify so as to recover the form that is presented in MCALLESTER (2003, see §(8) in BÉGIN et al., 2016). We leverage results in MAURER (2004, §Lemma 3 and §Theorem 1), exploiting the properties of Bernoulli random variables and convex functions, to adapt BÉGIN et al. (2016) and the bound that is presented therein to allow for heterogeneous cost (in lieu of the standard binary loss function that is prevalent in the classification literature and that is studied in BÉGIN et al., 2016).

⁵ See BÉGIN et al. (2016) for the implications of different convex functions.

by the Kullback-Leibler divergence, which is increasing in the dimension of the support of Π and Π_0 , and so is non-decreasing in the complexity of \mathfrak{G} .

—| SUBSECTION 3.B |—

OPTIMAL UPDATING RULE

In a standard Bayesian setting, unknown parameters index the distribution of data and inference on parameters is conducted with respect to the posterior distribution. Typically, the posterior distribution is constructed from a well-defined likelihood function via Bayes’ theorem. We leverage Theorem D.1 to construct Π from Π_0 and R_S^Π . This approach is valid since Theorem D.1 holds for all $\Pi \in \mathfrak{P}$. We emphasise that the posterior distribution that we construct is over assignment rules rather than over the data generating processes, which is distinct from BISSIRI et al. (2016) and CSABA and SZOKE (2020).

Following MCALLESTER (2003) and GERMAIN et al. (2009), we define an optimal posterior distribution, which we denote by Π^* , as a distribution over \mathfrak{G} that minimises the right hand side of Equation D.20. That is, Π^* minimises

$$R_S^\Pi + \sqrt{\frac{1}{2n} \cdot \left[d_\ell(\Pi, \Pi_0) + \ln \left(\frac{2\sqrt{n}}{\varepsilon} \right) \right]} \tag{D.21}$$

Theorem D.2. *The optimal posterior Π^* over \mathfrak{G} satisfies*

$$d\Pi^*(g) \doteq \frac{\exp(-\chi \cdot R_S(g))}{\int_{\mathfrak{G}} \exp(-\chi \cdot R_S(g)) \cdot d\Pi_0(g)} \cdot d\Pi_0(g) \tag{D.22}$$

where

$$\chi \doteq 4n \cdot \sqrt{\frac{1}{2n} \cdot \left[d_\ell(\Pi, \Pi_0) + \ln \left(\frac{2\sqrt{n}}{\varepsilon} \right) \right]} \tag{D.23}$$

We present a proof of Theorem D.2 in the appendices.

The posterior distribution that we derive is analogous to the optimal posterior in MCALLESTER (2003) with the difference that our observations are mapped to the unit interval rather than to positive-negative, which is the standard support in classification. This particular distribution is common in the statistical mechanics literature and is a Boltzmann (or Gibbs) distribution, and has the form of exponential tilting of the prior, where the exponential tilting term involves the negative empirical welfare risk. The degree of tilting depends upon the magnitude of $\chi > 0$, which is the inverse of the Lagrange multiplier of the

associated minimisation problem and that corresponds to the root of Equation D.23. The Lagrange multiplier controls the extent to which Π_0 is updated by empirical welfare risk in minimising the upper bound for R^Π .

Although Theorem D.2 offers an analytical characterisation of the optimal posterior, we are unable to obtain a closed-form expression for the optimal posterior density. Whilst Equation D.23 does suggest a means to compute this density, it is likely that, in practice, this computation is difficult to perform with any precision. Given that we have, however, established that this density exists, we can consider approximating it using the variational approximation method, as is considered in ALQUIER et al. (2016).

—| SUBSECTION 3.C |—

VARIATIONAL APPROXIMATION OF THE OPTIMAL STOCHASTIC ASSIGNMENT RULE

We develop a variational approximation of the optimal posterior density. Variational approximation is useful in situations where Gibbs distributions are difficult to sample from directly, such as is the case for graphical models where Markov Chain Monte Carlo (MCMC) sampling is costly (WAINWRIGHT and JORDAN, 2008). In variational approximation, we choose to approximate the optimal posterior distribution via a family of distributions of our choice, $\mathfrak{V} \subset \mathfrak{P}$. This allows us to develop an analytically tractable upper bound for the welfare regret attained by the resulting stochastic assignment rule.

Aside from guaranteeing tractability in estimation, variational approximation can also be motivated as a convenient way to impose constraints on the set of policies. For instance, a fairness criterion requiring that individuals with similar characteristics have similar probabilities of treatment (DWORK et al., 2012) can be enforced by specifying that \mathfrak{V} is a continuous family, and by limiting the concentration of the density. A similar approach can be used if \mathfrak{V} is a parametric family to limit how much policy can deviate from the status quo, by fixing the parameters of the posterior distribution or restricting them to some set, say.

We approximate Π^* (implicitly defined in Theorem D.2) by minimising the right-hand side of Equation D.20 with respect to posterior distributions in \mathfrak{V} , defining

$$\tilde{\Pi} \doteq \operatorname{argmin}_{\Pi \in \mathfrak{V}} \left(R_S^\Pi + \sqrt{\frac{1}{2n} \cdot \left[d_\ell(\Pi, \Pi_0) + \ln \left(\frac{2\sqrt{n}}{\varepsilon} \right) \right]} \right) \tag{D.24}$$

We then use this optimal variational posterior distribution to define a new bound for welfare regret from which we can characterise its convergence rate.

Lemma D.1. *With probability $1 - \varepsilon$ under P^n , for all ε such that $0 < \varepsilon < 1$, the expected welfare risk under the optimal variational posterior satisfies*

$$R^{\tilde{\Pi}} \leq \inf_{\Pi \in \mathfrak{V}} \left(R^{\Pi} + \frac{a(\lambda, n)}{\lambda} + \frac{d_{\ell}(\Pi, \Pi_0)}{\lambda} + \frac{\ln(2/\varepsilon)}{\lambda} + \sqrt{\frac{1}{2n} \cdot \left[d_{\ell}(\Pi, \Pi_0) + \ln\left(\frac{4\sqrt{n}}{\varepsilon}\right) \right]} \right) \quad \text{D.25}$$

where $\lambda > 0$ is an arbitrary constant.

In what follows, we let \mathfrak{V} be a variational family of distributions that assigns positive density to assignment rules in \mathfrak{L} only. Equivalently, we let \mathfrak{V} be a directional family that assigns positive density to unit vectors on the hypersphere, recalling that every policy in the LES class can be uniquely associated with a unit vector (see Equation D.15 and surrounding discussion). A particularly tractable directional family, and one that we use, is the von Mises-Fisher family of distributions, which is characterised by a probability density function satisfying, for all $\kappa > 0$ and m -dimensional unit vectors μ ,

$$d\Pi(\beta; \kappa, \mu) \doteq \frac{\kappa^{m/2-1} \cdot \exp(\kappa \cdot \mu^{\top} \beta)}{[2\pi]^{m/2} \cdot I_{m/2-1}(\kappa)} \cdot d\beta \quad \text{D.26}$$

We refer the reader to Chapter C for further details about the von Mises-Fisher family of distributions.

We now introduce some further notation that facilitates our analysis. First, we let $\bar{R} \doteq \inf_{\Pi \in \mathfrak{P}} R^{\Pi}$ denote the minimum expected welfare risk amongst assignment rules in \mathfrak{L} and $\bar{\beta}$ denote the vector that parametrises the assignment rule that induces \bar{R} . We also add the following assumption that restricts the marginal distribution of \mathbf{X} .

Assumption D.5 (Margin assumption). *There exists a constant $\eta > 0$ such that, for all m -dimensional unit vectors β^{\dagger} and β^{\ddagger} ,*

$$P(\text{sign}(\mathbf{X}^{\top} \beta^{\dagger}) \neq \text{sign}(\mathbf{X}^{\top} \beta^{\ddagger})) \leq \eta \cdot \|\beta^{\dagger} - \beta^{\ddagger}\|_2 \quad \text{D.27}$$

Assumption D.5 is satisfied whenever \mathbf{X} has bounded density on the unit hypersphere and is also present in the analysis of ALQUIER et al. (2016). An interpretation is that the proportion of individuals whose treatment status switches is continuous with respect to the LES, with

$$R(\beta) - \bar{R} \leq 2\eta \cdot \|\beta - \bar{\beta}\|_2 \quad \text{D.28}$$

being an implication of Assumption D.5.

We now present a high-probability uniform upper bound for the welfare regret of the stochastic assignment rule obtained by variational approximation via the von Mises-Fisher family of distributions.

Theorem D.3. *Suppose that Assumptions D.1 to D.5 are satisfied, that Π_0 is a uniform distribution over the unit hypersphere, that Π is a von Mises-Fisher distribution, and that $n \geq 8$. Then, with probability at least $1 - \varepsilon$ under P^n ,*

$$R^{\tilde{\Pi}} - \bar{R} \leq \text{universal constant} \cdot \frac{\ln(n)}{\sqrt{n}} \tag{D.29}$$

where the universal constant is a function of ε .

We present proof of Theorem D.3 and provide an analytical expression for the universal constant in the appendices.

The uniform upper bound on welfare regret that is defined by Equation D.29 decays at a rate of $\ln(n)/\sqrt{n}$. This rate is slightly slower than the welfare regret convergence rate of the EWM (deterministic) assignment rule studied in KITAGAWA and TETENOV (2018b). A simple comparison of these rates, however, is not quite meaningful for the following reason: we do not know if the convergence rate of $\ln(n)/\sqrt{n}$ that is obtained in Theorem D.3 is sharp or not. Theorem D.3 requires Assumption D.5, whilst KITAGAWA and TETENOV (2018b) does not impose this assumption in showing that $1/\sqrt{n}$ is the minimax optimal rate of welfare regret convergence.⁶ We do not know what the minimax optimal rate of welfare regret convergence is when Assumption D.5 is additionally imposed and, hence, cannot rule out the possibility that the convergence rate of Theorem D.3 can be improved upon and made faster than $1/\sqrt{n}$. We leave further investigation of this matter for future research.

It is worth emphasising that the regret convergence result of Theorem D.3 imposes weak restrictions on the distribution of data (Assumptions D.3 - D.5) and does not require the specification of a likelihood function or of regression equations. This contrasts with other approaches such as a Bayesian approach, where misspecification of likelihood can lead to non-convergence of the welfare regret even when Bayes optimal policies are constrained to deterministic ones in \mathfrak{L} .

Our approach is similar to ALQUIER et al. (2016) in which the families of distributions for variational approximation are multivariate Gaussian distributions on Euclidean space with flexible covariance ma-

⁶ KITAGAWA and TETENOV (2018b, §Theorem 2.3 and §Theorem 2.4) establishes that the minimax-optimal rate under a stronger condition than our Assumption D.5 is $1/n^{2/3}$. This stronger condition embeds a margin assumption that implies our Assumption D.5 but also embeds the requirement that the first-best treatment rule is contained in the set of admissible decision rules—that \mathfrak{L} contains the deterministic assignment rule that minimises expected welfare risk amongst all assignment rules. We do not make any assumption about whether \mathfrak{L} (or indeed \mathfrak{G} in earlier parts of our analysis) contains the first-best assignment rule, or whether this rule is deterministic or stochastic.

trices. In our approach, we stipulate a class of von Mises-Fisher distributions. Since the empirical welfare criterion for LES rules is invariant to the scale of β , it is natural to consider von Mises-Fisher distributions rather than Gaussian ones. The scale invariance of von Mises-Fisher distributions can simplify optimisation of the variational approximation by reducing the set of optima to a singleton.

—| SECTION 4 |—
IMPLEMENTATION

To implement our procedure, we restrict attention to \mathfrak{L} and let \mathfrak{V} be the von Mises-Fisher family of distributions. Our goal is then to minimise the objective function,

$$R_S^\Pi + \sqrt{\frac{1}{2n} \cdot \left[d_\ell(\Pi, \Pi_0) + \ln \left(\frac{2\sqrt{n}}{\varepsilon} \right) \right]} \quad \text{D.30}$$

with respect to κ and μ , which are the parameters of our chosen variational family (the concentration parameter and the mean direction, respectively). Here, we reiterate that ε relates to the probability with which our high probability bounds hold. We assume that Π_0 is the uniform distribution over the sphere and set ε equal to the 5% level throughout.

We propose numerically minimising the objective function, approximating R_S^Π using Monte Carlo draws of β from the von Mises-Fisher distribution for a given realisation of data and for fixed values of the parameters of the von Mises-Fisher distribution. We let β^j be one such pseudo-random draw that we obtain, and compute

$$\hat{R}_S^\Pi \doteq \frac{1}{n} \sum_{i=1}^n h_i \cdot [t_i \cdot [1 - \hat{\Pi}_i] + [1 - t_i] \cdot \hat{\Pi}_i] \quad \text{D.31}$$

where

$$\hat{\Pi}_i \doteq \frac{1}{J} \sum_{j=1}^J 1(\mathbf{x}_i^\top \beta^j \geq 0) \quad \text{D.32}$$

Fast pseudo-random sampling of von Mises-Fisher random vectors is possible using the rejection sampling method of WOOD (1994) or the inversion method of KURZ and HANEBECK (2015). The analogue of Equation D.30 that we then minimise is

$$\hat{R}_S^\Pi + \sqrt{\frac{1}{2n} \cdot \left[\left[\frac{m}{2} - 1 \right] \cdot \ln \left(\frac{\kappa}{2} \right) - \ln \left(I_{m/2-1}(\kappa) \right) - \ln \left(\Gamma \left(\frac{m}{2} \right) \right) + \frac{I_{m/2}(\kappa)}{I_{m/2-1}(\kappa)} \cdot \kappa + \ln \left(\frac{2\sqrt{n}}{\varepsilon} \right) \right]} \quad \text{D.33}$$

which follows from our assumption of a uniform prior and the form that the Kullback-Leibler divergence takes under that assumption.

As is shown in Chapter C, the Kullback-Leibler divergence of the von Mises-Fisher distribution from the uniform distribution over the hypersphere does not depend upon μ and increases at the logarithmic rate in the concentration. In contrast, although R_S^Π is a function of the parameters of the von Mises-Fisher distribution, the nature of this relationship is not clear. That several combinations of these parameters induce local minima of R_S^Π cannot be ruled out for instance. Moreover, since we construct \hat{R}_S^Π based upon random draws of β from Π , our objective function is not a smooth function of the concentration nor the mean direction. Given this, we suggest that a grid-based search for the minimum of the objective function is appropriate, with this search limited to values of the concentration parameter between zero and some specified upper limit. Further insight about the behaviour of the objective function is provided in an online appendix.

—| SECTION 5 |—

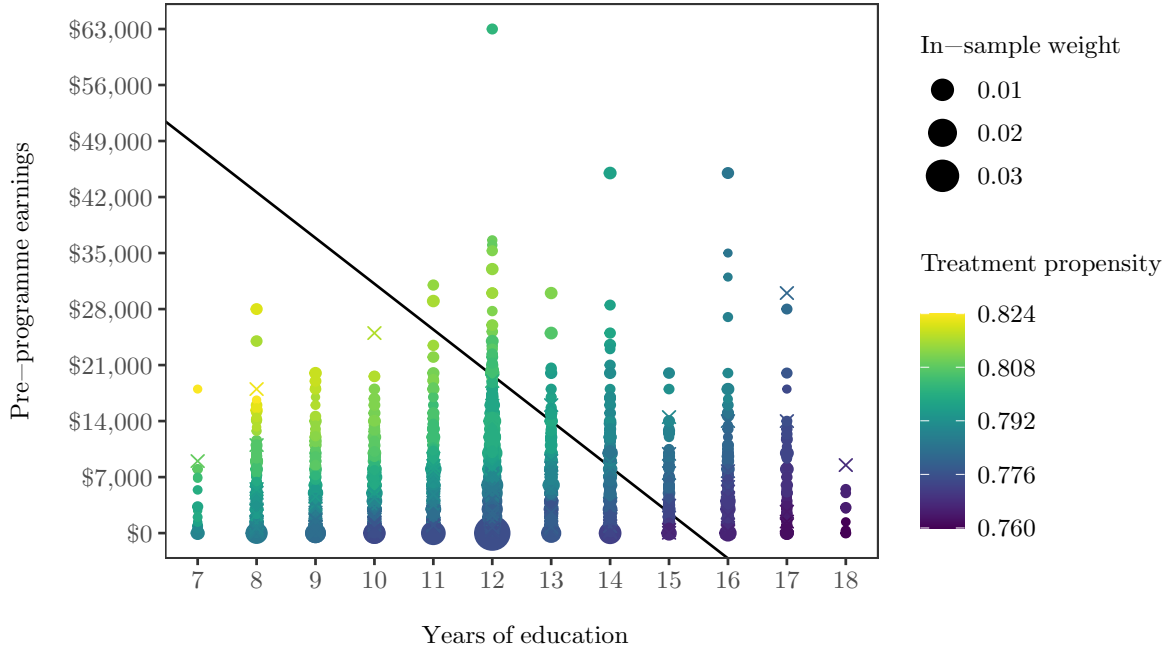
EMPIRICAL ILLUSTRATION

We illustrate our procedure using data from the National Job Training Partnership Act (JTPA) Study. Applicants to the Study were randomly allocated to one of two groups. Applicants allocated to the treatment group were extended training, job search assistance and other services provided by the JTPA over a period of 18 months. Applicants allocated to the control group were excluded from JTPA assistance. Along with information collected prior to the commencement of the intervention, the Study also collected administrative and survey data relating to applicants' earnings in the 30 months following its start. Further details about the data and the Study can be found elsewhere (see, for instance, BLOOM et al., 1997). We restrict attention to a sample of 9,223 observations for which data on years of education and pre-programme earnings amongst the sample of adults (aged 22 years and older) used in the original evaluation of the programme and in subsequent studies (ABADIE et al., 2002; BLOOM et al., 1997; HECKMAN et al., 1997) is available. Applicants in this sample were assigned to the treatment group with a probability of two thirds. Like KITAGAWA and TETENOV (2018b), we define T to be the initial assignment of treatment, rather than the actual take-up due to the presence of non-compliance in the experiment. We study stochastic assignment rules.

We follow KITAGAWA and TETENOV (2018b) in considering total individual earnings in the 30 months after programme assignment as our principal welfare measure. Moreover, we focus exclusively on the class of linear rules,

$$\begin{aligned} \mathcal{L} &= \{g : g(\mathbf{X}) = 1(\beta_0 + \beta_1 \cdot x_1 + \beta_2 \cdot x_2 \geq 0) \mid \|\{\beta_0, \beta_1, \beta_2\}\|_2 = 1\} \\ &= \{g : g(\mathbf{X}) = 1(\beta_0 + \beta_1 \cdot \text{prior earnings} + \beta_2 \cdot \text{education} \geq 0) \mid \|\{\beta_0, \beta_1, \beta_2\}\|_2 = 1\} \end{aligned} \tag{D.34}$$

FIGURE D.1
Variation in treatment propensity across individuals in the JTPA Study sample



JTPA Study sample. This figure illustrates the treatment propensity of individuals under the posterior assignment rule that is induced by $\{\kappa^, \mu^*\}$. Each point represents the individual characteristics of an individual or several individuals in the (crosses denote individuals with zero in-sample weight). For comparison, individuals to the left of the solid diagonal line are assigned treatment under the optimal deterministic assignment rule of KITAGAWA and TETENOV (2018b).*

that are studied in that paper.

To implement our procedure, we map prior earnings and education to the unit interval,⁷ and calculate H as outlined in Equation D.2. We perform this calculation without adjusting post-programme earnings by the average cost of JTPA assistance (\$774 per individual) for treated individuals.⁸ We then utilise a grid search approach over the parameters of the von Mises-Fisher distribution, specifying a reasonably fine grid over the unit sphere and over a finite subset of the reals.⁹

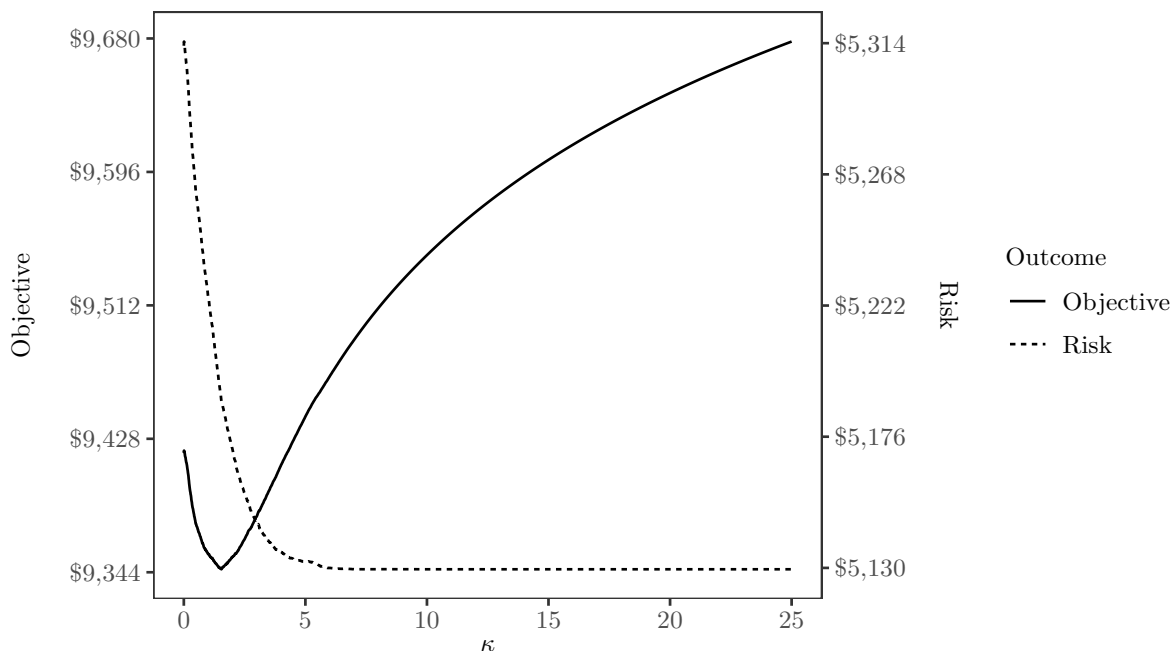
For each point on our grid, we draw 1,000 values of β from the corresponding von Mises-Fisher distribution and approximate empirical welfare risk as per Equation D.31. We then substitute these values into Equation D.33 to provide an estimate of the objective function.

⁷ We map each variable to the unit interval by dividing through by its maximum in the sample. KITAGAWA and TETENOV (2018b) also does this. Such a change of units is useful when the domain of one variable is much larger than the domain of another and the respective coefficients on the two variables reflect this. For instance, in our sample, every individual has between seven and 18 years of education, and no individual earned more than \$63,000 prior to the start of the intervention.

⁸ We adjust post-programme earnings by the average cost of JTPA assistance in an online appendix.

⁹ We design our grid so as to place an upper limit on the great-circle distance between any point on the sphere and its closest point on the grid. Our grid comprises a total of 10,116 directional vectors combined with a sequence of evenly-spaced concentrations on the zero to five interval. For reference, the surface area of the sphere is 4π , which means that our grid has a density of approximately 0.001.

FIGURE D.2
Behaviour of the objective function at μ^* given variation in κ



JTPA Study sample. This figure illustrates the shape of the objective function and its risk component at μ^* as κ is varied; high concentration is associated with low risk but incurs a large penalty for divergence from the uniform prior.

We find that the objective function is minimised (amongst the class of von Mises-Fisher distributed linear assignment rules) by the stochastic assignment rule with $\kappa = 1.550$ and $\mu = \{0.883, 0.442, 0.158\}$,¹⁰ which we label as κ^* and μ^* , respectively.

The optimal stochastic assignment rule, on average, assigns treatment to individuals in the JTPA Study sample around 83% of the time. This probability is not, however, uniform, and there is some variation in the probability with which distinct individuals are assigned treatment. This variation in assignment propensity can be seen in Figure D.1, which plots the individual characteristics of all individuals in the sample. The propensity with which individuals with distinct characteristics are assigned treatment is represented by the color of each point, and the weight given to individuals in the sample with particular characteristics is represented by the size of each point. The weight attached to a given point is proportional to the sum of post-programme earnings over all individuals with those characteristics.¹¹ Figure D.1 shows that the optimal stochastic assignment rule is more likely to assign treatment to an individual with few years of education and high prior earnings than an individual with more years of education and lower prior earnings, with the assignment probability ranging from 78% to 84%. The deterministic assignment rule of KITAGAWA and TETENOV (2018b), in contrast, assigns only individuals

¹⁰ This directional vector can be represented by an azimuth of 27° and an inclination of 81° using spherical coordinates.

¹¹ To simplify Figure D.1, we scale the weights such that they sum to one.

with few years of education and low prior earnings to treatment, with around 93% of individuals assigned treatment. We plot this deterministic rule as a useful benchmark for comparison in Figure D.1.

It is important to emphasise that it is the regularisation term and, in particular, the Kullback-Leibler divergence that limits the value of the concentration parameter at the optimum, and leads to an interior probability of assignment for all individuals. This can be seen in Figure D.2, which plots the USD equivalent of the objective function (left-hand axis, solid line) and of empirical welfare risk (right-hand axis, dashed line) for a range of values of the concentration parameter, holding fixed $\mu = \mu^*$. We observe that empirical welfare risk decreases as the value of the concentration parameter increases, remaining low and constant once its value is sufficiently large.¹² The intuition here is that large values of the concentration parameter lead to stochastic assignment rules that mimic deterministic ones; we expect the mean direction to eventually coincide with the deterministic assignment rule of KITAGAWA and TETENOV (2018b) as the value of the concentration parameter approaches infinity, since there does not exist a (linear) deterministic rule that can improve upon this. Tempering this preference towards large values of the concentration parameter is the Kullback-Leibler divergence of the von Mises-Fisher distribution from the uniform distribution, which is increasing at the logarithmic rate in the concentration and generates the difference between the objective function and empirical welfare risk in Figure D.2. As the value of the concentration parameter increases, the regularisation term begins to dominate.

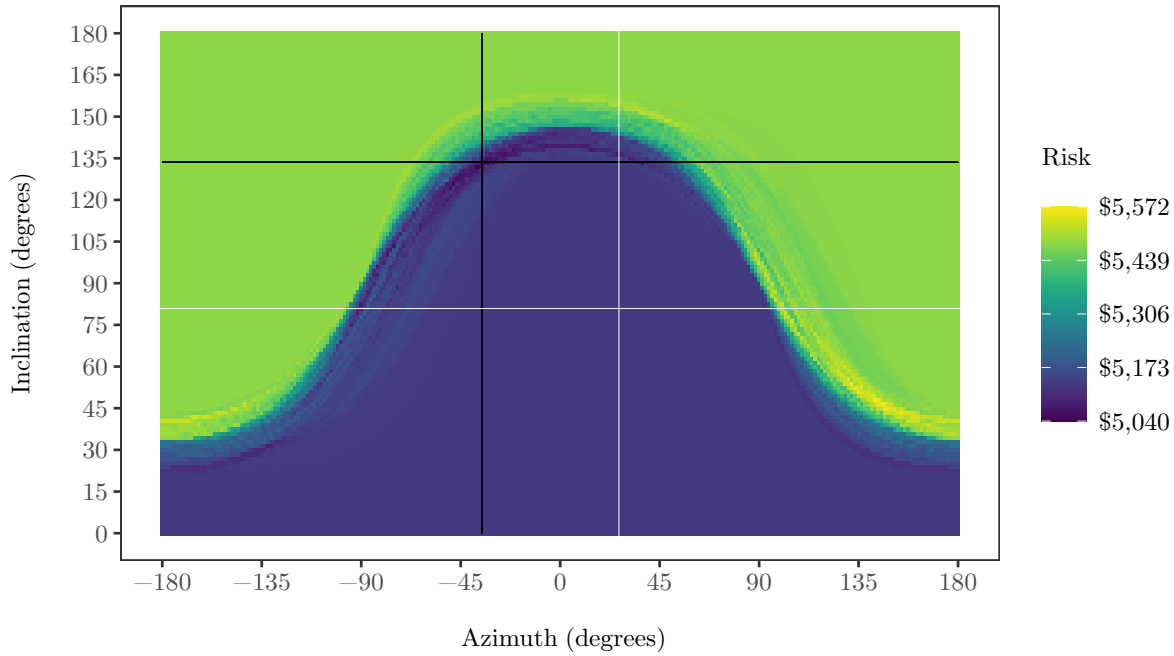
To better understand why μ^* does not coincide with the deterministic assignment rule of KITAGAWA and TETENOV (2018b) holding fixed $\kappa = \kappa^*$, we refer to Figure D.3, which plots empirical welfare risk for all vectors on the unit sphere—i.e., the empirical welfare risk associated for each deterministic assignment rule in \mathfrak{L} . Figure D.3 utilises the spherical coordinate system

$$\{\beta_0, \beta_1, \beta_2\} = \{\cos(\theta) \cdot \sin(\phi), \sin(\theta) \cdot \sin(\phi), \cos(\phi)\} \tag{D.35}$$

where $-180^\circ \leq \theta < 180^\circ$ is the azimuth and $0^\circ \leq \phi \leq 180^\circ$ is the inclination. It is perhaps convenient to think of the azimuth as related to longitude and the inclination as related to latitude. For non-trivial values of the concentration parameter, the von Mises-Fisher distribution allocates probability mass to the sphere in such a way that its density contours are concentric about the mean direction, with points closer to this direction more likely to occur. The deterministic assignment rule of KITAGAWA and TETENOV (2018b) can be seen from Figure D.3 to be located on the boundary between a high risk region (no-one treated) and a moderate risk region (everyone treated). As such, a stochastic assignment rule with a

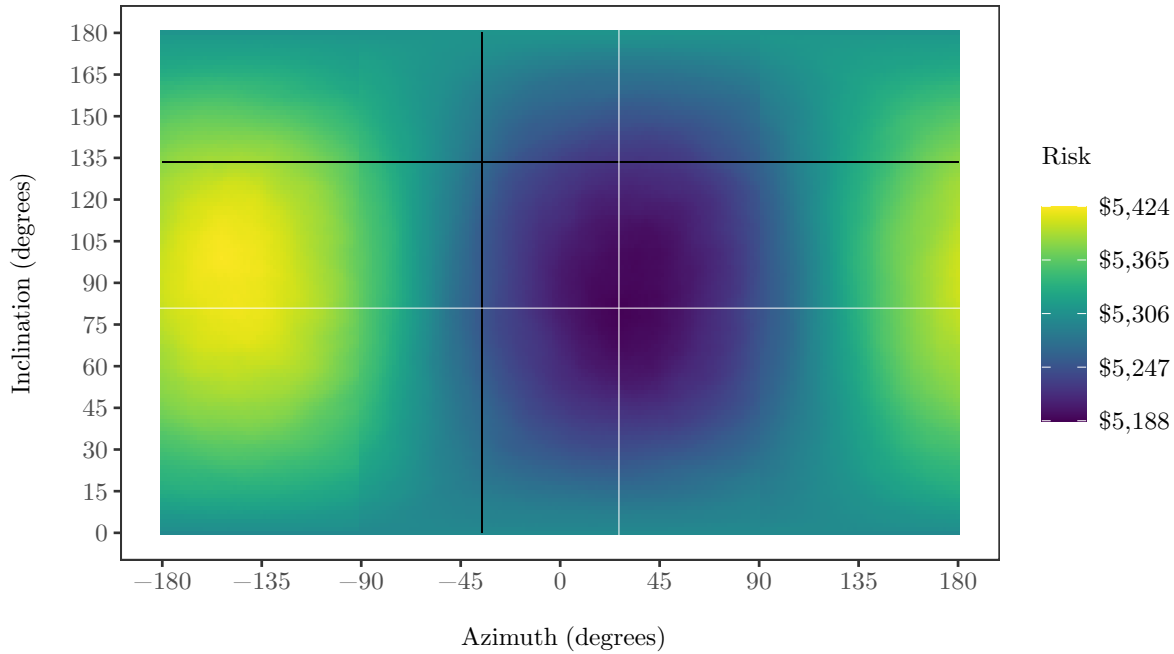
¹² That empirical welfare risk decreases for small to moderate values of the concentration parameter is specific to the data and chosen mean direction, and is arguably also attributable to the lack of consideration given to the cost of treatment.

FIGURE D.3
Deterministic assignment rules and empirical welfare risk



JTPA Study sample. This figure illustrates the risk that is associated with (deterministic) assignment rules in \mathcal{L} . A spherical coordinate mapping is implemented. The intersection of the two white lines is located at μ^ . The intersection of the two black lines is located at the optimal deterministic assignment rule of KITAGAWA and TETENOV (2018b), which attains the minimal regret amongst all deterministic linear rules.*

FIGURE D.4
Behaviour of the objective function at κ^ given variation in the mean direction μ*



JTPA Study sample. This figure illustrates the risk that is associated with (stochastic) assignment rules in \mathfrak{Q} ; the concentration parameter is fixed at κ^ whilst μ is varied. The intersection of the two white lines is located at μ^* . The intersection of the two black lines is located at the optimal deterministic assignment rule of KITAGAWA and TETENOV (2018b), which attains the minimal regret amongst all deterministic linear rules.*

non-trivial value of the concentration parameter with its mean direction located at this point would approximately allocate probability mass to each of these regions in equal amounts. By shifting the mean direction towards the centre of the moderate risk region, we allocate relatively more mass to rules that induce moderate risk and less mass to rules that induce high risk, which reduces empirical welfare risk overall.

This pattern underlies what we observe in Figure D.4, which plots empirical welfare risk for all directions on the sphere holding fixed $\kappa = \kappa^*$. Despite the apparent discontinuity of risk over deterministic assignment rules, empirical welfare risk (and the objective function) appear to vary smoothly.

—| APPENDIX D.1 |—

PROOFS

\hookrightarrow *Proof of Theorem D.2.* An optimal posterior minimises

$$E_{\Pi}(R_S(g)) + \sqrt{\frac{1}{2n} \cdot \left[E_{\Pi} \left(\ln \left(\frac{d\Pi(g)}{d\Pi_0(g)} \right) \right) + \ln \left(\frac{2\sqrt{n}}{\varepsilon} \right) \right]} \text{ subject to } \int_{\mathfrak{G}} d\Pi(g) = 1, \quad \text{D.36}$$

$$\inf_{g \in \mathfrak{G}} d\Pi(g) \geq 0$$

Provided that the solution satisfies the non-negativity constraints (the second constraint of Equation D.36), this is equivalent to minimising

$$\int_{\mathfrak{G}} R_S(g) \cdot d\Pi(g) + \sqrt{\frac{1}{2n} \cdot \left[\int_{\mathfrak{G}} \ln \left(\frac{d\Pi(g)}{d\Pi_0(g)} \right) \cdot d\Pi(g) + \ln \left(\frac{2\sqrt{n}}{\varepsilon} \right) \right]} + \xi \cdot [\Pi(\mathfrak{G}) - 1] \quad \text{D.37}$$

where ξ is the Lagrange multiplier. We separate this minimisation into two parts, by minimising Equation D.37 over Π subject to its Kullback-Leibler divergence from Π_0 being equal to $c \geq 0$ and, subsequently, by searching for the value of c that minimises the objective function.

We can write the constrained minimisation that forms the first part of the problem as

$$\int_{\mathfrak{G}} R_S(g) \cdot d\Pi(g) + \sqrt{\frac{1}{2n} \cdot \left[c + \ln \left(\frac{2\sqrt{n}}{\varepsilon} \right) \right]} + \xi \cdot [\Pi(\mathfrak{G}) - 1] + \frac{1}{\chi} \cdot \left[\int_{\mathfrak{G}} \ln \left(\frac{d\Pi(g)}{d\Pi_0(g)} \right) \cdot d\Pi(g) - c \right] \quad \text{D.38}$$

where ξ and $1/\chi$ are Lagrange multipliers, and where $c \geq 0$ is a constant, provided that the omitted non-negativity constraints are satisfied at the solution. The associated first order condition of the minimand with respect to $d\Pi(g)$ is

$$R_S(g) + \xi + \frac{1}{\chi} \cdot \left[\ln \left(\frac{d\Pi(g)}{d\Pi_0(g)} \right) + 1 \right] = 0 \quad \text{D.39}$$

Rearranging, we obtain

$$d\Pi^*(g) = \frac{\exp(-\chi \cdot R_S(g))}{\exp(1 + \xi \cdot \chi)} \cdot d\Pi_0(g) \quad \text{D.40}$$

where we emphasise that Π^* is a function of χ . In view of the first constraint of Equation D.37,

$$\exp(1 + \xi \cdot \chi) = \int_{\mathfrak{G}} \exp(-\chi \cdot R_S(g)) \cdot d\Pi_0(g) \quad \text{D.41}$$

and so, for all $g \in \mathfrak{G}$,

$$d\Pi^*(g) = \frac{\exp(-\chi \cdot R_S(g))}{\int_{\mathfrak{G}} \exp(-\chi \cdot R_S(g)) \cdot d\Pi_0(g)} \cdot d\Pi_0(g) \quad \text{D.42}$$

which integrates to one as is required. We reiterate that Equation D.42 is derived for an arbitrary (but non-negative) c , and so holds for any feasible values of χ .

For Π^* to satisfy $d_\ell(\Pi^*, \Pi_0) = c$, we require that

$$\begin{aligned} c &= \int_{\mathfrak{G}} \ln \left(\frac{d\Pi^*(g)}{d\Pi_0(g)} \right) \cdot d\Pi^*(g) \\ &= - \int_{\mathfrak{G}} \chi \cdot R_S(g) \cdot d\Pi^*(g) - \ln \left(\int_{\mathfrak{G}} \exp(-\chi \cdot R_S(g)) \cdot d\Pi_0(g) \right) \end{aligned} \quad \text{D.43}$$

This relationship between the radius of the Kullback-Leibler ball and the inverse of the Lagrange multiplier shows that c is strictly monotonically increasing in χ provided that $R_S(g)$ is not constant over those g supported by Π_0 , since

$$\frac{dc}{d\chi} = - \int_{\mathfrak{G}} R_S(g) \cdot d\Pi^*(g) - \int_{\mathfrak{G}} \chi \cdot R_S(g) \cdot \left[\frac{d}{d\chi} \frac{d\Pi^*(g)}{d\Pi_0(g)} \right] \cdot d\Pi_0(g) + \int_{\mathfrak{G}} R_S(g) \cdot d\Pi^*(g) \quad \text{D.44}$$

$$= \chi \cdot \int_{\mathfrak{G}} [R_S(g) - R_S^{\Pi^*}]^2 \cdot d\Pi^*(g) \quad \text{D.45}$$

$$\geq 0 \quad \text{D.46}$$

where Equation D.46 is strict if $\chi > 0$ and $R_S(g) \neq R_S^{\Pi^*}$ for at least some g supported by Π_0 . That χ is non-negative is evident upon further consideration of Equation D.42: posterior distributions associated with negative values of χ assign more mass to policies that are associated with high levels of risk and so cannot be optimal. Accordingly, in the second part of the optimisation, we substitute Π^* into Equation D.38 in place of Π and solve the unconstrained minimisation problem with respect to χ . That

is, we solve

$$\min_{\chi} \left(R_S^{\Pi^*} + \sqrt{\frac{1}{2n} \cdot \left[-\chi \cdot R_S^{\Pi^*} - \ln \left(\int_{\mathfrak{G}} \exp(-\chi \cdot R_S(g)) \cdot d\Pi_0(g) \right) + \ln \left(\frac{2\sqrt{n}}{\varepsilon} \right) \right]} \right) \quad \text{D.47}$$

Penalty

We note that

$$\begin{aligned} \frac{d}{d\chi} R_S^{\Pi^*} &= \frac{d}{d\chi} \int_{\mathfrak{G}} R_S(g) \cdot \frac{\exp(-\chi \cdot R_S(g))}{\int_{\mathfrak{G}} \exp(-\chi \cdot R_S(g)) \cdot d\Pi_0(g)} \cdot d\Pi_0(g) \\ &= \int_{\mathfrak{G}} R_S(g) \cdot \frac{d}{d\chi} \frac{\exp(-\chi \cdot R_S(g))}{\int_{\mathfrak{G}} \exp(-\chi \cdot R_S(g)) \cdot d\Pi_0(g)} \cdot d\Pi_0(g) \\ &= - \int_{\mathfrak{G}} [R_S(g) - R_S^{\Pi^*}]^2 \cdot d\Pi^*(g) \end{aligned} \quad \text{D.48}$$

which we interpret as the variance of empirical welfare risk under Π^* . Using this result, we further note that

$$\begin{aligned} \frac{d}{d\chi} \text{Penalty} &= \frac{1}{4n \cdot \text{Penalty}} \cdot \left(-R_S^{\Pi^*} - \chi \cdot \frac{d}{d\chi} R_S^{\Pi^*} + \frac{\int_{\mathfrak{G}} R_S(g) \cdot \exp(-\chi \cdot R_S(g)) \cdot d\Pi_0(g)}{\int_{\mathfrak{G}} \exp(-\chi \cdot R_S(g)) \cdot d\Pi_0(g)} \right) \\ &= \frac{1}{4n \cdot \text{Penalty}} \cdot -\chi \cdot \frac{d}{d\chi} R_S^{\Pi^*} \\ &= \frac{1}{4n \cdot \text{Penalty}} \cdot \chi \cdot \int_{\mathfrak{G}} [R_S(g) - R_S^{\Pi^*}]^2 \cdot d\Pi^*(g) \end{aligned} \quad \text{D.49}$$

Together Equations D.48 and D.49 imply that the associated first order condition of Equation D.47 with respect to χ is

$$\frac{1}{4n \cdot \text{Penalty}} \cdot \chi \cdot \int_{\mathfrak{G}} [R_S(g) - R_S^{\Pi^*}]^2 \cdot d\Pi^*(g) = \int_{\mathfrak{G}} [R_S(g) - R_S^{\Pi^*}]^2 \cdot d\Pi^*(g) \quad \text{D.50}$$

Provided—trivially, we might add—that Π^* is not degenerate and there is variation in empirical welfare risk (the conditions under which Equation D.46 is strict), then this condition reduces to

$$\chi = 4n \cdot \sqrt{\frac{1}{2n} \cdot \left[d_{\ell}(\Pi^*, \Pi_0) + \ln \left(\frac{2\sqrt{n}}{\varepsilon} \right) \right]} \quad \text{D.51}$$

which is exactly the condition that appears in Theorem D.2.

Although Equations D.42 and D.51 characterise a possible interior solution of the optimisation, we have yet to guarantee that this proposed solution is a global optimum. To address this, we show that the two possible corner solutions are sub-optimal, such that the first order conditions from which Equation D.51 are derived are applicable. Continuity and differentiability of Equation D.38 are then sufficient¹³ to

¹³ It is possible to show that the difference between the left- and right-hand sides of Equation D.51 is negative at $\chi = 0$

guarantee that a fixed point satisfying Equation D.51 exists.

Neither Π_0 nor a degenerate distribution are optimal. To establish that Π_0 is not optimal, we note that Π^* coincides with Π_0 when $\chi = 0$. By marginally increasing χ , such that we move in the direction of Π^* , we obtain a probability distribution that is in the interior. It suffices to show that such a probability distribution reduces the value of the objective function relative to Π_0 . Evaluating Equations D.48 and D.49 at $\chi = 0$, we obtain

$$\begin{aligned} \frac{d}{d\chi} [R_S^{\Pi^*} + \text{Penalty}]|_{\chi=0} &= \frac{d}{d\chi} R_S^{\Pi_0} \\ &< 0 \end{aligned} \tag{D.52}$$

as we require, with any χ satisfying

$$0 < \chi < 2n \cdot R^{\Pi_0} \tag{D.53}$$

yielding a strictly lower value than $\chi = 0$. To establish that a degenerate distribution is not optimal, we note that such a distribution implies infinite divergence from Π_0 (if Π_0 were itself degenerate then our analysis would be meaningless since the prior and posterior distributions would always coincide; and if Π_0 is atomic then the discrete case would apply). Given that Π_0 attains a finite value of the objective function, however, Π_0 is always preferred to a degenerate distribution, and so a degenerate distribution cannot be optimal ■

To prove Theorem D.3, we rely on several intermediate results.

Lemma D.2. *The circular variance of a m -variate von Mises-Fisher random vector with concentration $\kappa \geq 0$ is bounded from above by*

$$\frac{m-1 + \sqrt{4\kappa^2 + [m+1]^2} - 2\kappa}{m-1 + \sqrt{4\kappa^2 + [m+1]^2}} \leq \frac{2m}{m-1 + \sqrt{4\kappa^2 + [m+1]^2}} \tag{D.54}$$

and is an $O(1/\kappa)$ function

Lemma D.2 is trivially established by subtracting the lower bound of Equation B.11 from one, and so is stated without proof. That the circular variance is an $O(1/\kappa)$ function is otherwise established in Chapter C using a Hankel expansion.

and positive if χ is sufficiently large, with existence then established using extensions of the intermediate value theorem and its corollary, Bolzano's theorem.

Lemma D.3. *The Kullback-Leibler divergence of a m -variate von Mises-Fisher random vector with concentration $\kappa \geq 0$ is bounded from above by*

$$\frac{m-1}{2} \cdot \left[\ln \left(\frac{m-1}{2} + \sqrt{\kappa^2 + \left[\frac{m+1}{2} \right]^2} \right) - \ln(m) \right] + \sqrt{\kappa^2 + \left[\frac{m-1}{2} \right]^2} - \sqrt{\kappa^2 + \left[\frac{m+1}{2} \right]^2} + 1 \quad \text{D.55}$$

and is an $O(\ln(\kappa))$ function

\hookrightarrow *Proof of Lemma D.3.* Chapter C shows that, when Π is a von Mises-Fisher distribution and Π_0 is the uniform distribution over the hypersphere,

$$d_\ell(\Pi, \Pi_0) = \left[\frac{m}{2} - 1 \right] \cdot \ln \left(\frac{\kappa}{2} \right) - \underbrace{\ln \left(\text{I}_{m/2-1}(\kappa) \right)}_{\text{Bessel fn.}} - \ln \left(\Gamma \left(\frac{m}{2} \right) \right) + \underbrace{\frac{\text{I}_{m/2}(\kappa)}{\text{I}_{m/2-1}(\kappa)} \cdot \kappa}_{\text{Ratio fn.}} \quad \text{D.56}$$

To derive an upper bound on the Kullback-Leibler divergence that does not involve modified Bessel functions or their ratios, we replace the terms labelled Bessel fn. and Ratio fn. in Equation D.56 with appropriate lower and upper bounds, respectively. To do so, we rely on results (and notation) in Chapter B.

First, the term labelled Bessel fn. is bounded from below (recall that the term enters negatively) by

$$\frac{1}{2} \cdot \ln \left(\frac{2}{\kappa} \right) - \ln \left(\Gamma \left(\frac{m}{2} \right) \right) + \frac{m-1}{2} \cdot \ln \left(\frac{\kappa \cdot m}{m-1 + \sqrt{4\kappa^2 + [m+1]^2}} \right) + \frac{2\kappa^2}{m-1 + \sqrt{4\kappa^2 + [m+1]^2}} \quad \text{D.57}$$

Second, the term labelled Ratio fn. is bounded from above by

$$\frac{2\kappa^2}{m-1 + \sqrt{4\kappa^2 + [m-1]^2}} \quad \text{D.58}$$

Substituting Equations D.57 and D.58 into Equation D.56, cancelling terms and noting that

$$\left[\frac{m}{2} - 1 \right] \cdot \ln \left(\frac{\kappa}{2} \right) - \frac{1}{2} \cdot \ln \left(\frac{2}{\kappa} \right) = \frac{m-1}{2} \cdot \ln \left(\frac{\kappa}{2} \right) \quad \text{D.59}$$

we obtain

$$\frac{m-1}{2} \cdot \left[\ln \left(\frac{\kappa}{2} \right) - \ln \left(\frac{\kappa \cdot m}{\underline{a} + \sqrt{4\kappa^2 + \underline{a}^2}} \right) \right] - \frac{2\kappa^2}{\bar{a} + \sqrt{4\kappa^2 + \bar{a}^2}} + \frac{2\kappa^2}{\underline{a} + \sqrt{4\kappa^2 + \underline{a}^2}} \quad \text{D.60}$$

where $\underline{a} \doteq m-1$ and $\bar{a} \doteq m+1$. To obtain the required result, we simplify the first two terms of Equation D.60 and use the fact that, for any real a (in this case, equal to either \underline{a} or \bar{a} , which are

positive reals),

$$\frac{4\kappa^2}{a + \sqrt{4\kappa^2 + a^2}} = \sqrt{4\kappa^2 + a^2} - a \quad \text{D.61}$$

which makes use of the formula for the difference of two squares ■

\hookrightarrow *Proof of Theorem D.3.* With Lemmata D.2 and D.3 to hand, we begin by noting that, together, Assumption D.5 and Lemma D.1 imply that

$$R^{\bar{\Pi}} \leq \bar{R} + \inf_{\kappa} \left(\inf_{\mu} \left(\text{First term} + \text{Second term} + \sqrt{\text{Third term}} \right) \right) \quad \text{D.62}$$

where

$$\begin{aligned} \text{First term} &= \int_{\mathbb{S}^{m-1}} 2\zeta \cdot \|\beta - \bar{\beta}\|_2 \cdot d\Pi(\beta; \kappa, \mu) \\ \text{Second term} &= \frac{1}{\lambda} \cdot \left[a(\lambda, n) + d_{\ell}(\Pi, \Pi_0) + \ln\left(\frac{2}{\varepsilon}\right) \right] \\ \text{Third term} &= \frac{1}{2n} \cdot \left[d_{\ell}(\Pi, \Pi_0) + \ln\left(\frac{4\sqrt{n}}{\varepsilon}\right) \right] \end{aligned} \quad \text{D.63}$$

and which we note are functions (either explicit or implicit) in κ and μ .

Focusing on the First term, we use Jensen's inequality to show that

$$2\zeta \cdot \int_{\mathbb{S}^{m-1}} \|\beta - \bar{\beta}\|_2 \cdot d\Pi(\beta; \kappa, \mu) \leq 2\zeta \cdot \sqrt{\int_{\mathbb{S}^{m-1}} \sum_{i=0}^{m-1} (\beta_i - \bar{\beta}_i)^2 \cdot d\Pi(\beta; \kappa, \mu)} \quad \text{D.64}$$

$$= 2\zeta \cdot \sqrt{\text{tr}\left(\mathbf{E}_{\Pi}\left([\beta - \bar{\beta}][\beta - \bar{\beta}]^{\top}\right)\right)} \quad \text{D.65}$$

Rearranging, we obtain

$$\mathbf{E}_{\Pi}\left([\beta - \bar{\beta}][\beta - \bar{\beta}]^{\top}\right) = \mathbf{E}_{\Pi}\left([\beta - \bar{\beta} + \mathbf{E}_{\Pi}(\beta) - \mathbf{E}_{\Pi}(\beta)][\beta - \bar{\beta} + \mathbf{E}_{\Pi}(\beta) - \mathbf{E}_{\Pi}(\beta)]^{\top}\right) \quad \text{D.66}$$

$$= \mathbf{E}_{\Pi}\left([\beta - \mathbf{E}_{\Pi}(\beta)][\beta - \mathbf{E}_{\Pi}(\beta)]^{\top}\right) + [\mathbf{E}_{\Pi}(\beta) - \bar{\beta}][\mathbf{E}_{\Pi}(\beta) - \bar{\beta}]^{\top} \quad \text{D.67}$$

$$= \text{Variance}(\beta | \beta \sim \Pi) + \left[\frac{\mathbf{I}_{m/2}(\kappa)}{\mathbf{I}_{m/2-1}(\kappa)} \cdot \mu - \bar{\beta} \right] \left[\frac{\mathbf{I}_{m/2}(\kappa)}{\mathbf{I}_{m/2-1}(\kappa)} \cdot \mu - \bar{\beta} \right]^{\top} \quad \text{D.68}$$

which relies on results from Chapter C relating to the first two moments of the von Mises-Fisher family of distributions. Restricting the set of von Mises-Fisher distributions to those satisfying $\mu = \bar{\beta}$, the

First term is bounded by

$$\begin{aligned}
2\zeta \cdot \sqrt{\text{tr} \left(\mathbb{E} \left([\beta - \bar{\beta}] [\beta - \bar{\beta}]^\top \right) \right)} &= 2\zeta \cdot \sqrt{\text{tr} \left(\text{Variance}(\beta | \beta \sim \Pi) + \left[\frac{I_{m/2}(\kappa)}{I_{m/2-1}(\kappa)} - 1 \right]^2 \cdot \mu\mu^\top \right)} \\
&= \sqrt{8\zeta^2} \cdot \sqrt{1 - \frac{I_{m/2}(\kappa)}{I_{m/2-1}(\kappa)}}
\end{aligned} \tag{D.69}$$

which is proportional to the square root of the circular variance.¹⁴ Lemma D.2 details the behaviour of the circular variance. Using Lemma D.2, Equation D.69 can be bounded from above by

$$\sqrt{8\zeta^2 \cdot \frac{m-1 + \sqrt{4\kappa^2 + [m+1]^2} - 2\kappa}{m-1 + \sqrt{4\kappa^2 + [m+1]^2}}} \leq 4\zeta \cdot \sqrt{\frac{m}{2\kappa}} \tag{D.70}$$

where the right-hand side follows from concavity of the square root, and by decreasing the denominator via elimination of $m-1$ and $m+1$ (both are non-negative constants). We substitute the right-hand side of Equation D.70 into Equation D.62 in place of the First term.

We now focus on the Kullback-Leibler divergence, which appears in both the Second term and the Third term. We observe that the difference between the square roots in Equation D.55 is increasing in κ (i.e, the difference becomes less negative as κ increases and as κ increases in importance relative to m). As such, it suffices to omit the difference between the square roots in Equation D.55 and to simply bound the Kullback-Leibler divergence from above by

$$\frac{m-1}{2} \cdot \ln \left(\frac{m-1 + \sqrt{4\kappa^2 + [m+1]^2}}{2m} \right) + 1 \leq \frac{m-1}{2} \cdot \ln(\kappa+1) + 1 \tag{D.71}$$

which is at least one due to the non-negativity of κ .

Substituting Equations D.70 and D.71 as upper bounds on the infimand of Equation D.62, we obtain

$$4\zeta \cdot \sqrt{\frac{m}{2\kappa}} + \frac{1}{\lambda} \cdot \left[a(\lambda, n) + \frac{m-1}{2} \cdot \ln(\kappa+1) + \ln \left(\frac{2e}{\varepsilon} \right) \right] + \sqrt{\frac{1}{2n} \cdot \left[\frac{m-1}{2} \cdot \ln(\kappa+1) + \ln \left(\frac{4e \cdot \sqrt{n}}{\varepsilon} \right) \right]} \tag{D.72}$$

which we emphasise is an upper bound on the infimum and is not dependent of β . Our objective is to minimise Equation D.72 by appropriately choosing κ and λ alongside the functional form of a . Accordingly, we let $\kappa = n$ and $\lambda = \sqrt{n}$ alongside $a(\lambda, n)/\lambda = 1/\lambda$. Given these choices, we can write

¹⁴ Recall that the right-hand side of Equation D.62 is preceded by an infimum over μ ; replacing the infimum with a specific value delivers an upper bound.

Equation D.72 as

$$4\zeta \cdot \sqrt{\frac{m}{2n}} + \frac{m-1}{2\sqrt{n}} \cdot \ln(n+1) + \frac{1}{\sqrt{n}} \cdot \ln\left(\frac{2e^2}{\varepsilon}\right) + \sqrt{\frac{1}{2n} \cdot \left[\frac{m-1}{2} \cdot \ln(n+1) + \ln\left(\frac{4e \cdot \sqrt{n}}{\varepsilon}\right) \right]} \quad \text{D.73}$$

which, using the fact that $1/\ln(n)$ and $\ln(n+1)/\ln(n)$ are decreasing in n , we can upper bound by

$$\frac{\ln(n)}{\sqrt{n}} \cdot \frac{1}{\ln(8)} \cdot \underbrace{\left[4\zeta \cdot \sqrt{\frac{m}{2}} + \frac{m-1}{2} \cdot \ln(9) + \ln\left(\frac{2e^2}{\varepsilon}\right) + \sqrt{\frac{1}{2} \cdot \left[\frac{m-1}{2} \cdot \ln(9) + \ln\left(\frac{4e \cdot \sqrt{n}}{\varepsilon}\right) \right]} \right]}_{\text{Universal constant}} \quad \text{D.74}$$

where we rely on the maintained assumption of Theorem D.3 that $n \geq 8$. The Universal constant is a decreasing function of ε ■

—| APPENDIX D.2 |—

ACCOUNTING FOR THE COST OF TREATMENT IN THE JTPA STUDY SAMPLE

In what follows, we distinguish between the JTPA Study sample, which we sometimes refer to as the raw data, and the cost-adjusted JTPA Study sample, which we sometimes refer to as the costed data. The costed data subtracts the cost of treatment from the outcome of interest (post-programme earnings) of all treated individuals.¹⁵ Following Kitagawa & Tetenov (2018), we assume that the cost of treatment is \$774. We then proceed to search for the optimal stochastic assignment rule, applying the same grid search approach as we outlined in the main text for the raw data, using the costed data.

We find that the objective function is minimised by the stochastic assignment rule with $\kappa = 0.560$ and $\mu = \{+0.872, +0.490, +0.018\}$, which we label κ^a and μ^a , respectively.¹⁶ The value of the objective function and the empirical welfare risk induced by this assignment rule are equivalent to \$9,302 and \$5,155, respectively.

In comparison, KITAGAWA and TETENOV (2018b) estimates that the deterministic assignment rule defined by $\beta = \{+0.117, -0.990, -0.086\}$ minimises empirical welfare risk. Relative to the raw analysis, this assignment rule ascribes less weight to individual characteristics, and essentially determines that all individuals with low pre-programme earnings (earnings of around \$5,000 or less) should be treated, irrespective of education.

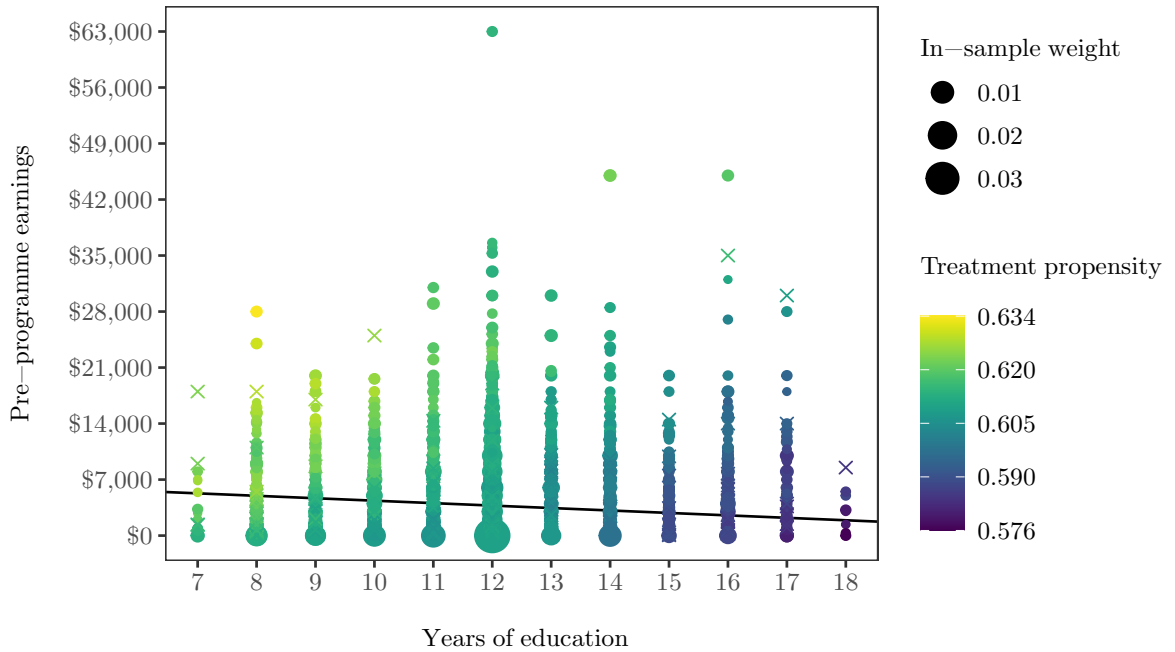
Figures D.5 to D.8 are intended to be comparable to the figures in the main text. There are two differences

¹⁵ We observe that this leads individuals with zero post-programme earnings to have a negative outcome. This violates Assumption 1, which requires that all outcomes be bounded and non-negative. We, nonetheless, proceed with this adjustment of post-programme earnings as is. An alternative would be to add \$774 to all post-programme earnings before subtracting the assumed cost of treatment.

¹⁶ This directional vector can be represented by an azimuth of 29° and an inclination of 89° using spherical coordinates.

FIGURE D.5

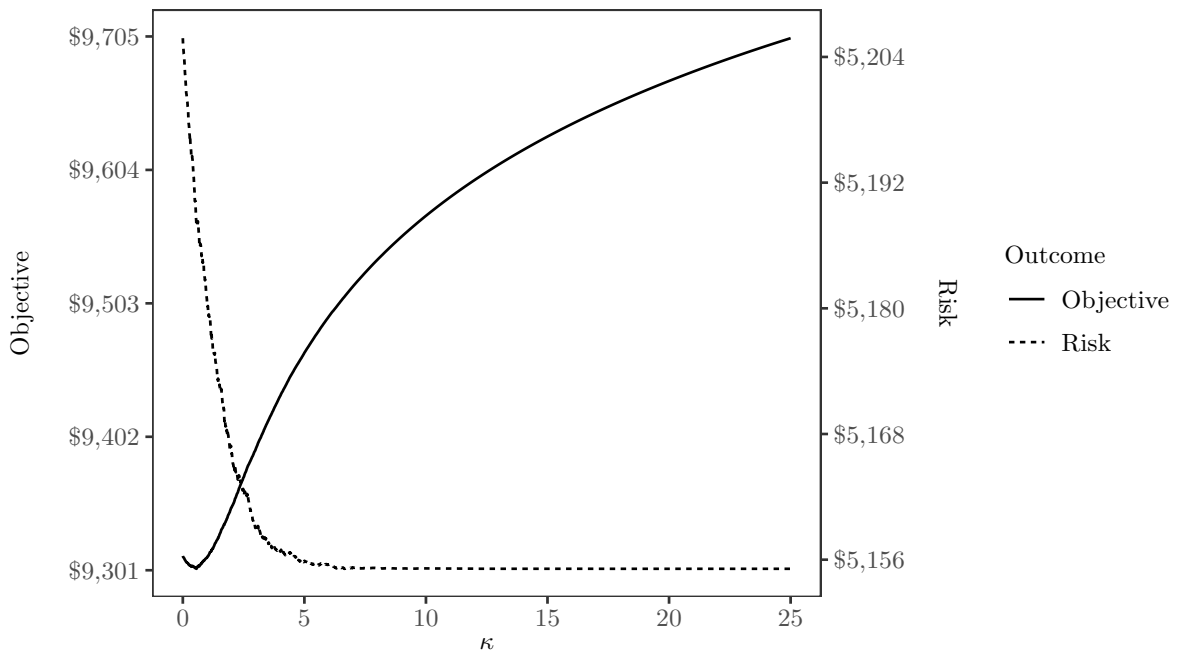
Variation in treatment propensity across individuals in the cost-adjusted JTPA Study sample



Cost-adjusted JTPA Study sample. This figure illustrates the treatment propensity of individuals under the posterior assignment rule that is induced by $\{\kappa^a, \mu^a\}$. Each point represents the individual characteristics of an individual or several individuals (crosses denote individuals with negative in-sample weight). For comparison, individuals to the left of the solid diagonal line are assigned treatment under the optimal deterministic assignment rule of KITAGAWA and TETENOV (2018b).

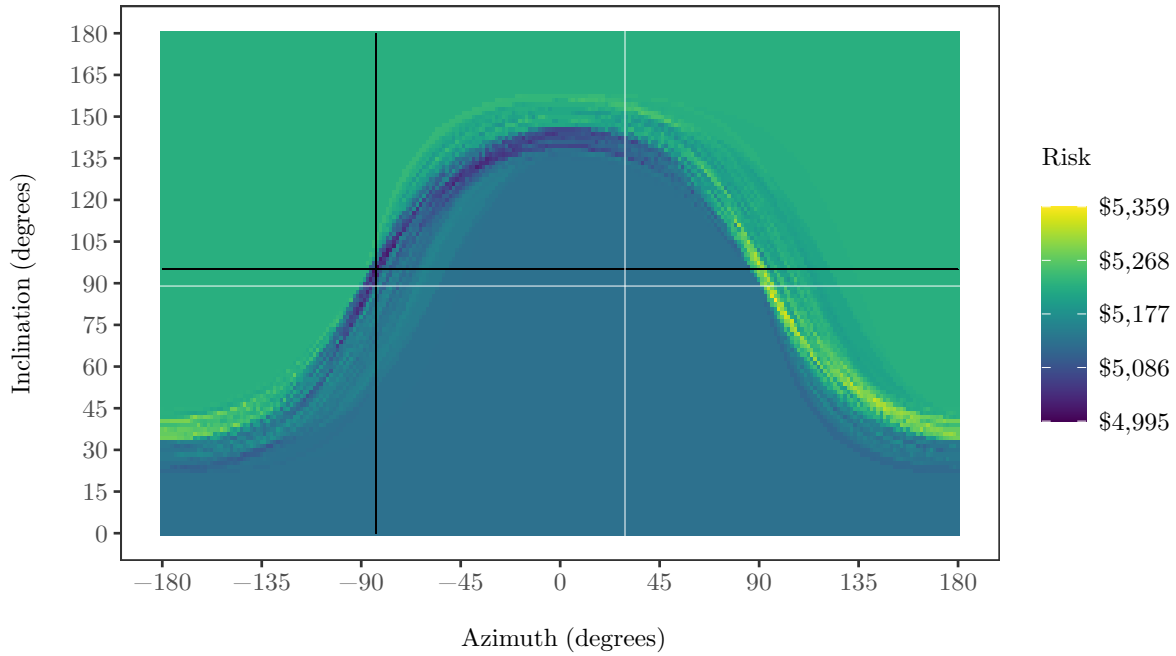
FIGURE D.6

Behaviour of the objective function at μ^* given variation in κ



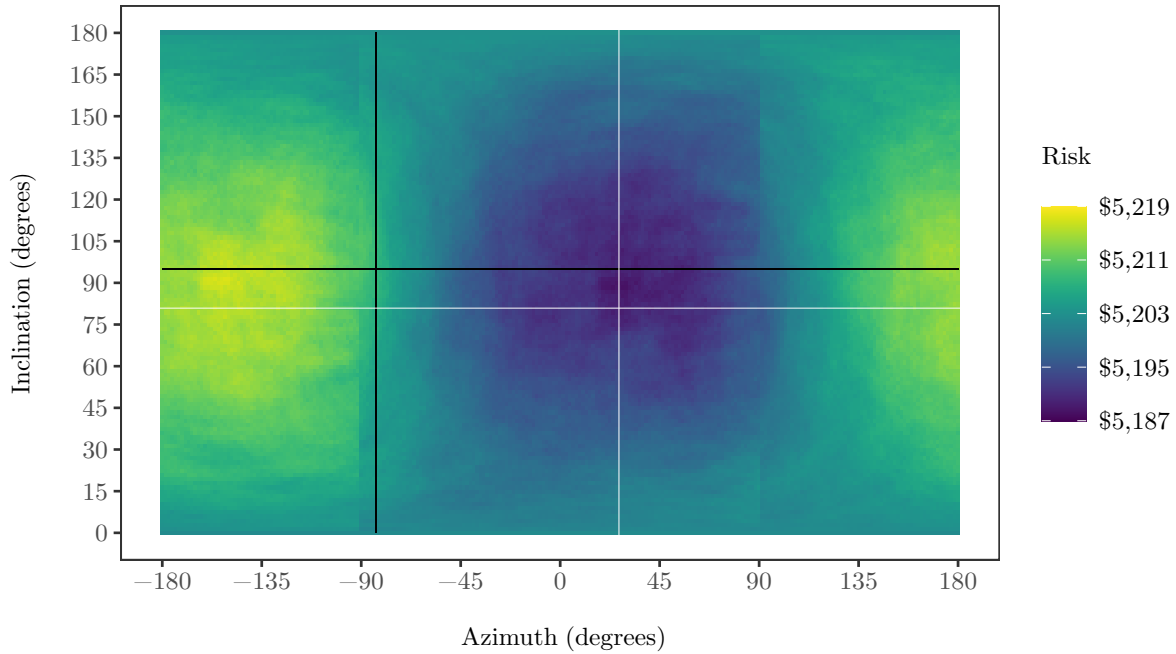
Cost-adjusted JTPA Study sample. This figure illustrates the shape of the objective function and its risk component at μ^a as κ is varied.

FIGURE D.7
Deterministic assignment rules and empirical welfare risk



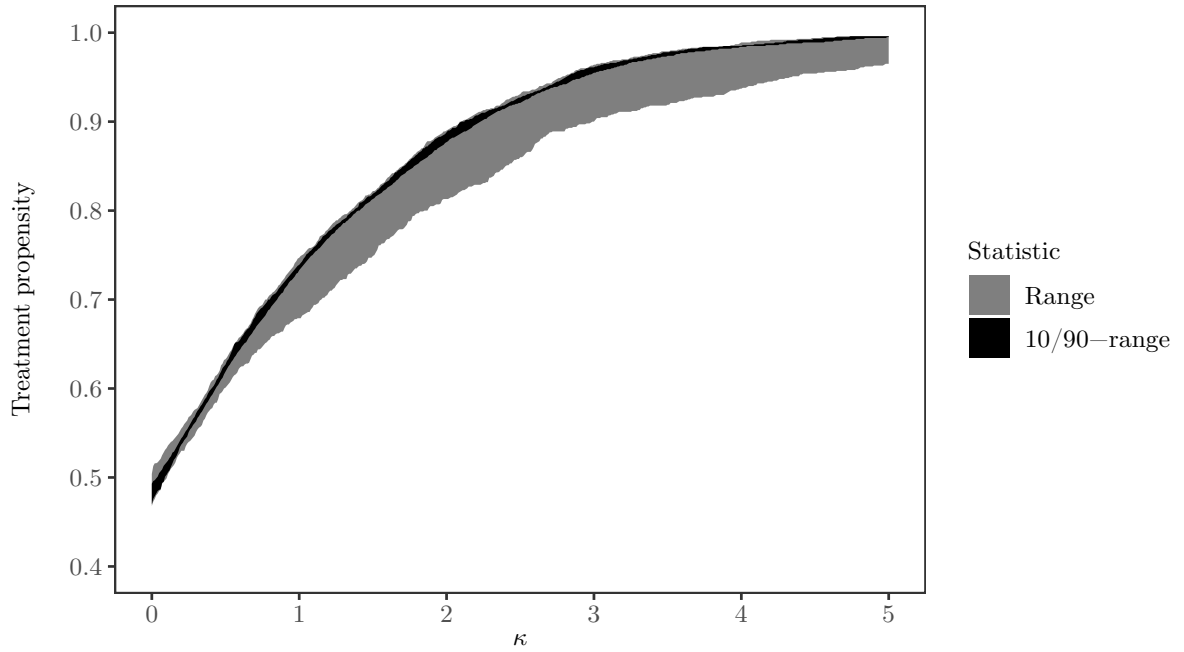
Cost-adjusted JTPA Study sample. This figure illustrates the risk that is associated with (deterministic) assignment rules in \mathfrak{F} . A spherical coordinate mapping is implemented. The intersection of the two white lines is located at μ^a . The intersection of the two black lines is located at the optimal deterministic assignment rule of KITAGAWA and TETENOV (2018b), which attains the minimal regret amongst all deterministic linear rules.

FIGURE D.8
Behaviour of the objective function at κ^a given variation in μ



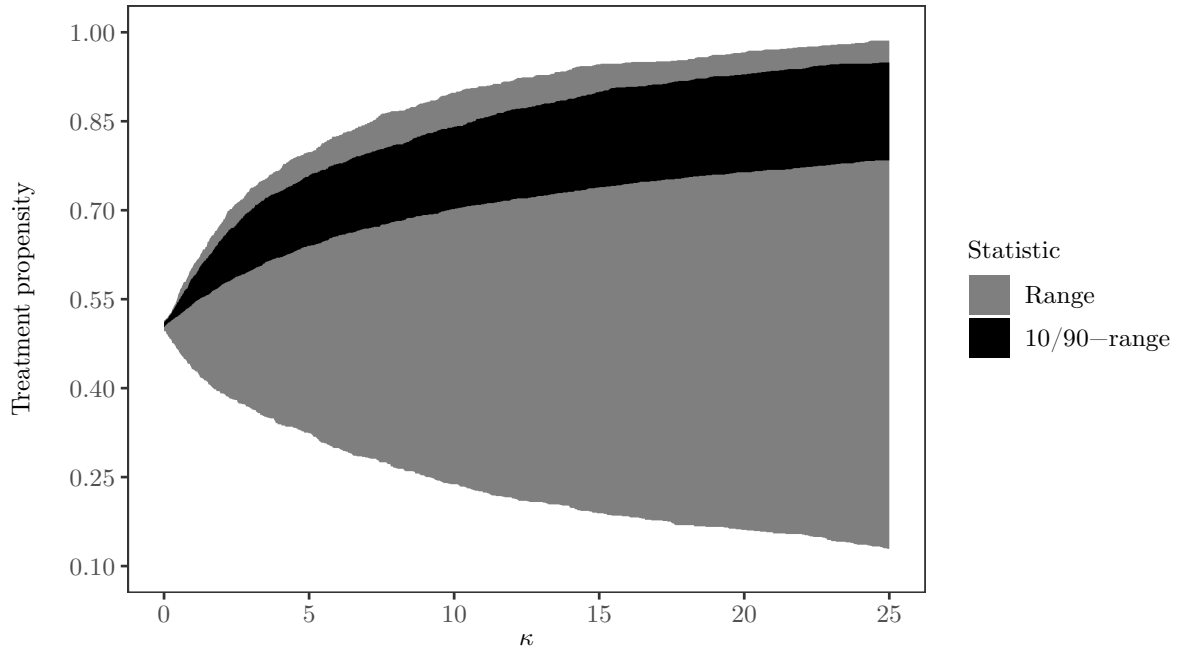
Cost-adjusted JTPA Study sample. This figure illustrates the risk that is associated with (stochastic) assignment rules in \mathfrak{V} ; the concentration parameter is fixed at κ^a whilst μ is varied. The intersection of the two white lines is located at μ^a . The intersection of the two black lines is located at the optimal deterministic assignment rule of KITAGAWA and TETENOV (2018b), which attains the minimal regret amongst all deterministic linear rules.

FIGURE D.9
Distribution of treatment propensity across individuals (adjusted) at $\{\kappa^a, \mu^a\}$



Cost-adjusted JTPA Study sample. This figure illustrates the variation in treatment propensity across individuals. We draw 1,000 directional vectors and count how many of these vectors, implemented as deterministic assignment rules, assign each individual to treatment. We sort individuals by how often they are assigned to treatment and summarise the implied distribution by the range and the 10th and 90th quantiles.

FIGURE D.10
Distribution of treatment propensity across individuals (raw) at $\{\kappa^, \mu^*\}$*



JTPA Study sample. This figure illustrates the variation in treatment propensity across individuals. We draw 1,000 directional vectors and count how many of these vectors, implemented as deterministic assignment rules, assign each individual to treatment. We sort individuals by how often they are assigned to treatment and summarise the implied distribution by the range and the 10th and 90th quantiles.

of note between the results of the raw and costed analyses. First, in Figure D.6, the relationship between the concentration parameter and empirical welfare risk is much more muted, with risk only decreasing by around \$50 as the value of the concentration parameter increases (over the interval that we study). In contrast, empirical welfare risk decreases by around \$200 when the cost of treatment is not accounted for. This reflects the second difference, which is apparent in Figure D.7. Specifically, that the difference in empirical welfare risk between those assignment rules that assign everyone to treatment versus those that assign no-one to treatment is much smaller. The benefit to an increase in the value of the concentration parameter, which is a reduction in the probability mass allocated to those assignment rules that assign no-one to treatment is, accordingly, smaller. The implication is that the penalty term starts to dominate the objective function for smaller values of the concentration parameter.

To supplement our analysis, we include two additional figures in Figures D.9 and D.10. These figures are intended to give some idea about how the concentration of the von Mises-Fisher distribution influences the propensity with which individuals are assigned treatment. For example, given a particular instance of the von Mises-Fisher distribution, one individual might have a high propensity of assignment whereas another might have a low propensity. In other words, the number of directional vectors that are drawn from the von Mises-Fisher distribution for which the first individual is assigned treatment is greater than for the second individual. As the value of the concentration parameter increases, so these directional vectors concentrate around the mean direction, and the propensity of assignment to treatment approaches zero or one for each individual. We see from Figure D.9 that, when $\mu = \mu^*$, the propensity of assignment in the costed data tends towards one for all individuals. In contrast, and to illustrate the possibility of complete dichotomy, for the raw data, we see from Figure D.10 that the propensity of assignment to treatment diverges as the value of the concentration parameter increases. We observe in both datasets that much of the variation in propensity is driven by a few individuals in the tail, who are likely those individuals whose individual characteristics make them outliers. In both cases, the median propensity of treatment (i.e., the probability with which the average individual in the sample is assigned treatment) tends towards one.

—| APPENDIX D.3 |—

NUMERICAL SIMULATIONS

We propose several experiments that investigate how various aspects of the sample data that is available to the social planner affect the posterior distribution and its shape when the specified prior distribution is uniform over the sphere. We conduct these experiments to better understand some of the empirical results that we obtain for the main paper.

Aspects of the sample data that we vary include the number of observations, the outcome of interest, and individual characteristics. We conduct these experiments with simulated data and through manipulation of the existing empirical application. Each experiment is comparable to the existing empirical application in any case, in that observed individual characteristics are taken to be education and pre-programme earnings alongside an intercept term, and the outcome of interest is taken to be post-programme earnings. We investigate how these aspects and our variation of them affect the objective function.

The specific questions that we ask, and that inform the design of our experiments, are as follows.

- How does the number of observations influence the shape of the objective function?
- How does the relative influence of education and pre-programme earnings on post-programme earnings influence empirical welfare risk?
- How does the distribution of individual characteristics influence the shape of the objective function?

To address these questions we propose a series of linear specifications that satisfy the bounded outcomes assumption that we require. We concede that these specifications are somewhat contrived, reflecting our need to balance tractability with the requirements of the bounded outcomes assumption. We suppose that

$$\begin{aligned} 0 \leq X_{\text{earn}} \leq 1 \\ 0 \leq X_{\text{educ}} \leq 1 \end{aligned} \tag{D.75}$$

which can always be maintained via an appropriate affine map of the individual characteristics (a step that we undertake in any case), and propose that

$$\begin{aligned} Y_1 &= \mathbf{X}^\top \alpha_1 + V_1 \text{ with } V_1 \sim N(0, \sigma_1^2) \text{ such that } 0 \leq V_1 \leq c_v \\ Y_0 &= \mathbf{X}^\top \alpha_0 + V_0 \text{ with } V_0 \sim N(0, \sigma_0^2) \text{ such that } 0 \leq V_0 \leq c_v \end{aligned} \tag{D.76}$$

with $\alpha_1 \in \mathbb{S}^2$ and $\alpha_0 \in \mathbb{S}^2$. The advantage of this specification is that it provides a clear interpretation of the influence of education and pre-programme earnings on post-programme earnings, and facilitates addressing all of the questions that we pose. Moreover, the potential outcomes satisfy

$$\begin{aligned} 0 \leq Y_1 \leq \alpha_1 + c_v \\ 0 \leq Y_0 \leq \alpha_1 + c_v \end{aligned} \tag{D.77}$$

We note that the conditional expectations of the potential outcomes are

$$\begin{aligned} E(Y_1|\mathbf{X}) &= \mathbf{X}^\top \alpha_1 + \frac{\Phi'(0) - \Phi'(c_v/\sigma_1)}{\Phi(c_v/\sigma_1) - \Phi(0)} \cdot \sigma_1 = \mathbf{X}^\top \tilde{\alpha}_1 \\ E(Y_0|x) &= \mathbf{X}^\top \alpha_0 + \frac{\Phi'(0) - \Phi'(c_v/\sigma_0)}{\Phi(c_v/\sigma_0) - \Phi(0)} \cdot \sigma_0 = \mathbf{X}^\top \tilde{\alpha}_0 \end{aligned} \tag{D.78}$$

which are both affine functions of individual characteristics.¹⁷ This is an attractive property that we exploit.

Knowledge of the data generating process (Equation D.76 in our framework) is sufficient to determine the optimal assignment rule.¹⁸ Specifically, the optimal assignment rule satisfies

$$g(\mathbf{X}) = 1(E(Y_1|\mathbf{X}) \geq E(Y_0|\mathbf{X})) \tag{D.79}$$

When the conditional expectation of the potential outcomes are affine functions of individual characteristics (and assignment to treatment is at random, which is an assumption that we implicitly maintain throughout) then the optimal assignment rule belongs to the LES class. Hence, Equation D.78 guarantees that the optimal policy has the specific form

$$g(\mathbf{X}) = 1(\mathbf{X}^\top [\tilde{\alpha}_1 - \tilde{\alpha}_0] / \|\tilde{\alpha}_1 - \tilde{\alpha}_0\|_2 \geq 0) \tag{D.80}$$

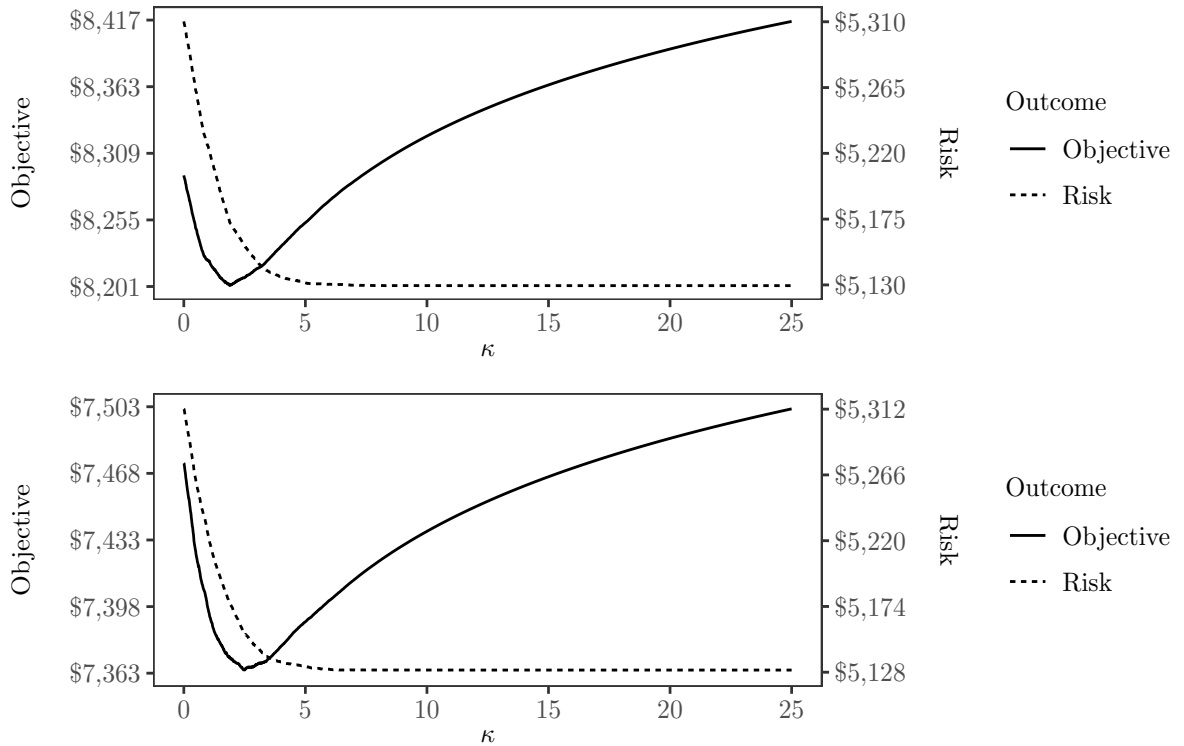
which is not only a member of the LES class but emphasises that the optimal assignment rule can be summarised by a vector on the sphere (provided that the potential outcomes have distinct process). This result provides motivation for our reliance on spherical distributions and, in particular, the von Mises-Fisher distribution. We observe that, abstracting from the issue of sampling variation, empirical welfare risk is minimised when the mean direction of the von Mises-Fisher distribution coincides with the policy defined in Equation D.80. This insight allows us to focus exclusively on the influence of the concentration parameter on empirical welfare risk and the objective function.

Throughout, we are careful to sample using inversion-based pseudo-random sampling methods where possible. Compared to rejection sampling-based pseudo-random sampling methods, inversion-based methods are able to guarantee comparability across experiments despite differences in parameter values. We also note that where we use data from the JTPA Study sample, this data is not adjusted for the

¹⁷ Other commonly invoked models fail to meet our test of tractability or the requirements of the bounded outcomes assumption: the standard censored outcome model generates a rectified normal distribution that is non-linear in individual characteristics and is unbounded from above; logarithmic transformation of the outcome implies a non-linear transformation of individual characteristics.

¹⁸ We do not make any claims about uniqueness in what follows and, indeed, presented with a finite sample of individual characteristics, it is likely that several assignment rules can attain the same partition of individuals as what we refer to as the optimal assignment rule.

FIGURE D.11
Behaviour of the objective function at μ^* given variation in κ



JTPA Study sample (artificially inflated by copying observations). This figure illustrates how the objective function and its risk component change as the number of observations is doubled (upper panel) and quadrupled (lower panel). We hold μ fixed at μ^* in each case, which differs according to the sample size, and vary κ .

cost of treatment (i.e., we use the raw data). We continue to label those values of the parameters that minimise the objective function by κ^* and μ^* , which we emphasise can vary across the various experiments that we conduct.

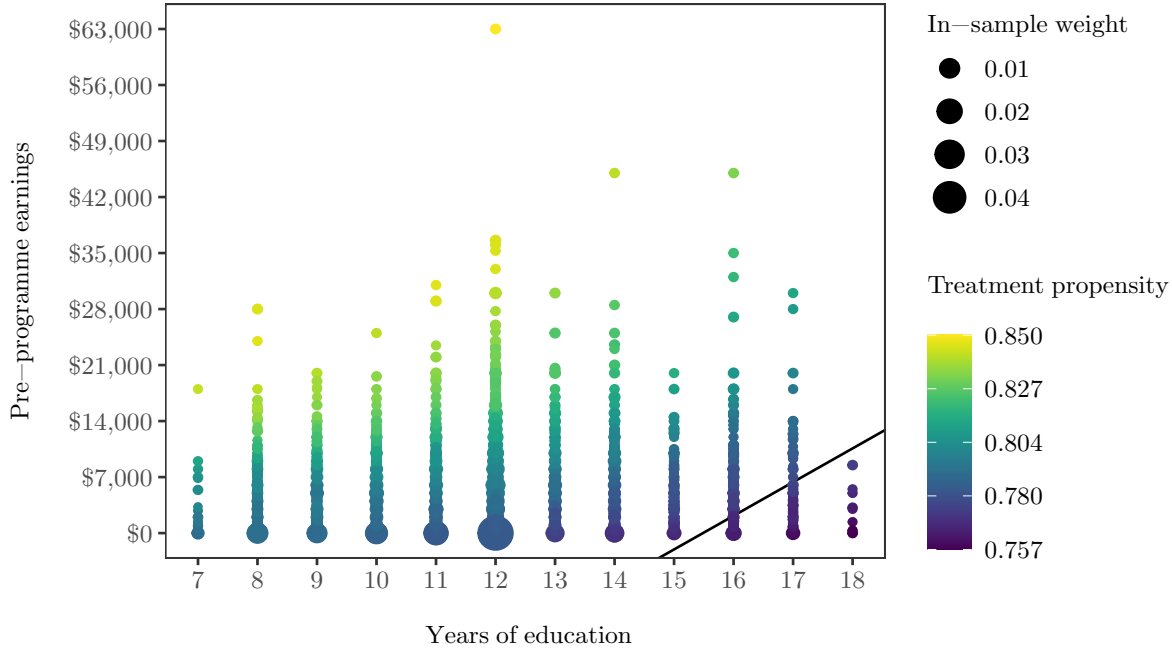
—| SUBAPPENDIX D.3.A |—

VARYING THE NUMBER OF OBSERVATIONS

To investigate the influence that the number of observations has on the objective function and on empirical welfare risk, we copy each observation in the JTPA Study sample either once or three times, thereby doubling and quadrupling the number of individuals in the sample. The mechanical effect of this change is to mute the influence of the penalty term in the objective function, without affecting empirical welfare risk. The immediate implication is that larger values of the concentration parameter can be sustained, since the penalty due to deviating from the uniform distribution is relatively smaller for every such deviation.

We plot the objective function and empirical welfare risk for the doubled and quadrupled JTPA Study samples in Figure D.11. We observe that increasing the number of observations has the direct effect of

FIGURE D.12
 (Experiment 1) Variation in treatment propensity across individuals



Simulated data generated using Equation D.82 in conjunction with data from the JTPA Study sample. This figure illustrates the treatment propensity of individuals in a simulated sample under the posterior assignment rule that is induced by $\{\kappa^*, \mu^*\}$. Each point represents the individual characteristics of an individual or several individuals. For comparison, individuals to the left of the solid diagonal line are assigned treatment under the oracle assignment rule of Equation D.83.

decreasing the magnitude of the objective function. Moreover, increasing the number of observations also leads to an increase in κ^* , which increases from 1.550 to 1.890 and then to 2.490. We emphasise that μ^* is also not the same across the two cases. The intuition here is that an increase in the value of the concentration parameter leads to the concentric contour map of the density function becoming more tightly arranged around the mean direction (whatever that may be). Locating the mean direction closer to the boundary between high and moderate regret regions say, such as where the deterministic assignment rule of KITAGAWA and TETENOV (2018b) is located, incurs less of a penalty in this instance since the density function assigns less probability mass to the high regret region than it would for a smaller value of the concentration parameter. For the doubled sample we find that $\mu^* = \{+0.812, +0.577, +0.088\}$, whilst for the quadrupled sample we find that $\mu^* = \{+0.917, +0.394, +0.053\}$.¹⁹

—| SUBAPPENDIX D.3.B |—

THE VARIANCE OF POST-PROGRAMME EARNINGS

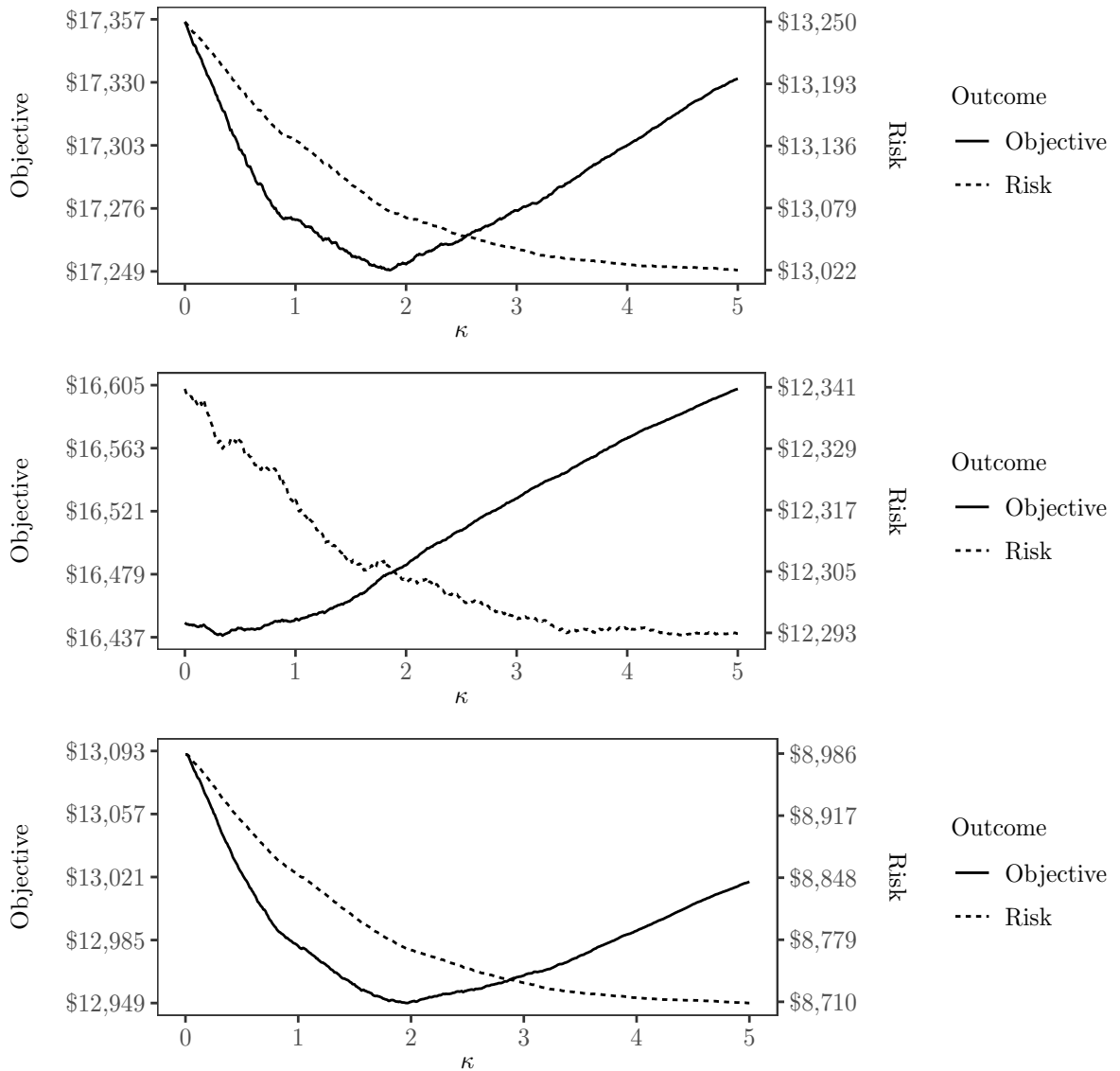
In the following three experiments, we investigate the influence of the variance of post-programme earnings on the results of our method. Specifically, we investigate how altering σ_1 and σ_0 in the specification

¹⁹ These directions translate (azimuth:inclination) to $35^\circ : 85^\circ$ and $23^\circ : 87^\circ$, respectively, as compared to $\mu^* = \{+0.883, +0.442, +0.158\}$ or $27^\circ : 81^\circ$ in the original sample.

FIGURE D.13

(Experiments 1–3) Behaviour of the objective function at μ^* given variation in κ

Experiment	μ^*	κ^*
1. Baseline (top panel)	{+0.822, +0.565, +0.078}	1.850
2. High variance (middle panel)	{+0.641, +0.703, -0.309}	0.340
3. Low variance (bottom panel)	{+0.831, +0.547, +0.105}	1.970



Simulated data generated using Equation D.82 in conjunction with data from the JTPA Study sample. This figure illustrates how the shape of the objective function and its risk component changes as the variance of post-programme earnings is increased and decreased. We hold μ fixed at μ^* in each case, which varies according to the sample size, and vary κ . Parameter estimates are presented above.

outlined in Equation D.76 affects our results. In particular, we are interested in whether there is a fundamental change in how empirical welfare risk varies with the parameters of the von Mises-Fisher distribution.

The JTPA Study sample consists of 9,223 observations, and we extract treatment status and individual characteristics directly from this dataset (we refer to the group of individuals in the sample who are treated as the experimental group and to the group of individuals who are not as the control group in line with the terminology surrounding randomised control trials). We then generate potential outcomes according to Equation D.76. We implement the affine map of Equation D.75 via the transformations

$$\begin{aligned} X_{\text{earn}} &\mapsto X_{\text{earn}} / \max(X_{\text{earn}}) \\ X_{\text{educ}} &\mapsto X_{\text{educ}} / \max(X_{\text{educ}}) \end{aligned} \tag{D.81}$$

We then regress post-programme earnings in the experimental group and the control group on these characteristics, separately and together, so that the baseline experiment somewhat mimics the JTPA Study sample. Using our simple regressions as a rough guide, we let

$$\begin{aligned} \alpha_1 &= \{+3,040; +86,446; +14,008\} \\ \alpha_0 &= \{-1,086; +82,458; +18,804\} \end{aligned} \quad \text{with } \sigma_1 = \sigma_0 = 15,914 \text{ and } c_v = 5\sigma_1 = 5\sigma_0 \tag{D.82}$$

so that the optimal policy is

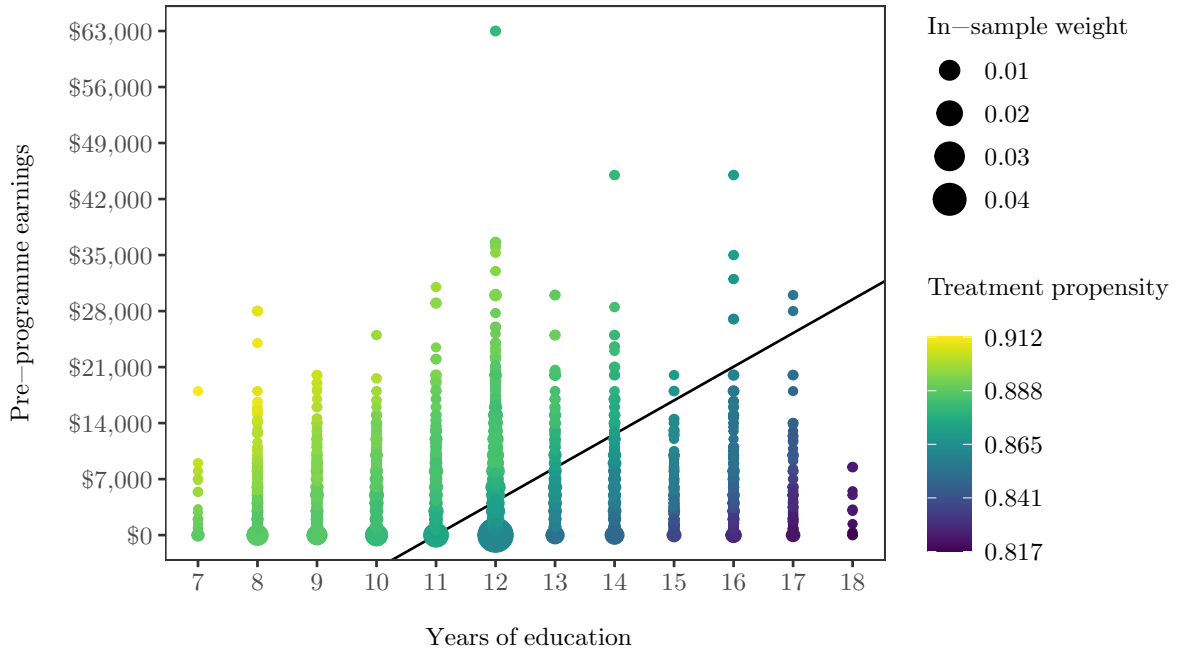
$$g(\mathbf{X}) = 1(\mathbf{X}^\top \{+0.552, +0.533, -0.641\} \geq 0) \tag{D.83}$$

which we plot in Figure D.12. We note that the coefficients of the linear specification are so large because our simple regressions scale individual characteristics but do not scale post-programme earnings, which inflates the effect of pre-programme earnings and years of education. This constitutes our first experiment.

We then propose two further experiments. Our second experiment increases the variance of post-programme earnings by inflating σ_1 , σ_0 and c_v in Equation D.82 by a factor of five. Our third experiment similarly reduces the variance of post-programme earnings by deflating σ_1 , σ_0 and c_v in Equation D.82 by a factor of five. We note that both experiments leave the oracle policy unchanged from Equation D.83. We present the corresponding estimates of the parameters of the posterior distribution in Figure D.13.

An immediate conclusion that we can draw from these results is that μ^* does not align with the oracle

FIGURE D.14
 (Experiment 4) Variation in treatment propensity across individuals



Simulated data generated using Equation D.84 in conjunction with data from the JTPA Study sample. This figure illustrates the treatment propensity of individuals in a simulated sample under the posterior assignment rule that is induced by $\{\kappa^*, \mu^*\}$. Each point represents the individual characteristics of an individual or several individuals. For comparison, individuals to the left of the solid diagonal line are assigned treatment under the oracle assignment rule of Equation D.85.

assignment rule. That being said, it is apparent from Figure D.12 that individuals who are assigned treatment under the oracle assignment rule are more likely to be assigned treatment under the posterior distribution that we obtain.

We plot the behaviour of the objective function and of empirical welfare risk in Figure D.13. We emphasise the non-smoothness of empirical welfare risk for the second experiment, and we suggest that increasing the variance of post-programme earnings makes the problem of finding the optimal assignment rule more difficult.

—| SUBAPPENDIX D.3.C |—

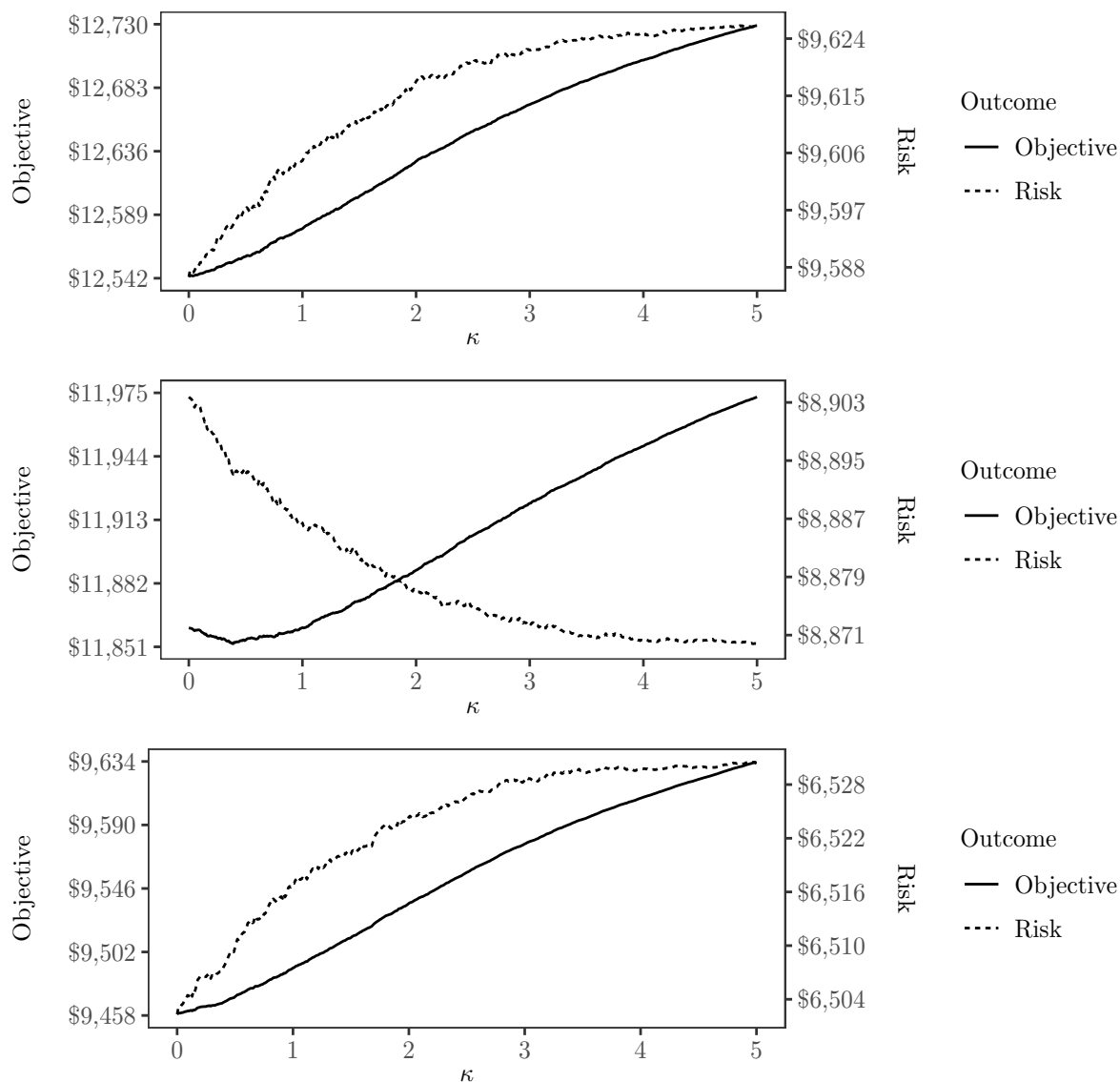
THE LOCATION OF THE ORACLE ASSIGNMENT RULE

A common feature of the data generating process outlined in Equation D.82 for Experiments 1 through 3 is the oracle assignment rule, which is located in the south-east corner of the covariate space (Figure D.12). Two things can be inferred from this. First, that only a small minority of individuals are not assigned treatment under the oracle assignment rule; specifically, those individuals who are educated to

FIGURE D.15

(Experiments 4–6) Behaviour of the objective function at μ^* given variation in κ

Experiment	μ^*	κ^*
4. Baseline (top panel)	{+0.932, +0.362, -0.035}	0.000
5. High variance (middle panel)	{-0.810, -0.526, -0.259}	0.390
6. Low variance (bottom panel)	{+0.899, +0.425, +0.105}	0.000



Simulated data generated using Equation D.84 in conjunction with data from the JTPA Study sample. This figure illustrates how the shape of the objective function and its risk component changes as the variance of post-programme earnings is increased and decreased. We hold μ fixed at μ^* in each case, which differs according to the sample size, and vary κ . Parameter estimates are presented above.

post-graduate level.²⁰ In other words, only a small subset of the data with extreme characteristics are not assigned treatment. Second, that, all else being equal, average outcomes in the experimental group are higher than average outcomes in the control group, and that this baseline difference is relatively important as compared to pre-programme earnings and years of education. These features, we suggest, are indicative of a relatively easy problem (of determining the optimal assignment rule).

The following three experiments largely replicate Experiments 1 through 3 but shift the oracle assignment rule to the left, narrowing the difference between the baseline average outcomes in the experimental and control groups. In Experiment 4, we assume that

$$\begin{aligned} \alpha_1 &= \{+2, 442; +86, 446; +14, 008\} \\ \alpha_0 &= \{-489; +82, 458; +18, 804\} \end{aligned} \quad \text{with } \sigma_1 = \sigma_0 = 15, 914 \text{ and } c_v = 5\sigma_1 = 5\sigma_0 \quad \text{D.84}$$

so that the optimal policy is

$$g(\mathbf{X}) = 1(\mathbf{X}^\top \{+0.425, +0.579, -0.696\} \geq 0) \quad \text{D.85}$$

which we plot in Figure D.14. Experiments 5 and 6 then mirror Experiments 2 and 3 in that they inflate and deflate, respectively, the variance of post-programme earnings. We suggest that Equation D.84 is a more difficult problem than Equation D.82.

We specifically design these experiments so that the oracle assignment rule is such that as close to 50% of the sample is assigned treatment as is possible without altering the importance of pre-programme earnings relative to years of education in the earnings process. We present the corresponding estimates of the parameters of the posterior distribution immediately in Figure D.15.

A curious feature of these results is that, for Experiments 4 and 6, the posterior distribution is uniform. We note that the mean direction is irrelevant in this case, but we include it anyway because knowledge of it is necessary to interpret Figure D.15. This uniformity of the posterior distribution is apparently driven by how empirical welfare risk increases alongside the concentration in Figure D.15. A possible explanation for this feature is that this data generating process is indeed hard. We suggest that the narrow gap between the baseline average outcomes of the experimental and control groups makes treating everyone versus treating no-one equally appealing (or unappealing). For instance, contrasting Figure D.7 with its analogue in the main text (for the raw data), we see that the difference in empirical welfare risk between the treat everyone rules and treat no-one rules that occupy the southern and northern

²⁰ Or, more precisely, that are educated for an equivalent amount of time as would be required to obtain a post-graduate degree, since we do not observe education level.

regions of the heatmaps, respectively, narrows. A similar effect happens here. Given the fact that the density contours of the von Mises-Fisher distribution are concentric, there is then no natural region of the sphere to locate in order to minimise empirical welfare risk unless the value of the concentration parameter is particularly large and the posterior distribution allocates substantial probability mass to an extremely localised area on the frontier between the collections of treat everyone rules and treat no-one rules. High concentration is heavily penalised by the Kullback-Leibler divergence though, and so uniformity cannot be improved upon. The fact that a different pattern is observed for Experiment 5 is compatible with this argument, and could be achieved if the high variance of post-programme earnings makes a few observations in the sample pivotal (recall that the weight attached to each observation is proportional to post-programme earnings; inflating the variance of post-programme earnings can lead to greater concentration of weight on a few observations, since extreme outliers are more likely).

—| SUBAPPENDIX D.3.D |—

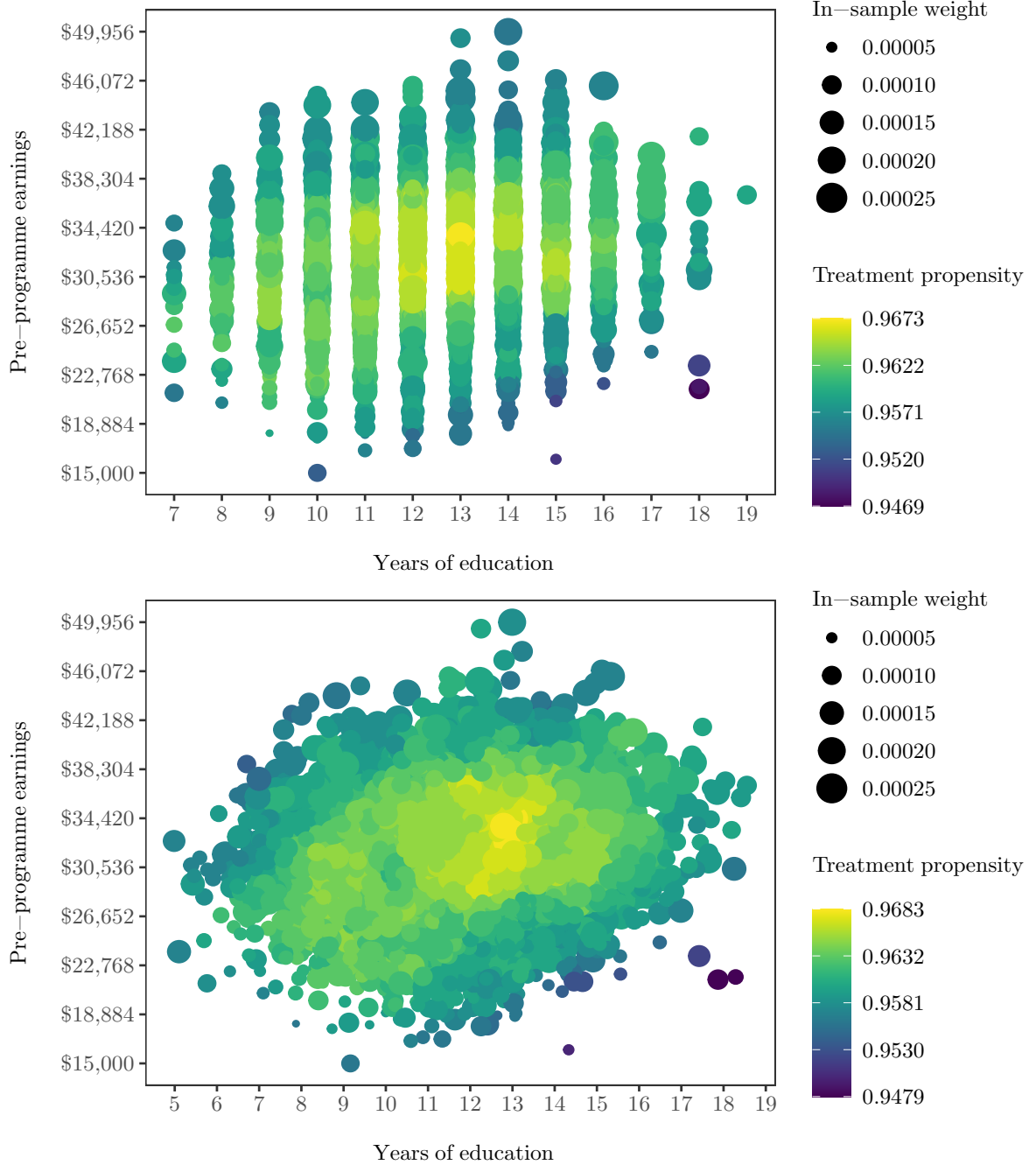
INDIVIDUAL CHARACTERISTICS AND THEIR DISTRIBUTION

The final experiments that we undertake vary the distribution of individual characteristics. Our intention is to understand whether the JTPA Study sample is, in some sense, special and whether altering individual characteristics substantially alters our results.

Experiments 7 and 8 mirror Experiment 1 in how the outcome is generated: both experiments rely on Equation D.76 and the parameter values that we outline in Equation D.82 for Experiment 1. The distinction between Experiments 7 and 8 and Experiment 1 is how individual characteristics are generated. Whereas Experiment 1 uses data taken from the JTPA Study sample, Experiments 7 and 8 generate individual characteristics according to a bivariate normal distribution with the mean vector equal to half of the maximum of pre-programme earnings and the average number of years of education, respectively. We estimate the covariance matrix of pre-programme earnings and years of education in the JTPA Study sample, and set the covariance of the bivariate normal distribution equal to this (Pearson correlation coefficient of 0.126). In the case of Experiment 7, we discretise years of education by assigning each observation to one of 12 equal-sized bins. We then map each characteristic to the unit interval by means of the aforementioned linear transformation.

We plot the individual characteristics that we use in Experiments 7 and 8 in Figure D.16. Absent from either plot is the oracle assignment rule. This is not an oversight. Rather, the oracle assignment rule is such that it recommends that all individuals be assigned treatment (i.e., it to the south-east of the plotting area). We plot the objective function and empirical welfare risk for Experiments 7 and 8 in Figure D.17. Due to the similarity of Figure D.17 to the other figures that we have presented, we

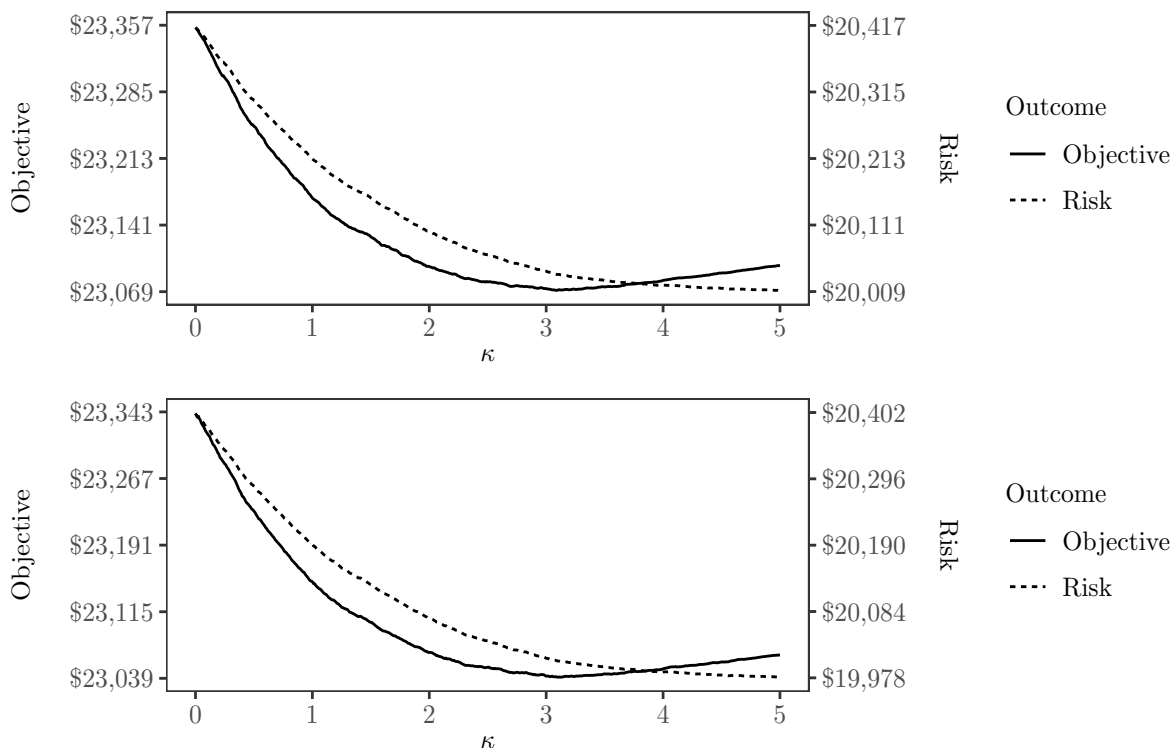
FIGURE D.16
 (Experiments 7–8) Variation in treatment propensity across individuals



Simulated data generated using Equation D.82 and a bivariate normal distribution. This figure illustrates the treatment propensity of individuals in simulated samples under the posterior assignment rule that is induced by $\{\kappa^*, \mu^*\}$. Each point represents the individual characteristics of an individual or several individuals. Every individual is assigned to treatment under the oracle assignment rule.

FIGURE D.17
 (Experiments 7–8) Behaviour of the objective function at μ^* given variation in κ

Experiment	μ^*	κ^*
7. Discrete (top panel)	{+0.794, +0.405, +0.454}	3.090
8. Continuous (bottom panel)	{+0.801, +0.408, +0.438}	3.120

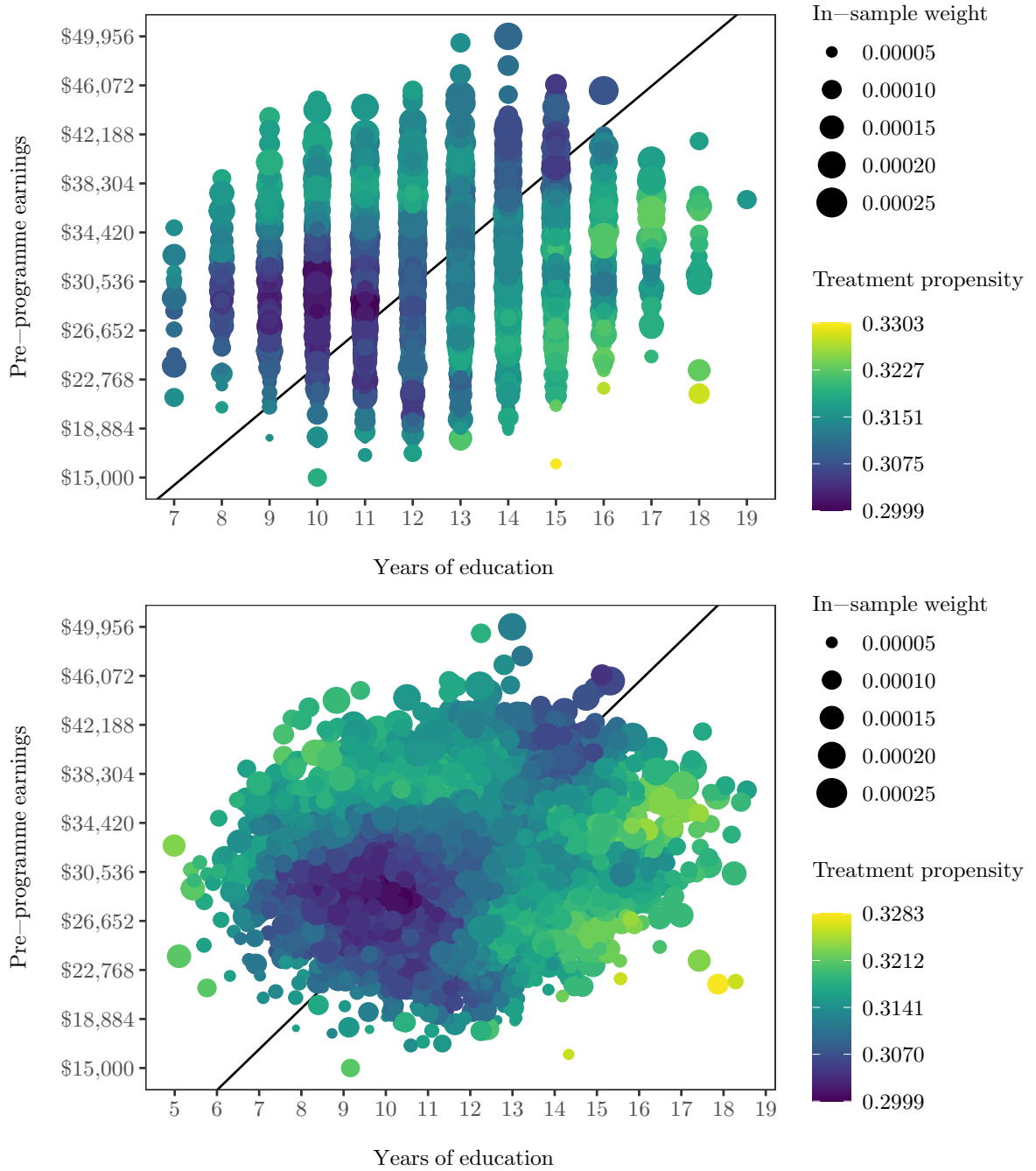


Simulated data generated using Equation D.82 and a bivariate normal distribution. This figure illustrates how the shape of the objective function and its risk component at μ^* changes as κ increases. Parameter estimates are presented above.

do not provide any further discussion of this figure. What is somewhat interesting though is how the propensity of treatment is highest for those individuals with individual characteristics located around their mean values. We present the corresponding estimates of the parameters of the posterior distribution in Figure D.17. We find that μ^* does not align with the oracle assignment rule, but does dictate that all individuals are assigned treatment (as they are under the oracle assignment rule); κ^* is also relatively large, as compared to its value in the previous experiments.

So as to make the assignment problem harder—at least, what we understand to be harder—we propose Experiments 9 and 10. These experiments follow Experiments 7 and 8 in how individual characteristics are generated, but differ slightly in how they generate the outcome. Whilst Experiments 9 and 10 broadly follow Experiments 7 and 8 with respect to how the outcome is generated, they shift the process for the potential outcomes and, thereby, the oracle assignment rule. Specifically, in Experiments 9 and

FIGURE D.18
 (Experiments 9–10) Variation in treatment propensity across individuals

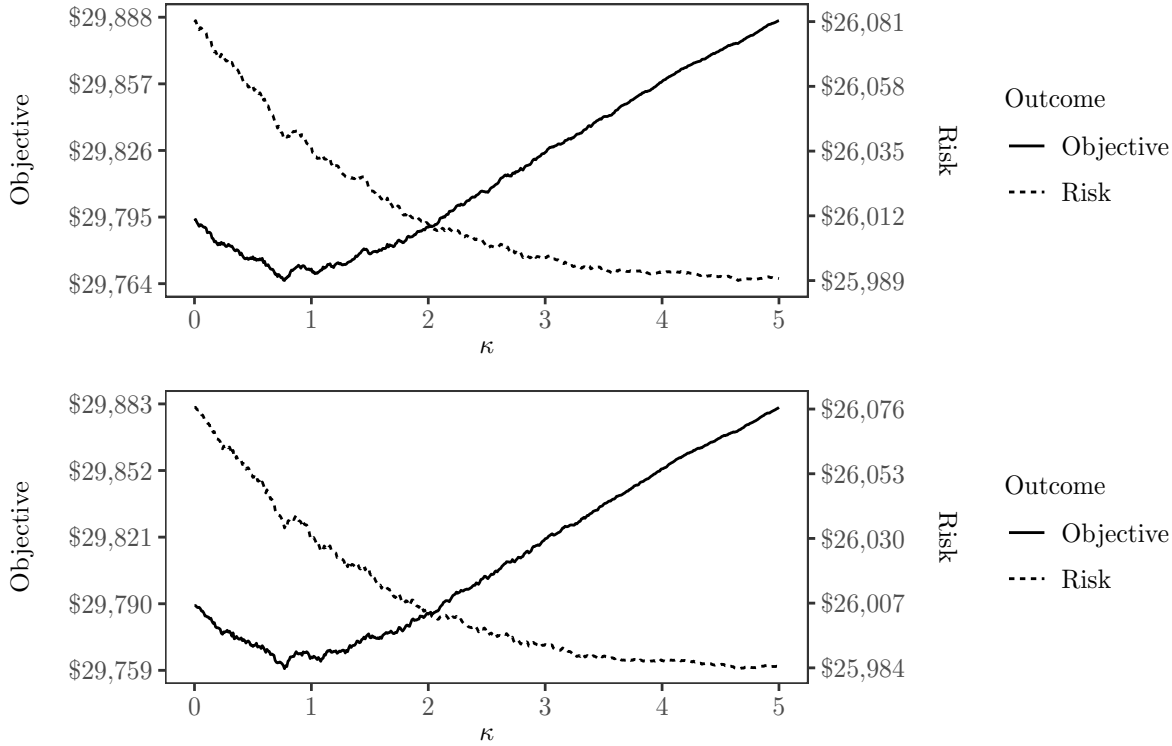


Simulated data generated using Equation D.86 and a bivariate normal distribution. This figure illustrates the treatment propensity of individuals in simulated samples under the posterior assignment rule that is induced by $\{\kappa^*, \mu^*\}$. Each point represents the individual characteristics of an individual or several individuals. For comparison, individuals to the left of the solid diagonal line are assigned treatment under the oracle assignment rule, which partitions the sample into two approximately equal-sized groups.

FIGURE D.19

(Experiments 9–10) Behaviour of the objective function at μ^* given variation in κ

Experiment	μ^*	κ^*
9. Discrete (top panel)	$\{-0.766, -0.470, -0.438\}$	0.770
10. Continuous (bottom panel)	$\{-0.766, -0.470, -0.438\}$	0.770



Simulated data generated using Equation D.86 and a bivariate normal distribution. This figure illustrates how the shape of the objective function and its risk component at μ^* changes as κ increases. Parameter estimates are presented above.

10, we assume that

$$\begin{aligned} \alpha_1 &= \{+1, 286; +86, 446; +14, 008\} \\ \alpha_0 &= \{-668; +82, 458; +18, 804\} \end{aligned} \quad \text{with } \sigma_1 = \sigma_0 = 15, 914 \text{ and } c_v = 5\sigma_1 = 5\sigma_0 \quad \text{D.86}$$

such that the oracle policy is

$$g(\mathbf{X}) = 1(\mathbf{X}^\top \{+0.098, +0.636, -0.765\} \geq 0) \quad \text{D.87}$$

The effect of this change is to maintain the slope of the oracle assignment rule (i.e., the contribution of pre-programme earnings relative to years of education) but to alter the intercept. In other words, to narrow the difference in baseline average outcomes between the experimental and control groups. In this regard, Experiments 9 and 10 are similar in intent to Experiment 4.

We plot the individual characteristics that we use in Experiments 9 and 10 in Figure D.18. The oracle assignment rule is designed to partition the sample into approximately two groups of equal size. We plot the objective function (left-hand panel) and empirical welfare risk (right-hand panel) for Experiments 9 and 10 in Figure D.19. We draw attention to several differences between the results that we obtain for Experiments 9 and 10 versus Experiments 7 and 8. First, we highlight the fall in the average propensity of treatment, as is apparent in Figure D.18. Whereas in Experiments 7 and 8 (and the preceding experiments too), the propensity of treatment is close to one (i.e., everyone is assigned treatment), here the propensity is closer to one third. There is more randomisation. Second, individuals who have a relatively high propensity of treatment as compared to their peers in Experiments 7 and 8 have a relatively low propensity of treatment in Experiments 9 and 10 (visually, the colours are inverted), as is apparent in Figure D.18. That is, individuals whose characteristics are close to the mean values have a low propensity of treatment. Third, the value of the concentration parameter decreases from above three to below one (i.e., the posterior distribution is less concentrated and more uniform), as is apparent in Figure D.19. We present the corresponding estimates of the parameters of the posterior distribution in Figure D.19. Interestingly, the mean direction is, in both cases, almost the exact opposite of the corresponding mean direction in Experiments 7 and 8 (i.e., the mean direction is approximately the negative of the mean direction in Experiments 9 and 10).

THE EFFECT OF ADDITIONAL CHILDREN ON MATERNAL LABOUR SUPPLY

The principal objective of this chapter is to estimate the effect of additional children on maternal employment under weak restrictions on the assumed data generating process. I frame this question in the context of the U.S. labour market, exclusively studying mothers with at least two children. I exploit natural variation in the characteristics (specifically, the birthdate and sex) of existing children to identify the employment response of mothers to additional children. The model that I employ is a minimally restrictive single equation model that incorporates an instrumental variable restriction; the model features a dichotomous outcome and a dichotomous treatment, and is partially identifying.

Understanding the sign and magnitude of the employment response of mothers to additional children is important for the design and evaluation of policy. Although family size is often, in itself, not a stated objective of policy and is not necessarily (directly) manipulable in any case, the effect of additional children on maternal employment is highly relevant to policymakers for its influence on present and future aggregate outcomes. Do labour supply responses amongst working-age mothers exacerbate or offset the long-term effect of falling fertility on public finances? Children are an important feature of the tax code and tax policy is often designed with families with children in mind—recent budget announcements indicate an intention by government to provide or increase childcare provision with the intention of improving the incentive to work amongst parents.¹ Yet, do children increase or reduce maternal employment? Understanding how children affect maternal labour supply is necessary to understand both how the tax and benefit schedule affects the household budget constraint, and whether these policies are the most suitable for achieving their stated aims. Whilst low maternal employment relative to non-maternal employment (see OECD, 2007) is a concern of policymakers, it is unclear whether this effect is driven by

1 The budget for the 2018 fiscal year included the provision of tax relief to help families struggling with child and dependent care expenses, on top of the existing Earned Income Credit that lowers marginal tax rates for working-age mothers with children in low-income groups. Since then, the American Rescue Plan and the President’s (proposed) Budget for the 2024 fiscal year include an expansion of child tax credits and of access to affordable childcare. Similarly, in the U.K., the 2017 Spring Budget included the provision of financial support to parents for each child under the age of 12 years and of free childcare allowance of up to 30 hours per week to working parents with a three or four year-old child, with the express intention of providing support for ordinary working families, and for women in the workplace. The 2023 Spring Budget extended the childcare allowance to children aged between nine months and two years of age.

children or by the pre-existing opportunities of mothers for employment relative to non-mothers.

Existing empirical evidence as to the effect of children on maternal employment is mixed, despite apparent similarities in sample construction—restricting attention to married mothers with at least two children so as to exploit variation in family size that is associated with the birthdate and sex of these children—across works.

ANGRIST and EVANS (1998), IACOVOU (2001), and AL-KHAJA (2016) adopt a reduced form perspective, recovering a Local Average Treatment Effect (IMBENS and ANGRIST, 1994) that is the effect of children on maternal employment amongst an unidentifiable subset of the sample population. Importantly, this subset is not invariant and is sensitive to the institutions and cultural norms that are present, and can lead to estimates that are not necessarily representative of the effect of children across a population. For instance, ANGRIST and EVANS (1998) finds that additional children reduce the probability that a mother is employed in the context of the U.S. labour market, whereas IACOVOU (2001) and AL-KHAJA (2016) find the opposite in the context of the British and Egyptian labour markets, respectively. It is unclear whether these reported differences represent an actual difference in maternal labour supply across the three countries, or whether they are due to differences in preferences over family composition and who is affected by the plurality or sex of a child.

I, instead, adopt a structural perspective, and focus on identifying and estimating an Average Treatment Effect—and its primitives, the Average Structural Functions (BLUNDELL and POWELL, 2004)—that can inform both *ex-ante* policy design and *ex-post* policy evaluation.^{2,3} A similar approach is taken by CHESHER and ROSEN (2020)—that subsumes CHESHER and ROSEN (2013)—with the stated aim of illustrating how to apply minimally restrictive models to data and to demonstrate what such models deliver.

The model that I study generalises the non-parametric model of treatment non-compliance that is studied by BALKE and PEARL (1997; hereafter, the Balke–Pearl model). The Balke–Pearl model permits any form of dependence between the employment decision and the fertility decision in the context of maternal labour supply, and constitutes a useful benchmark. Put differently, the Balke–Pearl model does not preclude the existence of any response types. This property contrasts favourably with control function approaches (see WOOLDRIDGE, 2015 for a discussion of control functions) and special regressor methods (see RACINE et al., 2014 §I.2 for a discussion of special regressors) that either restrict these unobservable components of the economic environment or rely upon continuous variation that may not be present.

2 The Average Treatment Effect is sometimes called the Average Causal Effect.

3 The Average Structural Functions are the means of the so-called potential outcome distributions introduced in D. B. RUBIN (1974)—and arguably preceding that in SPLAWA-NEYMAN (1923).

The inherent freedom of the Balke–Pearl framework does, however, mean that multiple behaviours—the measures of each of which combine to constitute the object of interest—are associated with a single outcome, which is the issue at the core of the anatomy of the selection problem (MANSKI, 1989). The Balke–Pearl model only partially identifies the Average Treatment Effect and its primitives, and this property is shared by the model that I study.

I generalise the Balke–Pearl model in two ways. First, I allow covariates to directly influence both response and treatment. For example, I allow mothers of different ethnicities to hold different preferences over how many children to have and whether to work. Second, I allow the instrument to have discrete—rather than dichotomous—support. For example, I allow for the birthdate and sex of children to be used in combination. The model—or, more precisely, models—that CHESHER and ROSEN (2020) studies also embeds these features, but maintains a monotonicity restriction that I do not. By relaxing the assumption of weak separability (or monotonic response) that is made in that paper,⁴ I allow for a richer set of behaviours. In particular, I allow some mothers to move into employment given an increase in family size, and for others to move out-of employment given a similar increase. Such behaviour is ruled-out by monotonic response, but is an important feature of the model of GRONAU (1977); there, differences in the endowments of and prices that are faced by families generate different employment responses at the extensive margin. To accommodate both generalisations, I extend the statistical independence restriction that the Balke–Pearl model embeds.

Aside from relaxing the assumption of weak separability, I also extend CHESHER and ROSEN (2020) by providing more recent evidence of the effect of additional children on maternal employment. Whereas CHESHER and ROSEN (2020) exclusively uses data from 1980—the original sample considered by ANGRIST and EVANS (1998)—I additionally use data from 2010 and 2015, and so can provide some insight into whether this effect has changed over time. I do, however, differ from CHESHER and ROSEN (2020) in which covariates I include, focusing on age and ethnicity or race rather than education level.

An interesting property of the Balke–Pearl model is the lack of information content that the treatment equation provides; the response equation and the statistical independence restriction constitute a sufficient representation of the model. In this sense, both the Balke–Pearl model and the model that is studied in CHESHER and ROSEN (2020) can be expressed as single-equation models, with the only difference between them being the assumption of weak separability that is present in the analysis of CHESHER and ROSEN (2013). Monotonicity is a powerful restriction, in that it effectively guarantees identification of a non-zero effect and, in many cases, the sign of a treatment effect. MOURIFIÉ (2015) also studies a tri-

⁴ I reiterate that the stated aim of CHESHER and ROSEN (2020) is to illustrate how to apply minimally restrictive models to data and to demonstrate what such models deliver. Imposing weak separability—irrespective of its merits—is compatible with this objective, since it leads to a simpler and more intuitive identification analysis.

angular model (i.e., a model featuring both a response equation and a treatment equation, as laid-out in STROTZ and WOLD, 1960), and imposes weak separability on both of these equations. MOURIFIÉ (2015) states sharp bounds on the Average Treatment Effect and its primitives that are robust to the failure of the support condition that is required in SHAIKH and VYTLACIL (2011). Unlike for the Balke-Pearl model, however, the treatment equation is not redundant and has information content.

RICHARDSON and ROBINS (2014) extends BALKE and PEARL (1997) in allowing instruments to be discrete rather than dichotomous, and is the closest paper to this work. RICHARDSON and ROBINS (2014) also varies the strength of the statistical independence restriction therein, with the presented bounds coinciding with those that I obtain when covariates are assumed to have a single point of support. KITAGAWA (2021) extends BALKE and PEARL (1997) in allowing for continuous response and varies the strength of the statistical independence restriction therein. The model that is studied in KITAGAWA (2021) is, like the Balke-Pearl model that it nests, partially identifying and is falsifiable via the generalised instrumental inequalities of KÉDAGNI and MOURIFIÉ (2020). CHESHER and ROSEN (2017) provides a general characterisation of sharp bounds on the identified set for a large of class of models and a wide choice of criteria using results from random set theory, with the analysis therein able to replicate the findings of BALKE and PEARL (1997) and KITAGAWA (2021). Although this characterisation equally applies when treatment is discrete or continuous, GUNSILIUS (2020) advances a path-sampling approach as a more tractable alternative for triangular models featuring large support.

A common issue for the application of the aforementioned models and their associated methods is that naïve estimation of the Average Treatment Effect and other partially identified objects is typically biased. I adopt a frequentist perspective and a method described in CHERNOZHUKOV et al. (2013) that is half-median unbiased. CHERNOZHUKOV et al. (2013) is but one paper that discusses inference on moment inequalities—other papers in the literature include ANDREWS and SOARES (2010), and BUGNI (2010), and references therein. One paper of particular note is KAIDO et al. (2019), which proposes a bootstrap procedure to correct for projection conservatism. Projection conservatism—that confidence regions do not have the correct size, and contain the truth in more than the pre-specified proportion of samples—is of concern in partially identifying models where structural characteristics of interest are not identified directly but as the projection of an identified set of several parameters. For instance, the Average Treatment Effect is identified as a projection of the Average Structural Functions (their difference).

An alternative approach is to undertake Robust Bayesian inference, which is the subject of GIACOMINI and KITAGAWA (2021b) and KLINE and TAMER (2016). The principal distinction between robust Bayesian inference and Bayesian inference is that the former specifies multiple priors. MOON and

SCHORFHEIDE (2012) shows that Bayesian credible regions differ substantially from frequentist confidence regions, even asymptotically, in partially identifying settings. Robust Bayesian credible regions though do not share this property; KITAGAWA (2021) shows that robust Bayesian inference using the lower envelope (the posterior lower probability) has a frequentist interpretation under certain conditions.⁵ Attractive features of robust Bayesian inference are that they do not suffer from projection conservatism and are often less expensive than frequentist approaches that rely on grid-based inversion of a test-statistic. I do not adopt a Bayesian perspective (robust or otherwise) in this chapter.

—| SECTION 1 |—

FRAMEWORK

I suppose that the social planner (or the econometrician acting on her behalf) observes experimental or observational data that comprises a probability distribution over $\{Y, T, Z, X\}$ that I denote by P^n and that is constructed from n independent and identically distributed draws. Here, $Y \in \{0, 1\}$ denotes response, which is some measured binary outcome of interest; $T \in \{0, 1\}$ denotes treatment, which is an indicator for some measured binary choice; $Z \in \mathbb{R}$ denotes an instrument, which is an indicator for the occurrence of some event; and $X \in \mathbb{R}$ denotes covariates, which are some measured individual characteristics. The population from which the experimental or observational data is drawn comprises a probability distribution over $\{U, Z, X\}$ that I denote by P . Here, U is unobserved heterogeneity that relates to Y and T via the relationship

$$\begin{aligned} y &= 1(s_y(t, z, x, u) \geq 0) \\ t &= 1(s_t(y, z, x, u) \geq 0) \end{aligned} \tag{E.1}$$

I refer to P^n as the empirical distribution and to P (together with Equation E.1) as the data-generating process.

Throughout my analysis, I maintain several assumptions that I collectively refer to as the generalised Balke-Pearl model for their resemblance to the assumptions that the eponymic model embeds.

Assumption E.1 (Exclusion). *The data generating process satisfies*

$$\begin{aligned} s_y &= h_y(t, x, u) \\ s_t &= h_t(z, x, u) \end{aligned} \tag{E.2}$$

such that the generalised Balke-Pearl model is triangular in Y , T and Z .

⁵ KITAGAWA (2021) notes that at least one of these conditions is fragile; equality of bounds, such as when an instrumental variable has no influence over treatment, invalidates this frequentist interpretation.

The Balke-Pearl model is a structural equation model, such that each structural equation specifies a causal relationship. Given its recursive nature, the Balke-Pearl model induces a probability distribution that is Markovian; together with the statistical independence restrictions that it embeds, I note that the Balke-Pearl model can instead be presented as a (Pearlian) Directed Acyclic Graph (DAG; see PEARL, 2009 for a comprehensive treatment of the equivalence of structural equation models and DAGs, and DAWID, 2008 for a discussion of the difference between Pearlian DAGs and probabilistic DAGs). The generalised Balke-Pearl model also has this dual representation.

Assumption E.2 (Discrete support). *The data generating process satisfies $\mathfrak{z} < \infty$ and $\mathfrak{r} < \infty$, with $\mathfrak{z} \leftarrow x \geq 2$ holding over \mathfrak{X} .*

Assumption E.3 (Statistical independence). *The data generating process satisfies*

$$U | \{Z, X\} \sim \text{Uniform}(0, 1)$$

E.3

for which h is appropriately normalised.

The generalised Balke-Pearl model exhibits a number of features that make it ideally suited to the study of the effect of additional children on female employment.

First, the generalised Balke-Pearl model permits endogenous selection of T . That is, the model allows for dependence between the employment and fertility decisions of mothers. In particular, it is likely that mothers factor their opportunities for employment into their decision about whether to have additional children, and *vice versa*. For example, mothers who face poor employment opportunities (and a low opportunity cost of pregnancy) relative to others may choose to have an additional child. This story hints at mothers with three or more children also being mothers who are in any case less likely to be employed; whilst this may or may not be true, the existence of some indirect mechanism that generates correlation between response and treatment is widely accepted.

Second, the generalised Balke-Pearl model permits rich heterogeneity of response both across and within sub-populations labelled by X . Heterogeneity of response across sub-populations can be generated by parametric models and so does not warrant much discussion.^{6,7} Of greater interest, however, and warranting more-detailed discussion, is heterogeneity of response to treatment within sub-populations. The

6 As a simple example, the generalised Balke-Pearl model can allow Hispanic and non-Hispanic mothers (if ethnicity, or race, is an observable individual characteristic) to have different employment rates and to exhibit different effects of additional children on maternal employment (as an extreme case, if the effect of additional children on maternal employment amongst Hispanic mothers has the opposite sign to that for non-Hispanic mothers, say).

7 The non-parametric specification does, however, mean that X can always be written as a scalar without loss of generality, a restatement that is with loss in a parametric specification.

flexible nature of both the response equation and the treatment equation means that the generalised Balke-Pearl model can allow different mothers to experience different effects of additional children. Broadly speaking, a mother may increase her labour supply in response to additional children, may decrease her labour supply, or may leave her labour supply unchanged. The generalised Balke-Pearl model can generate all three forms of response. In particular, I emphasise that response can be non-monotonic, in that some mothers may increase their labour supply whilst others decrease their labour supply. That mothers may exhibit opposite responses at the extensive margin is due to non-separability of the response equation in unobserved heterogeneity.⁸ Such behaviour is ruled-out by monotonic treatment response, but is an important feature of the model of GRONAU (1977); there, differences in the endowments of and prices that are faced by families generate different employment responses at the extensive margin. For example, it is intuitive to think of children as both affecting the cost of employment (the need to provide childcare) and the value of employment (children increase the marginal utility of income); if the first effect dominates the second then it is feasible that a mother increases her labour supply in response to additional children, but the opposite will occur if the reverse is true.

Third, the generalised Balke-Pearl model features a lack of restrictions on both the dimension and the distribution of unobserved heterogeneity. Underpinning this feature of the generalised Balke-Pearl model is the flexibility of the structural equations, which I discuss in Section 2; specifically, that unobserved heterogeneity may be partitioned into a finite number of equivalence classes. Generally, restrictions on unobserved heterogeneity are not verifiable, and may impose non-trivial constraints on the economic environment—particularly in parametric models. These constraints may include restrictions on the distributions of endowments and of prices. The credibility of the generalised Balke-Pearl model is that it imposes no constraint of the shape of these distributions, and can allow for a large class of features of the economic environment to influence the employment decision and the fertility decision.

—| SUBSECTION 1.A |—

THE MEANING AND CREDIBILITY OF THE RESTRICTIONS

Although the generalised Balke-Pearl model is minimally restrictive, it is not entirely non-restrictive. In particular, the model embeds several explicit restrictions—Assumptions E.1 to E.3—on the data generating process and a further restriction that is implicit—a rank condition. I discuss each of these assumptions in turn, suggesting an interpretation for each and scenarios that would lead to their violation.

The exclusion restriction—Assumption E.1—can be restated as $Y \perp\!\!\!\perp Z \mid \{T, X, U\}$, so precluding the

⁸ VYTLACIL (2002) shows that monotonicity is equivalent to weak separability, which precludes one of these responses.

instrument from directly influencing response (contingent upon its antecedents, response is deterministic). In the context of maternal employment, the exclusion restriction requires that having twins or higher plurality siblings does not directly influence the employment decision. Violation of the restriction occurs if twins incur additional costs versus non-plural siblings; for example, if having children of the same age requires duplication, versus handing-down or re-using. This reasoning similarly applies if the sex of existing children is used in place of plurality of birth; for example, if female and male children participate in different activities at distinct venues, or if they attend different schools. The argument in both scenarios is that the instrument—whether a plural second birth or the sex of existing children— influences the employment decision via its effect on the budget constraint, and not simply through its influence on the decision of whether to have additional children.

The rank condition can be stated as $T \not\perp Z$, so that the treatment equation is a non-trivial function of the instrument. In the context of maternal employment, a plural second birth means that a mother has three or more children, and so automatically satisfies this requirement by definition. If the sex of existing children is used in place of plurality of birth, violation of the restriction occurs if parents either do not hold a preference over the sexes of their children or if parents are equally balanced over all possible combinations (i.e., as many parents prefer mixtures of female and male children as prefer just female or just male).

The discrete support restriction—Assumption E.2—is a statement about both the amount of variation in the instrument and the sort of covariates that can be incorporated. The discrete support restriction does not, however, otherwise impose any constraint on the marginal distribution of the instrument and the covariates—pairs of realisations can exhibit dependence or not occur at all. As a restriction on the observable exogenous economic variates, the discrete support restriction is verifiable (i.e., it can be tested and either accepted or rejected). In the context of maternal employment, the discrete support restriction requires that there are at least some mothers for every combination of age, age at first birth, and ethnicity and race that experience a plural second birth. Violation of the restriction occurs if a plural second birth is too rare an event; for example, if the aforementioned characteristics are so finely granulated that there are fewer plural second births than there are groups. If the sex of existing children is used in place of plurality of birth, violation of the restriction is unlikely due to the natural sex ratio (i.e., it is nearly—although not equally—as likely that a mother has a female child as a male child).

The statistical independence restriction—Assumption E.3—can be restated as $U \perp\!\!\!\perp \{Z, X\}$, since the distributional assumption (i.e., that unobserved heterogeneity is uniformly distributed) is simply a normalisation accorded by the non-parametric nature of the structural equations. In the context of maternal employment, the statistical independence restriction implies that mothers do not base their age, age at

first birth, and ethnicity and race on their opportunities for employment, nor is the probability of a plural second birth dependent upon these opportunities. Given the immutable and predetermined nature of these characteristics, there is strong reason to suppose that this restriction is satisfied in this regard. Of greater concern is whether any of these economic variates are statistically associated with a further variable that influences Y but that is not included explicitly in the model. For example, the statistical independence restriction is violated if education is statistically associated with age at first birth or ethnicity and race, since education likely influences the probability of employment (through spousal income if more highly educated mothers marry equally educated men, and through the offered wage). This reasoning similarly applies if the sex of existing children is used in place of plurality of birth.

Together, these restrictions define Z as an instrumental variable (for the effect of treatment on response).⁹ I note that the statistical independence restriction may be weakened; specifically, the statistical independence restriction may be weakened by partitioning U into a part that enters the response equation and a part that enters the treatment equation. A sufficient restriction is then that Z is statistically independent of that part of U that enters the response equation, although it can be statistically associated with the part that enters the treatment equation. KITAGAWA (2021) further weakens this restriction to marginal statistical independence, but it is unclear what the interpretation of this restriction is in a structural equation model.

—| SECTION 2 |—

EQUIVALENCE CLASSES AND LATENT TYPES

The generalised Balke-Pearl model is a structural equation model, such that each structural equation specifies a causal relationship. The response equation and the treatment equation then fully describe behaviour; in particular, the outcome equation describes maternal employment for both the realised treatment and the unrealised treatment. The potential outcome notation that is developed in D. B. RUBIN (1974) advances this idea by defining both the realised outcome and the unrealised outcome as variates, and is complementary to the structural approach that I take in this chapter. I adopt the potential outcome notation briefly so as to illustrate an important point that is made in PEARL (2009).

$$y_t(x) \doteq 1(h_y(t, x, u) \geq 0)$$

$$t_z(x) \doteq 1(h_t(z, x, u) \geq 0)$$

E.4

9 This definition of an instrumental variable is consistent with the definition that is stated in PEARL (2009).

From the definition of the potential outcomes, it is evident that the potential outcomes share the support of the response and of treatment—the potential outcomes are dichotomous random variables. Moreover, all variation in the potential outcomes is due to variation in U , since each potential outcome designates t , z and x . It immediately follows that unobserved heterogeneity may be collected into sets that deliver common values of the potential outcomes—equivalence classes.

An equivalence class is a set in u -space such that any two members of a given class induce the same potential outcomes as each other for every possible designation of the observable exogenous random variates. I note that this definition of an equivalence class does not rely on a support condition (for unobserved heterogeneity)—the definition is valid regardless of the dimension and distribution of unobserved heterogeneity. This observation is consistent with the claim that is made in BALKE and PEARL (1997) that the Balke-Pearl model can allow for unobserved heterogeneity of arbitrary dimension. Each equivalence class collects values of unobserved heterogeneity that deliver common values of the potential outcomes. Given that each potential outcome may take a finite number of values and there are a finite number of potential outcomes, unobserved heterogeneity may be collected into a finite number of equivalence classes (PEARL, 2009). The number of equivalence classes increases exponentially in \mathfrak{z} and \mathfrak{r} .

It is convenient to think of each equivalence class as a latent type (PEARL, 2009), in that each class describes a common set of behaviours. By appropriately defining the structural equations as index functions, the Balke-Pearl model and its generalisation can be restated in terms of response types. Doing so establishes an equivalence between the potential outcome framework and the structural equation framework that is often disregarded by proponents of each approach. Moreover, these latent types can always be projected onto the unit continuum using a probability transformation, thereby confirmation the normalisation that I impose.

The appropriate interpretation of the associated statistical independence restriction that arises in this potential outcome–*cum*–structural equation framework is that latent types are statistically independent of covariates and the instrument. In the context of maternal employment, this restriction is effectively equivalent to the statement that mothers are *ex-ante* identical up to a finite number of latent types and that covariates and the instrument are randomly allocated. Put differently, mothers of the same latent type would make the same employment and fertility decisions but make (possibly) different decisions due to *ex-post* variation in their covariates and the instrument.

IDENTIFICATION

I extend BALKE and PEARL (1997) in allowing covariates to influence the economic system, as a cause of the response and of treatment. A contribution of this chapter is that I state sharp bounds on structural parameters of interest for this case, which represents a simple extension of BALKE and PEARL (1997). These bounds coincide with those that are stated in BALKE and PEARL (1997) for the special case in which the analysis of BALKE and PEARL (1997) is applied to a sub-population x rather than the full population. Such a refinement of the population does not invalidate the analysis of BALKE and PEARL (1997) due to the independence restriction that I invoke (that covariates are statistically independent of unobserved heterogeneity). However, a refinement of the population also alters the object that BALKE and PEARL (1997) bounds; specifically, I state sharp bounds on the Average Treatment Effect in the sub-population x and its primitives (as opposed to the ATE and its primitives). I follow HECKMAN et al. (2006) in shortening the Average Treatment Effect in the sub-population x to $ATE(x)$.¹⁰

An intuitive interpretation of $ATE(x)$ is the change in the employment rate between two counterfactuals, which I write as

$$ATE(x) \doteq \bar{Y}_1(x) - \bar{Y}_0(x) \tag{E.5}$$

The first counterfactual is where all mothers in the sub-population x have three or more children.

$$\bar{Y}_1(x) \doteq \Pr(Y_1(x) = 1|x) \tag{E.6}$$

The second counterfactual is where all mothers in the sub-population x have two children.

$$\bar{Y}_0(x) \doteq \Pr(Y_0(x) = 1|x) \tag{E.7}$$

If the employment rate of the first counterfactual is greater than that of the second counterfactual then $ATE(x)$ is positive, such that additional children increase maternal employment.

Broadly speaking, a mother may increase her labour supply in response to additional children, may decrease her labour supply, or may leave her labour supply unchanged. This motivates a further intuitive interpretation of $ATE(x)$: as the difference in the measures of the first two types of behaviour in the sub-population x . Although these behaviours are distinct from the compliers and defiers that are introduced

¹⁰ This parameter is otherwise referred to as $CATE(x)$ in ABREVAVA et al. (2015) and is closely related to the Average Structural Functions in BLUNDELL and POWELL (2004).

in IMBENS and D. B. RUBIN (1997)—in that they relate to the response equation rather than the treatment equation—they, nonetheless, have a similar flavour to those two concepts. It follows that $\text{ATE}(x)$ is positive if the proportion of mothers who move into employment exceeds the measure of proportion of mothers who move out of employment, such that there is a net increase in employment. These behaviours are exactly defined by the response type. I note that this interpretation of $\text{ATE}(x)$ applies only when the outcome and treatment are dichotomous.

I let $Q_{yt.z}(x) \doteq \Pr(Y = y, T = t|x, z)$, which extends the notation in BALKE and PEARL (1997). The generalised Balke-Pearl bounds are, for all $\{a, b\} \in \mathfrak{Z} \mapsto x$,

$$\left\{ \begin{array}{c} Q_{10.a}(x) \\ (Q_{10.a} + Q_{11.a} - Q_{00.b} - Q_{11.b}) \circ (x) \end{array} \right\} \leq \bar{Y}_0(x) \leq \left\{ \begin{array}{c} 1 - Q_{00.a}(x) \\ (Q_{01.a} + Q_{10.a} + Q_{10.b} + Q_{11.b}) \circ (x) \end{array} \right\} \quad \text{E.8}$$

and

$$\left\{ \begin{array}{c} Q_{11.a}(x) \\ (Q_{00.a} + Q_{11.a} - Q_{00.b} - Q_{01.b}) \circ (x) \end{array} \right\} \leq \bar{Y}_1(x) \leq \left\{ \begin{array}{c} 1 - Q_{01.a}(x) \\ (Q_{00.a} + Q_{11.a} + Q_{10.b} + Q_{11.b}) \circ (x) \end{array} \right\} \quad \text{E.9}$$

These bounds are sharp.¹¹

Although the $\text{ATE}(x)$ and its primitives are informative, they may be of limited policy relevance; policy-makers are often unable to discriminate on the basis of covariates. For instance, equality legislation often prevents policy-makers from discriminating on the basis of ethnicity or race. Of greater importance to policy-makers in such an environment is ATE and its primitives. The policies that were recently announced by the U.K. government, and that are discussed earlier, are examples of policies that are applied uniformly. What is of interest to the policy-maker then is whether the effect of the policy is positive or negative on average across the population. To obtain the aggregate outcome or effect, it is sufficient to take the weighted Minkowski sum over $\text{ATE}(x)$ or its primitives (i.e., take the weighted average over the minimum of each set to obtain the aggregate lower bound, and the weighted average over the maximum of each set to obtain the aggregate upper bound).¹² Here, each weight corresponds to the relative size of the sub-population x in the population, and can be calibrated to the population that is observed or else to some other hypothetical population for *ex-ante* analysis.

11 Proof of this immediately follows from CHESHER and ROSEN (2017). I, nonetheless, include a discussion of the result in the appendices, and a formal derivation in Section 6.

12 The continuous equivalent of the Minkowski sum is the Aumann integral (AUMANN, 1965).

DATA AND SAMPLE COMPOSITION

I use data from the Integrated Public Use Microdata Series (IPUMS; RUGGLES et al., 2015); the data consist of individual-level observations of households in the U.S., and were collected as part of the decennial census of 1980, and the American Community Survey (ACS) of 2010 and of 2015.

A notable difference between the earlier census and the ACS is the timing of collection—whereas the earlier census presents a snapshot of society on April 1st of that year, the ACS is collected throughout a calendar year. The ACS is less informative, in that the birth year that is imputed by the U.S. Census Bureau is incorrect for all individuals who have a birthday that falls after the date of survey. Although this error is random, and averages across the sample, an unfortunate consequence is that it is not necessarily possible to distinguish two children that are born ten or 11 months apart from twins. Furthermore, the ACS does not contain information about the number of children ever born (a measure of fertility and infant mortality), nor does it contain data on the age at or quarter of marriage.

I refine the census and ACS data according to a number of rules. First, I include only women who are aged between 21 and 35 years that are married and whose spouse is present. I exclude women who belong to a married couple where either the woman or her spouse have been married more than once. Second, I include only women who have at least two children. I exclude women who are identified as the adopted- or step-mother of any child, or whose spouse is identified as the adopted- or step-father of any child. I exclude women whose first child is 18 years old or older, or whose second child is less than one year old. Third, I exclude women whose first child was born before her 15th birthday or that of her spouse. Fourth, I exclude women whose age or sex, or that of her children, has been imputed by the U.S. Census Bureau.

I restrict attention to married women for two reasons. First, to retain comparability with ANGRIST and EVANS (1998), which follows many of the same sample selection rules.¹³ Second, by restricting attention to married women, I am able to investigate the effect of children on maternal employment for a stable sub-population. If I do not impose this focus, then it is possible that any observed change in the effect of children on maternal employment over time is not due to a change in the employment response of mothers, but due to changes in household composition.

I report summary statistics for the sample in Table E.1 and in Table E.2. Married women represent a

¹³ I differ from ANGRIST and EVANS (1998) in that I do not exclude women whose first child was born before marriage, and I do not exclude women whose reported number of children ever born differs from their reported number of children. The ACS does not contain the information that is required for these rules.

particularly high fertility group in the population, and have exhibited a slight upward trend in fertility since 1980 that has partially offset declining fertility amongst unmarried women. The most striking statistic as evidence of this overall decline in fertility is a fall in the number of children that all women aged between 21 and 35 years report having, from 1.153 children in 1980 to 0.844 in 2015. Set against this decline in fertility is an increase in the female employment rate, both amongst married and unmarried women. The principal objective of this chapter is to discover whether this increase in female employment is (in part) due to the documented decline in fertility, or is despite it; and in the context of maternal employment, whether the observed increase in the employment rate of mothers has been dampened by an increase in family size.

I remark upon two further trends that are evident in Figure E.2 . First, that a husband’s labour supply remains stable across the study period, both at the extensive margin and at the intensive margin. This trend is well-documented elsewhere, and suggests that male labour supply is inelastic. If children do influence the employment decision then it is likely that this effect accrues on the employment decision of the mother. Second, (deflated) household income increases from (1999) \$50,900 in 1980 to (1999) \$55,034 in 2015 amongst married couples. This change in household income is of the order of ten percentage points, and is much more subdued than historic growth rates would imply. Indeed, decomposing household income into wife’s labour income and husband’s labour income, it is apparent that this increase is driven purely by an increase in wife’s labour income—real wage growth is otherwise stagnant. It is reasonable to suggest that this increase in wife’s labour income is due to the documented increase in maternal employment.

—| SUBSECTION 4.A |—

VARIABLE DEFINITION

In presenting the Balke-Pearl model in the context of maternal employment, I have stated particular definitions for each of the variables. I note that I subject some of these definitions to change—I propose various definitions of covariates and of the instrument. I emphasise that I continue to define the outcome as a dichotomous indicator for whether a mother worked for pay in the previous year; and I define treatment as a dichotomous indicator for whether a mother has two children, or three or more children.

I continue to define age, age at first birth, and ethnicity and race as covariates. However, in some instances I restrict the vector of covariates that I include to be a proper subset of these covariates. Including the full set of covariates is not feasible in all cases, and leads to degeneracy of the estimator. To guarantee non-degeneracy of the estimator whilst retaining as much detail of covariates as possible,

FIGURE E.1
Fertility and labour supply measures

	Means (and standard deviations)					
	PUMS 1980		ACS 2010		ACS 2015	
Women aged 21-35						
Children ever born	1.238	(1.295)	-	-	-	-
Reported number of children	1.153	(1.236)	0.929	(1.216)	0.844	(1.186)
Has two or more children	0.374	(0.484)	0.292	(0.455)	0.263	(0.440)
Worked in previous year	0.732	(0.443)	0.773	(0.419)	0.787	(0.409)
Observations	1,380,870		263,009		271,322	
Women aged 36-50						
Children ever born	2.818	(1.921)	-	-	-	-
Reported number of children	1.725	(1.454)	1.344	(1.247)	1.386	(1.258)
Has two or more children	0.531	(0.499)	0.439	(0.496)	0.455	(0.498)
Worked in previous year	0.665	(0.472)	0.768	(0.422)	0.772	(0.419)
Observations	885,405		313,578		285,436	
Women aged 21-35 with two or more children and no adopted children						
Children ever born	2.588	(0.905)	-	-	-	-
Reported number of children	2.500	(0.797)	2.532	(0.815)	2.538	(0.828)
Has three or more children	0.359	(0.480)	0.380	(0.485)	0.380	(0.485)
Worked in previous year	0.560	(0.496)	0.668	(0.471)	0.679	(0.467)
Observations	486,736		67,073		60,216	
Married women aged 21-35 with two or more children and no adopted children						
Children ever born	2.530	(0.842)	-	-	-	-
Reported number of children	2.476	(0.770)	2.494	(0.785)	2.506	(0.811)
Has three or more children	0.349	(0.477)	0.361	(0.480)	0.361	(0.480)
Worked in previous year	0.532	(0.499)	0.631	(0.482)	0.628	(0.483)
Observations	333,334		38,958		35,036	

The table reports means of variables (that are equivalent to probabilities in some cases). Standard deviations are reported in parentheses. Children ever born is the number of biological children that a mother reports ever having, and is available only in the PUMS. Reported number of children is the number of children that a mother reports having in the household; Has two or more children and Has three or more children are defined in terms of Reported number of children. Worked in previous year indicates whether a mother worked for profit, pay, or as an unpaid family worker in the previous year. All samples exclude mothers whose age or gender have been imputed by the Census Bureau or have children whose age or gender have been imputed. The married sample consists of women who are married and have been married only once and are married to spouses that have been married only once.

I discretise both age and age at first birth, and partition ethnicity and race into several broad classes. I partition age and age at first birth into five sets: 15 to 18 years, 18 to 22 years, 22 to 25 years, 25 to 30 years, and 30 to 35 years. I design these partitions—specifically the first three partitions—to reflect the current U.S. education system, noting that individuals who are currently undertaking an education

FIGURE E.2
Fertility measures, child gender, and labour market outcomes

	Means (and standard deviations) for women aged 21-35 with two or more children								
	PUMS 1980			ACS 2010			ACS 2015		
	All women	Married couples		All women	Married couples		All women	Married couples	
	Wives	Husbands		Wives	Husbands		Wives	Husbands	
Children ever born	3.561 (0.919)	3.553 (0.887)	-	-	-	-	-	-	
Reported number of children	2.550 (0.811)	2.529 (0.795)	2.661 (0.864)	2.640 (0.847)	2.684 (0.887)	2.661 (0.876)	2.661 (0.876)	2.661 (0.876)	
Has more than two children	0.399 (0.490)	0.388 (0.487)	0.472 (0.499)	0.463 (0.499)	0.478 (0.500)	0.466 (0.499)	0.466 (0.499)	0.466 (0.499)	
Male first-born child	0.514 (0.500)	0.514 (0.500)	0.512 (0.500)	0.513 (0.500)	0.506 (0.500)	0.505 (0.500)	0.505 (0.500)	0.505 (0.500)	
Male second-born child	0.512 (0.500)	0.511 (0.500)	0.509 (0.500)	0.507 (0.500)	0.512 (0.500)	0.511 (0.500)	0.511 (0.500)	0.511 (0.500)	
First two children male	0.265 (0.442)	0.265 (0.441)	0.263 (0.440)	0.263 (0.440)	0.260 (0.439)	0.259 (0.438)	0.259 (0.438)	0.259 (0.438)	
First two children female	0.240 (0.427)	0.239 (0.427)	0.242 (0.428)	0.242 (0.429)	0.242 (0.428)	0.242 (0.428)	0.242 (0.428)	0.242 (0.428)	
First two children share gender	0.505 (0.500)	0.504 (0.500)	0.505 (0.500)	0.506 (0.500)	0.502 (0.500)	0.501 (0.500)	0.501 (0.500)	0.501 (0.500)	
Plural second birth	0.009 (0.093)	0.008 (0.091)	0.013 (0.112)	0.012 (0.109)	0.013 (0.112)	0.012 (0.110)	0.012 (0.110)	0.012 (0.110)	
Age	30.215 (3.497)	30.175 (3.521)	31.246 (3.143)	31.158 (3.187)	31.500 (3.006)	31.438 (3.039)	31.438 (3.039)	31.438 (3.039)	
Age at first birth	21.226 (3.028)	21.414 (2.991)	22.101 (3.726)	22.348 (3.729)	22.415 (3.729)	22.678 (3.714)	22.678 (3.714)	22.678 (3.714)	
Observations	376,161	297,124	41,089	30,258	36,196	26,805	26,805	26,805	

Continued ↔

FIGURE E.2
Fertility measures, child gender, and labour market outcomes

	Means (and standard deviations) for women aged 21-35 with two or more children									
	PUMS 1980		ACS 2010		ACS 2015					
	All women	Married couples	All women	Married couples	All women	Married couples				
	Wives	Husbands	Wives	Husbands	Wives	Husbands				
Worked for pay	0.545 (0.498)	0.535 (0.499)	0.974 (0.158)	0.634 (0.482)	0.628 (0.483)	0.945 (0.228)	0.626 (0.484)	0.619 (0.486)	0.958 (0.201)	
Weeks worked	17.8↔21.6	17.3↔21.0	45.5↔48.5	26.7↔29.5	26.3↔29.1	43.3↔46.4	26.7↔29.4	26.1↔28.8	45.5↔48.1	
Hours per week	13.3↔15.8	12.9↔15.5	39.2↔43.7	-	-	-	-	-	-	
Labour income	7,022 (11,138)	6,677 (10,846)	37,981 (25,498)	12,789 (18,036)	12,508 (17,858)	33,378 (34,418)	13,180 (19,914)	13,095 (20,077)	35,542 (36,735)	
Family income	50,903 (28,111)	50,900 (27,910)	-	52,723 (43,302)	52,730 (43,495)	-	54,803 (46,623)	55,034 (46,676)	-	
Non-wife income	43,881 (26,633)	44,222 (26,628)	-	39,934 (38,182)	40,222 (38,554)	-	41,623 (40,952)	41,939 (40,869)	-	
Observations	376,161	297,124	41,089	30,258	36,196	26,805				

The sample includes mothers aged between 21 and 35 years with at least two children, whose age or gender (or that of any children) has not been imputed by the Census Bureau. Married couples include only those couples that are married at the date of survey and have only been married to each other. Labour market outcome variables are reported for the previous year. Weeks worked and Hours per week are set-valued and so the reported means are set-valued (an upper bound of 84 hours is used for Hours worked). Incomes are adjusted to 1999 USD values using Census Bureau adjustment factors. Non-wife income is defined as Family income less the value of Labour income for a wife.

programme may postpone having children until after they have completed that programme. I define ethnicity as a dichotomous indicator for whether a mother is of Hispanic origin. I define race as a three-valued indicator for whether a mother is black, white, or other race.

I propose four definitions of the instrument. First, I define the instrument as a dichotomous indicator for a plural second birth (a mother gives birth to twins or higher plurality siblings). Second, I define the instrument as a dichotomous indicator for whether a mother's first two children share the same sex. Third, I propose separating the sex composition of a mother's first two children into its constituent pairs. That is, I define the instrument as a four-valued variable that takes different values depending upon the sex of the first child and of the second child. Fourth, I combine a plural second birth with the sex composition of a mother's first two children. I note that the influence on the fertility decision is different for a plural second birth and sex; in particular, a plural second birth acts mechanically, in that a mother that experiences a plural second birth has three or more children by definition. In contrast, sex influences the fertility decision through parental preferences. If a mother experiences a plural second birth though, parental preferences are irrelevant to the fertility decision. To reflect this asymmetry, I define the instrument as a five-valued variable with a hierarchical structure that assigns priority to a plural second birth; if a mother does not experience a plural second birth then the instrument takes either the second through fifth value, depending upon the sex composition of her first two children.

The generalised Balke-Pearl model imposes an exclusion restriction on the assumed data generating process. Such a restriction is violated if twins incur additional costs versus non-plural siblings; for example, if having children of the same age requires duplication, versus handing-down or reusing. These additional costs also include non-pecuniary costs that are especially high during early childhood, and that may extend beyond a year (recall that I exclude mothers whose second child is less than one year old) in the case of twins. Due to these early childhood considerations, I refine the sample to exclude mothers whose second child is aged under five years whenever I define the instrument as a plural second birth (either by itself or with sex composition). A plural second birth is a rare event, occurring in around 1-in-100 second births in the data that I study. Degeneracy of the estimator is of first-order concern whenever I define the instrument as a plural second birth, and is exacerbated by the additional sample refinement that I perform in this case.

FIGURE E.3
Parental choice over family size given sex of children

	Probabilities (and standard errors) of having another child					
	All women			Married women		
	PUMS 1980	ACS 2010	ACS 2015	PUMS 1980	ACS 2010	ACS 2015
Sex of first child in families with one or more children						
(1.) ♂	0.714 (0.001)	0.743 (0.003)	0.742 (0.003)	0.726 (0.001)	0.752 (0.003)	0.745 (0.003)
(2.) ♀	0.713 (0.000)	0.749 (0.003)	0.744 (0.004)	0.726 (0.001)	0.756 (0.003)	0.749 (0.004)
↔(1)-(2)	0.001 (0.001)	-0.006 (0.004)	-0.002 (0.005)	0.000 (0.001)	-0.004 (0.005)	-0.004 (0.005)
Observations	564,445	60,576	54,146	456,066	50,839	45,983
Sex of first two children in families with two or more children						
(3.) ♂♀ or ♀♂	0.360 (0.001)	0.419 (0.005)	0.423 (0.005)	0.352 (0.001)	0.421 (0.005)	0.419 (0.006)
♂♂	0.413 (0.002)	0.471 (0.006)	0.498 (0.007)	0.408 (0.002)	0.467 (0.007)	0.490 (0.008)
♀♀	0.432 (0.002)	0.488 (0.006)	0.490 (0.008)	0.430 (0.002)	0.491 (0.007)	0.494 (0.008)
(4.) ♂♂ or ♀♀	0.422 (0.001)	0.479 (0.004)	0.494 (0.005)	0.419 (0.001)	0.479 (0.005)	0.492 (0.005)
↔(3)-(4)	-0.062 (0.002)	-0.060 (0.006)	-0.071 (0.007)	-0.056 (0.002)	-0.058 (0.007)	-0.073 (0.007)
Observations	360,919	35,892	31,633	297,124	30,258	26,805

The table displays the probability of having an additional child given the sex of existing children in the household; specifically, the probability of having a second child given the sex of the first-born child, and the probability of having a third child given the sex of the first- and second-born child. The sample includes mothers aged between 21 and 35 years with at least one child, whose age or sex (or that of any children) has not been imputed by the Census Bureau. Married couples include only those couples that are married at the date of survey and have only been married to each other.

—| SECTION 5 |—

LOCAL ESTIMATION

Existing empirical evidence as to the effect of children on maternal employment is mixed. ANGRIST and EVANS (1998) finds that additional children reduce the probability that a mother is employed. IACOVOU (2001) and AL-KHAJA (2016) conduct comparable analyses in the context of the British and Egyptian labour markets respectively, instead finding that additional children increase the probability that a mother is employed. Common to all three papers is the parameter that is recovered; specifically, ANGRIST and EVANS (1998), IACOVOU (2001) and AL-KHAJA (2016) recover a Local Average Treatment Effect (LATE; IMBENS and ANGRIST, 1994) that is the effect of children on maternal employment

FIGURE E.4
Differences in means for demographic variables by instrumental variable definition

	Differences in means (and standard errors)					
	PUMS 1980		ACS 2010		ACS 2015	
	Same-sex	Plural birth	Same-sex	Plural birth	Same-sex	Plural birth
Age	-0.021 (0.015)	0.259 (0.071)	0.075 (0.047)	0.632 (0.198)	-0.006 (0.045)	-0.445 (0.334)
Age at first birth	0.009 (0.011)	0.297 (0.067)	0.110 (0.051)	0.291 (0.270)	0.063 (0.050)	-0.770 (0.336)
Black	0.000 (0.001)	0.019 (0.008)	0.003 (0.004)	0.014 (0.029)	0.000 (0.004)	-0.020 (0.021)
White	0.000 (0.001)	-0.019 (0.009)	-0.007 (0.006)	0.089 (0.038)	0.003 (0.007)	-0.017 (0.050)
Other race	-0.000 (0.001)	-0.000 (0.004)	0.004 (0.005)	-0.102 (0.029)	-0.102 (0.006)	0.037 (0.047)
Hispanic	-0.001 (0.001)	-0.011 (0.007)	-0.009 (0.006)	-0.115 (0.045)	-0.008 (0.008)	-0.089 (0.044)
Observations	297,124	171,515	30,258	13,784	26,805	12,232

The sample includes mothers aged between 21 and 35 years with at least two children, whose age or gender (or that of any children) has not been imputed by the Census Bureau. Married couples include only those couples that are married at the date of survey and have only been married to each other. The table displays the difference in the mean of each variable between each the instrument, for each definition of the instrument. The columns titled Same-sex state the difference in mean when the first two children share the same gender versus when they do not. The columns titled Plural birth state the difference in mean when a mother experiences a plural second birth versus when she does not. The variables in the left-hand column are demographic and ethnic and racial indicators.

amongst an unidentifiable subset of the sample population. In Figure E.5 I report comparable estimates for the census and ACS data that I study.¹⁴ Like ANGRIST and EVANS (1998), I find that the local effect of additional children in 1980 is negative and statistically significant. Estimates for 2010 and 2015, however, vary. In particular, I find that the local effect of additional children in 2010 is positive and statistically significant. Furthermore, the absolute magnitudes of these effects vary considerably, from four percentage points up to 25 percentage points.

Although the local effects that I report are empirically interesting, I include these estimates for two reasons. First, to emphasise that local effects may provide mixed evidence. I do not intend this as a criticism. Instead, I view this ambiguity as a natural implication of allowing for heterogeneity of response to treatment. I argue that it is interesting to think about how these local effects fit with evidence about an effect in the wider population. For example, does the sub-population that is affected by a policy represent a group that is particularly responsive to treatment? Second, to provide a numerical benchmark that may be useful for interpreting subsequent estimates.

¹⁴ I note that the Two-stage least squares estimates that I report relies on a linear model. I specify this model to include age and age at first birth (not discretised; in years), and a full set of interactions between ethnicity (Hispanic or non-Hispanic) and race (black, white, and other race).

FIGURE E.5
Estimates of local treatment effects

	Estimates and predictions (with standard errors)					
	PUMS 1980		ACS 2010		ACS 2015	
	Same-sex	Plural birth	Same-sex	Plural birth	Same-sex	Plural birth
Wald estimator						
Intercept	0.585 (0.011)	0.612 (0.014)	0.624 (0.053)	0.480 (0.069)	0.638 (0.048)	0.680 (0.150)
LATE	-0.131 (0.027)	-0.067 (0.027)	0.008 (0.118)	0.257 (0.111)	-0.043 (0.104)	-0.104 (0.094)
Two-stage least squares estimator						
Intercept	0.581 (0.010)	0.613 (0.014)	0.633 (0.050)	0.515 (0.065)	0.639 (0.045)	0.663 (0.100)
LATE	-0.121 (0.026)	-0.070 (0.027)	-0.015 (0.111)	0.199 (0.104)	-0.047 (0.099)	-0.080 (0.159)
Observations	297,124	171,515	30,258	13,784	26,805	12,232

The sample includes mothers aged between 21 and 35 years with at least two children, whose age or gender (or that of any children) has not been imputed by the Census Bureau. Married couples include only those couples that are married at the date of survey and have only been married to each other. Estimates of local treatment effects using the Wald estimator and the Two-stage least squares estimator are stated. The Two-stage least squares estimator includes as covariates all variables in Table E.4 plus interaction terms between the racial and ethnic indicators. Intercept is reported at the mean of the covariates for the Two-stage least squares estimator.

—| SECTION 6 |—

ESTIMATION

I focus on the ATE and its primitives, the so-called potential outcome distributions (D. B. RUBIN, 1974), that may inform both *ex-ante* policy design and *ex-post* policy evaluation. I report estimates of these objects in Figure E.6 and in Figure E.7. Each estimate is a non-trivial set, reflecting the fact that the Balke-Pearl model partially identifies the ATE and its primitives. Each estimate is obtained as the sample analogue of a set of bounds that are implied by the model, and are calculated using a half-median unbiased estimator that is described in CHERNOZHUKOV et al. (2013). Focusing on the ATE, the interpretation of an estimate is then as a set that contains the sample analogue of the ATE with probability one; the sample analogue of the ATE may then take any value in this set.

I separate the analysis into two parts, by estimating the ATE and its primitives with and without covariates. I note that the model without covariates implies the generalised Balke-Pearl bounds when the sub-population x is taken to be the population. By separating the analysis into these two parts, I

am able to study the influence that covariates have on estimates.¹⁵ If covariates and the instrument are not statistically independent then it is (theoretically) necessary to include covariates. Such a statement, however, does not preclude the possibility that there is no empirical benefit to including covariates.

I make the following observations. First, that estimates are broadly uninformative. In particular, I fail to recover the sign of the effect of additional children on maternal employment. Effects of 20 percentage points or more cannot be ruled-out in many cases. Second, I cannot reject the hypothesis that the effect of additional children on maternal employment is unchanged over the study period, although I note that I also am unable to reject the competing hypotheses that the effect is increasing or is decreasing. Third, that there is little difference between the estimates that are presented in Figure E.6 versus those that are presented in Figure E.7. One comparison where this observation does not hold is when I define the instrument as a plural second birth and allow age to influence the employment decision and the fertility decision. In this instance, there is considerable difference between the estimates that are delivered under each model. As can be seen in Figure E.4, there is strong correlation between the age of a mother and the probability of a plural second birth. I identified such a case as (theoretically) necessitating the inclusion of covariates, and I find that this is supported empirically. I note that the confidence regions that I report suffer from projection conservatism, and it is possible that this observation is driven simply by the inclusion of an exogenous characteristic with (relatively) many points of support. Fourth, that estimates that use the definition of the instrument as a plural second birth are more informative than those that use the definition of the instrument as the sex composition of a mother's first two children. I note that a plural second birth is a monotonic variable; plural second birth acts mechanically, in that a mother that experiences a plural second birth has three or more children by definition. Put differently, a mother that experiences a plural second birth cannot have two children. It follows that this definition of the instrument precludes some response types and has strong identifying power. Fifth, that the local effects that I report in Figure E.5 do not appear to be extreme in relation to the estimates of the ATE that I report.

—| APPENDIX E.1 |—

THE SUITABILITY OF AN AUXILIARY MODEL

An interesting question that naturally arises is whether it is necessary to include covariates. More precisely, does omitting relevant covariates constitute a form of model misspecification? Or, is an analysis *à la* BALKE and PEARL (1997) sufficient? The short answer is that it depends. If covariates and the instrument are statistically independent then it is not necessary to include covariates, although the

¹⁵ I address the question of whether it is necessary to include covariates in an appendix.

FIGURE E.6
Estimates of global treatment effects without covariates

	Half-median unbiased estimates (with 95% confidence sets)							
	PUMS 1980		ACS 2010		ACS 2015			
	$\bar{Y}_{0,n}$	ATE _n	$\bar{Y}_{0,n}$	ATE _n	$\bar{Y}_{0,n}$	ATE _n	$\bar{Y}_{0,n}$	ATE _n
Same-sex	0.271↔0.629 (0.269↔0.631)	-0.405↔0.535 (-0.410↔0.540)	0.099↔0.539 (0.094↔0.549)	-0.407↔0.573 (-0.423↔0.588)	0.113↔0.553 (0.107↔0.564)	-0.401↔0.570 (-0.419↔0.588)		
Gender decomposition	0.271↔0.629 (0.270↔0.631)	-0.403↔0.528 (-0.406↔0.531)	0.098↔0.534 (0.093↔0.542)	-0.403↔0.528 (-0.415↔0.578)	0.113↔0.554 (0.107↔0.562)	-0.403↔0.569 (-0.417↔0.583)		
Observations		297,124		30,258		26,805		
Plural second birth	0.175↔0.679 (0.164↔0.681)	-0.252↔0.291 (-0.279↔0.331)	0.064↔0.692 (0.056↔0.702)	-0.513↔0.180 (-0.530↔0.246)	0.089↔0.714 (0.070↔0.725)	-0.506↔0.261 (-0.525↔0.373)		
Plural second birth and gender decomposition	0.179↔0.653 (0.177↔0.656)	-0.224↔0.296 (-0.240↔0.314)	0.066↔0.647 (0.059↔0.657)	-0.467↔0.186 (-0.486↔0.237)	0.090↔0.681 (0.081↔0.692)	-0.466↔0.276 (-0.487↔0.348)		
Observations		171,515		13,784		12,232		

The sample includes mothers aged between 21 and 35 years with at least two children, whose age or gender (or that of any children) has not been imputed by the Census Bureau. Married couples include only those couples that are married at the date of survey and have only been married to each other. Half-median unbiased estimates and 95% confidence regions are reported; these estimates are calculated using a method outlined in CHENNOZHUKOV et al. (2013). Estimates are presented for various definitions of the instrument. Same-sex is defined as the event that the first two children share the same gender; gender decomposition decomposes gender further into the gender of the first child and the gender of the second child; plural second birth is defined as the event that a mother experiences a plural second birth; and Plural second birth and gender decomposition combines the two definitions of the instrument, giving priority to a plural second birth.

FIGURE E.7
Estimates of global treatment effects with covariates

	Half-median unbiased estimates and p-values (confidence sets)								
	PUMS 1980		ACS 2010		ACS 2015				
	$\hat{Y}_{0:n}$	ATE _n	χ^2	$\hat{Y}_{0:n}$	ATE _n	χ^2	$\hat{Y}_{0:n}$	ATE _n	χ^2
Samesex									
Hispanic	0.271↔0.628 (0.270↔0.630)	-0.407↔0.534 (-0.410↔0.537)	0.176	0.100↔0.538 (0.096↔0.544)	-0.406↔0.571 (-0.416↔0.581)	0.000	0.114↔0.552 (0.109↔0.559)	-0.399↔0.567 (-0.412↔0.579)	0.000
Hispanic, age, age at first birth	0.263↔0.638 (0.260↔0.638)	-0.426↔0.550 (-0.432↔0.555)	0.018	0.073↔0.542 (0.069↔0.584)	-0.433↔0.619 (-0.482↔0.646)	0.000	0.088↔0.590 (0.078↔0.605)	-0.465↔0.629 (-0.490↔0.653)	0.000
Observations		297,124			30,258			26,805	
Plural second birth									
Hispanic	0.166↔0.678 (0.164↔0.680)	-0.246↔0.305 (-0.261↔0.325)	0.000	0.059↔0.689 (0.056↔0.696)	-0.511↔0.216 (-0.522↔0.258)	0.000	0.073↔0.712 (0.068↔0.719)	-0.503↔0.311 (-0.516↔0.370)	0.671
Age	0.165↔0.679 (0.163↔0.681)	-0.260↔0.328 (-0.279↔0.347)	0.000	0.060↔0.691 (0.056↔0.698)	-0.514↔0.222 (-0.527↔0.279)	0.000		Degenerate	0.000
Observations		171,515			13,784			12,232	

The sample includes mothers aged between 21 and 35 years with at least two children, whose age or gender (or that of any children) has not been imputed by the Census Bureau. Married couples include only those couples that are married at the date of survey and have only been married to each other. Half-median unbiased estimates and 95% confidence regions are reported; these estimates are calculated using a method outlined in CHERNOZHUKOV et al. (2013). Estimates are presented for various definitions of the instrument. For each definition of the instrument, a vector of covariates is specified. Hispanic is a dichotomous indicator that a mother is Hispanic; Hispanic, age and age at first birth also includes the age of a mother and the age of her first birth (both in years). P-values of Pearson's χ -square test (that the instrument is random across exogenous covariates) are also reported.

bounds that are stated in BALKE and PEARL (1997) are not sharp and instead define an outer region. If many covariates are observed—so that there is a high-dimension vector of covariates—then there may be good reason to prefer to ignore covariates with the loss of some precision. If covariates and the instrument are not statistically independent then it is (theoretically) necessary to include covariates.

To develop a (hopefully) satisfactory explanation, I posit the following simple example. Imagine that h is a non-trivial function of x (i.e., h varies in x) so that covariates are in some sense relevant. I compare what can be learnt from the generalised Balke-Pearl model—Assumptions E.1 to E.3—with what can be learnt from Assumption E.2 combined with the following two assumptions.

Assumption E.4 (Simple exclusion). *The data generating process satisfies*

$$\begin{aligned} s_y &= g_y(t, u) \\ s_t &= g_t(z, u) \end{aligned} \tag{E.10}$$

such that the generalised Balke-Pearl model is triangular in Y , T and Z .

Assumption E.5 (Simple statistical independence). *The data generating process satisfies*

$$U | \{Z\} \sim \text{Uniform}(0, 1) \tag{E.11}$$

for which g is appropriately normalised.

I note that Assumptions E.2, E.4 and E.5 are equivalent to the Balke-Pearl model. I intend Assumptions E.2, E.4 and E.5 to constitute an auxiliary model for Assumptions E.1 to E.3. To fix ideas, suppose that Assumptions E.1 to E.3 identify the $\text{ATE}(x)$ and the $\text{ATE}(\tilde{x})$, for some $a, b \in \mathfrak{Z} \mapsto (x, \tilde{x})$, up to the half-lines

$$\left\{ \begin{array}{l} q(a) \\ q(b) \end{array} \right\} \leq \text{ATE}(x) \quad \text{and} \quad \left\{ \begin{array}{l} \tilde{q}(a) \\ \tilde{q}(b) \end{array} \right\} \leq \text{ATE}(\tilde{x}) \tag{E.12}$$

Of course, these bounds are overly simplistic and I do not claim that they reflect all of the identifying information that Assumptions E.1 to E.3 contains—they are a tool for exposition and serve only to convey a point. Given these bounds though, it follows that Assumptions E.1 to E.3 identify the ATE

up to the half-line

$$\left\{ \begin{array}{l} q(a) \cdot \Pr(X = x) + \tilde{q}(a) \cdot \Pr(X = \tilde{x}) \\ q(a) \cdot \Pr(X = x) + \tilde{q}(b) \cdot \Pr(X = \tilde{x}) \\ q(b) \cdot \Pr(X = x) + \tilde{q}(a) \cdot \Pr(X = \tilde{x}) \\ q(b) \cdot \Pr(X = x) + \tilde{q}(b) \cdot \Pr(X = \tilde{x}) \end{array} \right\} \leq \text{ATE} \quad \text{E.13}$$

In contrast, Assumptions E.2, E.4 and E.5 do not identify either the $\text{ATE}(x)$ or the $\text{ATE}(\tilde{x})$ but do identify the ATE—albeit incorrectly—up to the half-line

$$\left\{ \begin{array}{l} q(a) \cdot \Pr(X = x|Z = a) + \tilde{q}(a) \cdot \Pr(X = \tilde{x}|Z = a) \\ q(b) \cdot \Pr(X = x|Z = b) + \tilde{q}(b) \cdot \Pr(X = \tilde{x}|Z = b) \end{array} \right\} \leq \text{ATE} \quad \text{E.14}$$

There are two differences between the half-lines defined by Equations E.13 and E.14. First, the number of unique constraints is different. Second, each constraint assigns (possibly) different weights to each lower bound. If the conditional and unconditional probabilities are equal then it is immediately apparent that the first and fourth constraints of Equation E.13 coincide with the two constraints of Equation E.14. I note that this condition is equivalent to strict randomisation (of Z , independently of X). An obvious consequence is that the auxiliary model identifies an outer region of the identified set. If the conditional and unconditional probabilities are not, however, equal then it is not possible to determine how the two half-lines relate to each other. It is in this specific sense that the auxiliary model is misspecified—the auxiliary model does not necessarily identify an outer region. I suggest that it is sensible to test for randomisation—something that can be done—before deciding whether to use an auxiliary model.

—| APPENDIX E.2 |—

DISCUSSION OF THE GENERALISED BALKE-PEARL BOUNDS

Bounding the Average Structural Functions in the context of the generalised Balke-Pearl model is equivalent to constructing a cover in u -space and mapping this to observable probabilities—an insight that CHESHER and ROSEN (2017) makes and exploits in a more general framework. Even in simple settings though, there are many possible ways to partition a set of latent types. Fortunately, CHESHER and ROSEN (2017) describes an algorithm that reduces the number of possible sets in u -space that must be considered to a smaller class that GALICHON and HENRY (2011) names the core-determining sets. Importantly, any undominated cover in u -space can be expressed in terms of the core-determining sets

FIGURE E.8
Latent types and the identification problem

$Y_1(x)$	$Y_0(x)$	$T_1(x)$	$T_0(x)$	Observable probabilities	
1	1	1	1	$Q_{11.1}(x)$	$Q_{11.0}(x)$
1	1	1	0	$Q_{11.1}(x)$	$Q_{10.0}(x)$
1	1	0	1	$Q_{10.1}(x)$	$Q_{11.0}(x)$
1	1	0	0	$Q_{10.1}(x)$	$Q_{10.0}(x)$
1	0	1	1	$Q_{11.1}(x)$	$Q_{11.0}(x)$
1	0	1	0	$Q_{11.1}(x)$	$Q_{00.0}(x)$
1	0	0	1	$Q_{00.1}(x)$	$Q_{11.0}(x)$
1	0	0	0	$Q_{00.1}(x)$	$Q_{00.0}(x)$
0	1	1	1	$Q_{01.1}(x)$	$Q_{01.0}(x)$
0	1	1	0	$Q_{01.1}(x)$	$Q_{10.0}(x)$
0	1	0	1	$Q_{10.1}(x)$	$Q_{01.0}(x)$
0	1	0	0	$Q_{10.1}(x)$	$Q_{10.0}(x)$
0	0	1	1	$Q_{01.1}(x)$	$Q_{01.0}(x)$
0	0	1	0	$Q_{01.1}(x)$	$Q_{00.0}(x)$
0	0	0	1	$Q_{00.1}(x)$	$Q_{01.0}(x)$
0	0	0	0	$Q_{00.1}(x)$	$Q_{00.0}(x)$

The table associates latent types—as defined in Equation E.4—with observable probabilities in a wholly dichotomous framework in which covariates can take only one value.

and their intersection with those latent types that underpin the structural characteristic of interest.^{16,17} Moreover, the core-determining sets can be stratified according to covariates, which is especially useful in practice given that the support of covariates is responsible for a curse of dimensionality on the number of latent types. The question that I ask here—and that I attempt to provide an intuitive answer to—is what motivates and rationalises this stratification?

As a preliminary step, I focus on the case in which covariates can take only one value and—for illustration—study a wholly dichotomous framework, which is sufficient to introduce a result that I subsequently rely upon. As is stated in BALKE and PEARL (1997) and its preceding works, this framework is associated with a total of 16 latent types, each defined according to the potential outcomes that they induce for different values of the observable exogenous random variates (in this case, simply different values of the instrument). These latent types are listed in Figure E.8 alongside the observable probabilities that they can map to, as determined by the capacity functional that is defined in MOLCHANOV (2005). What is immediately apparent is that several latent types are associated with the same observable probabilities and it is these collections of latent types—that are disjoint for a single value of the instrument, and overlapping for different values—that form the basis of the core-determining sets. A consequence is that any cover of all 16 latent types must map to observable probabilities that

¹⁶ By undominated I mean any cover that maps to a sharp bound on the structural characteristic of interest.

¹⁷ Here I emphasise that the core-determining sets themselves are not necessarily proper subsets of those latent types that underpin the structural characteristic of interest. Instead, they are unions of latent types, only some of which need be latent types that underpin the structural characteristic of interest.

sum to one or more, with violation of this a condition for falsification of the Generalised Balke-Pearl model as is explored in KÉDAGNI and MOURIFIÉ (2020).

Supposing instead that covariates can take two values, there are a total of 256 latent types even in this simple setting, which serves to illustrate the aforementioned curse. Rather than extend Figure E.8, I instead ask the reader to imagine a matrix comprising 16 rows and 16 columns—with each row corresponding to specific behaviour given one value of the covariates, and each column corresponding to specific behaviour given the other. With this analogy in mind, the problem of bounding one of the Average Structural Functions amounts to constructing a cover over some of the rows or some of the columns. Since each row intersects every column though, the previously stated consequence implies that using columns to cover a row—or *vice versa*—necessarily maps to a trivial unit bound. Such a cover is clearly dominated by using groups of rows to cover rows and groups of columns to cover columns—what amounts to stratification. It is for this reason that each value of the covariates can be considered in isolation.

I remark that this argument does not specifically rely upon the dimensionality of the framework—whether that is with respect to the support of the instrument or the support of the covariates. In particular, the matrix can be expanded in any direction to accommodate the increased number of behaviours that are created by additional points of support of the instrument conditional on a given value of the covariates; the analysis of BERESTEANU et al. (2012) can easily be extended to derive Equations E.8 and E.9 in any direction in which the instrument has more than two points of support; and the dimension of the matrix can be increased to accommodate additional points of support of the covariates, with rows and columns transforming into planes and higher-dimensional equivalents.

—| APPENDIX E.3 |—

DERIVING THE GENERALISED BALKE-PEARL BOUNDS*

This appendix formalises the ideas that are advanced in the preceding appendix, presenting a formal proof of the statement that Equations E.8 and E.9 constitute sharp bounds on the Average Structural Functions.

The proof itself adapts other works. In particular, I lean heavily on results in BERESTEANU et al. (2012), CHESHER and ROSEN (2017), and RICHARDSON and ROBINS (2014), following the approaches taken in KÉDAGNI and MOURIFIÉ (2020) and RICHARDSON and ROBINS (2024) to proving sharpness. To facilitate this, I exploit the mapping that exists between the structural equation and potential outcome notations

* I thank my examiners—Professors Dennis Kristensen and Karim Chalak—for their helpful suggestions during my oral examination. This appendix is designed to commit their required amendments.

that is discussed in Section 2, using the notation that is more convenient for whichever result I wish to show. I mirror RICHARDSON and ROBINS (2024) in separating the proof of sharpness of Equations E.8 and E.9 into two parts—establishing that

$$P(Y_1(x) = y_1, Y_0(x) = y_0) \leq Q_{y_1 1.z}(x) + Q_{y_0 0.z}(x) \tag{E.15}$$

$$P(Y_t(x) = 1 - y) \leq 1 - Q_{yt.z}(x)$$

characterises $P(\{Y_1(x), Y_0(x) : x \in \mathfrak{X}\})$, which is RICHARDSON and ROBINS (2024, §Theorem 1); and then showing that Equations E.8 and E.9 are the vertices of the linear programme that is implied by Equation E.15, which is RICHARDSON and ROBINS (2024, §Theorem 2). The first part necessitates (i.) showing that the set of structures that are admitted by the generalised Balke-Pearl model and that are compatible with the observable distribution is contained within the set of structures that are implied by Equation E.15, and (ii.) proposing a structure that is compatible with Equation E.15 and that is admitted by the generalised Balke-Pearl model. To do (i.), I apply results from random set theory, as in BERESTEANU et al. (2012) and CHESHER and ROSEN (2017); to do (ii.), I propose a minor adaptation of the marginal distribution that is proposed in KÉDAGNI and MOURIFIÉ (2020)—a distribution that exhibits many of the properties that I require—and invoke RICHARDSON and ROBINS (2024, §Lemma 4) to establish the existence of a joint distribution with the required properties. The second part is somewhat redundant, as RICHARDSON and ROBINS (2024, §Theorem 2) can always be invoked upon noting that any cover of the Average Structural Functions only uses observable quantities with the same value of covariates—an implication of the idea that is pursued in the preceding appendix that only by using all columns can a row or rows of a matrix be covered.

In what follows, I adopt the following convention. I signify the power set over a collection by including a lemniscate in the superscript position (as I do for more general spaces). For example, \mathfrak{U}^∞ denotes the power set over \mathfrak{U} . Where the collection is the image of a correspondence, I position the lemniscate before the parentheses enclosing the arguments.

To facilitate this identification analysis, I sequentially define

$$\text{Contour}(y, t, z, x) \doteq \{u : y = 1(h_y(t, x, u) \geq 0) \text{ and } t = 1(h_t(z, x, u) \geq 0)\} \tag{E.16}$$

$$\text{Level}(\mathbf{u}, z, x) \doteq \{y, t : \text{Contour}(y, t, z, x) \cap \mathbf{u} \neq \emptyset\} \tag{E.17}$$

$$\text{Capacity}(\mathbf{u}, z, x) \doteq Q(\text{Level}(\mathbf{u}, z, x) | \mathbf{z}, \mathbf{x}) \tag{E.18}$$

which I emphasise are functionals that depend implicitly upon the population and structural equations, and where $\mathbf{u} \in \mathfrak{U}^\infty$.

The most general characterisation of those admissible structures that are compatible with a given observable distribution is

$$\text{Artstein}(Q, z, x) \doteq \{P, h : P(\mathbf{u}) \leq \text{Capacity}(\mathbf{u}, z, x) \text{ for all } \mathbf{u} \in \mathfrak{U}^\infty\} \quad \text{E.19}$$

which is sharp (ARTSTEIN, 1983). If the intersection of Equation E.19 is finite for every observable distribution then the model has complete identification power; and if the intersection of Equation E.19 is empty for some observable distribution then the model is falsifiable.

I recall that heterogeneity is finite, in the sense that it reduces to a finite number of equivalence classes. As such, the power set enumerates all closed sets and unions of closed sets on its support—as is required by the theory of random sets from which Equation E.19 originates (MOLCHANOV, 2005). Given that the capacity functional and population are non-decreasing and increasing in heterogeneity, respectively, there exists a proper subset of the power set that is sufficient to characterise the identification region that is elsewhere labelled the class of core-determining sets (GALICHON and HENRY, 2011). I let

$$\mathfrak{S}(z, x) \doteq \{\text{Contour}(y, t, z, x) : y, t \in \{1, 0\} \times \{1, 0\}\} \quad \text{E.20}$$

and

$$\Theta(Q, z, x) \doteq \{P, h : P(\mathbf{u}) \leq \text{Capacity}(\mathbf{u}, z, x) \text{ for all } \mathbf{u} \in \mathfrak{S}(z, x)\} \quad \text{E.21}$$

for convenience. The generalised Balke-Pearl model identifies

$$\Theta_Q \doteq \bigcap_{x \in \mathfrak{X}} \bigcap_{z \in \mathfrak{Z} \leftarrow x} \Theta(Q, z, x) \quad \text{E.22}$$

as the identification region. That is, Equation E.19 reduces to Equation E.21, where Equation E.20 is the appropriate class of core-determining sets (CHESHER and ROSEN, 2017). Importantly, the class of core-determining sets contains only images of the contour functional corresponding to the specified values of the observable exogenous random variates, which is a consequence of how the contour functional connects.²⁵

The structural equations enact a labelling of heterogeneity (i.e., which realisation of heterogeneity induces which responses and treatments for particular combinations of the observable exogenous random

²⁵ The capacity functional is the complement of the containment functional (MOLCHANOV, 2005). It is straightforward to translate results from one to the other by, for example, substituting intersections for unions. The class of core-determining sets contains only images of the contour functional and their unions (CHESHER and ROSEN, 2017, §Lemma 1) that cannot be partitioned into disjoint images (CHESHER and ROSEN, 2017, §Theorem 3).

variates). Fixing this association normalises the structural equations and formulates the identification region as a linear programme over the population—to be precise, over the marginal distribution of heterogeneity. The identification region is, therefore, convex and non-finite due to the abundance of equivalence classes relative to the number of possible realisations of instruments and covariates (i.e., the linear programme has more unknowns than equations). Hence, the generalised Balke-Pearl model has incomplete identification power and is partially identifying since the identification region is non-trivial (MANSKI, 2003).

I reiterate that Equation E.22 is a sharp characterisation of the identification region. Although the solving of linear programmes typically incur low computational cost, the formulation of the constraint matrix above need not be feasible (the constraint matrix comprising $2^{2 \cdot (\mathfrak{r}+1) + \sum_{x \in \mathfrak{X}} \mathfrak{z}^{\leftarrow x}}$ columns and $\sum_{x \in \mathfrak{X}} \mathfrak{z}^{\leftarrow x}$ rows). As such, it is of practical and theoretical importance to prove the earlier statement that Equations E.8 and E.9 constitute sharp bounds on the Average Structural Functions.

Given Equation E.22, the first result that I am required to show—the set of structures that are admitted by the generalised Balke-Pearl model and that are compatible with the observable distribution is contained within the set of structures that are implied by Equation E.15—is immediate upon addition of the inequalities. Specifically,

$$P(Y_t(x) = y_1, Y_0(x) = y_0) = P(Y_t(x) = y_1, Y_0(x) = y_0, T_z(x) = 1) + P(Y_t(x) = y_1, Y_0(x) = y_0, T_z(x) = 0) \tag{E.23}$$

$$\leq P(Y_t(x) = y_1, T_z(x) = 1) + P(Y_0(x) = y_0, T_z(x) = 0) \tag{E.24}$$

$$= P(u \in \text{Contour}(y_1, 1, z, x)) + P(u \in \text{Contour}(y_0, 0, z, x)) \tag{E.25}$$

$$\leq Q_{y_1 1.z}(x) + Q_{y_0 0.z}(x) \tag{E.26}$$

and

$$P(Y_t(x) = 1 - y) \leq P(Y_t(x) = 1 - y, T_z(x) = t) + P(T_z(x) = 1 - t) \tag{E.27}$$

$$= P(u \in \text{Contour}(1 - y, t, z, x)) + \sum_y P(u \in \text{Contour}(y, 1 - t, z, x)) \tag{E.28}$$

$$\leq 1 - Q_{yt.z}(x) \tag{E.29}$$

which establishes the result. I could otherwise derive Equation E.15 directly via Equation E.19, as in BERESTEANU et al. (2012).

The fact that the class of core-determining sets contains only images of the contour functional corre-

sponding to the specified values of the observable exogenous random variates (i.e., there is separation of the marginal distributions) is an ever-present feature of this identification analysis. As a trivial extension of RICHARDSON and ROBINS (2024, §Lemma 4), if marginal distributions of the form

$$P_{z,x}^* (\{Y_1(x), Y_0(x) : x \in \mathfrak{X}\}, T_z(x)) \tag{E.30}$$

such that

$$P_{z,x}^* (\{Y_1(x), Y_0(x) : x \in \mathfrak{X}\}) = P^* (\{Y_1(x), Y_0(x) : x \in \mathfrak{X}\}) \tag{E.31}$$

can be proposed then there exists a single joint distribution of the form

$$P^* (\{Y_1(x), Y_0(x) : x \in \mathfrak{X}\}, \{T_z(x) : z \in \mathfrak{Z} \leftarrow x, x \in \mathfrak{X}\}) \tag{E.32}$$

that can be obtained from the marginal distributions. It is this result that motivates continued separation of the marginal distributions in this identification analysis.

KÉDAGNI and MOURIFIÉ (2020) proposes a collection of marginal distributions that that exhibit many of the properties that I require; I adapt this collection to suit the framework of the generalised Balke-Pearl model. I let

$$p_{y_1 y_0 t | z}(x) \doteq P^*(Y_1(x) = y_1, Y_0(x) = y_0, T_z(x) = t | z, x) \tag{E.33}$$

and

$$q_{i+j.z}(x) \doteq \min_{z \in \mathfrak{Z} \leftarrow x} (Q_{i1.z} + Q_{j0.z}) \circ (x) \tag{E.34}$$

$$\bar{Q}_{yt.z}(x) \doteq \max_{z \in \mathfrak{Z} \leftarrow x} (Q_{yt.z}(x))$$

for convenience, which is similar notation (although not exactly the same) to that used in KÉDAGNI and MOURIFIÉ (2020, in particular, I transpose the position of the first two subscripts relative to that paper). I propose

$$p_{100|z}(x) = \min(A_x, Q_{00.z}(x), L_x - Q_{11.z}(x), U_x + A_x - Q_{11.z}(x) - Q_{10.z}(x)) \tag{E.35}$$

$$p_{101|z}(x) = A_x - p_{100|z}(x) \tag{E.36}$$

$$p_{110|z}(x) = L_x - Q_{11.z}(x) - p_{100|z}(x) \tag{E.37}$$

$$p_{111|z}(x) = Q_{11.z}(x) - A_x + p_{100|z}(x) \tag{E.38}$$

$$p_{010|z}(x) = Q_{11.z}(x) + Q_{10.z}(x) - L_x + p_{100|z}(x) \quad \text{E.39}$$

$$p_{011|z}(x) = U_x + A_x - Q_{11.z}(x) - Q_{10.z}(x) - p_{100|z}(x) \quad \text{E.40}$$

$$p_{000|z}(x) = Q_{00.z}(x) - p_{100|z}(x) \quad \text{E.41}$$

$$p_{001|z}(x) = 1 - U_x - A_x - Q_{00.z}(x) + p_{100|z}(x) \quad \text{E.42}$$

where

$$\begin{aligned} A_x &\doteq \max(0, 1 - q_{1+0.z}(x) - q_{0+0.z}(x) - q_{0+1.z}(x), \bar{Q}_{00.z}(x) - q_{0+0.z}(x), \bar{Q}_{11.z}(x) - q_{1+1.z}(x)) \\ U_x &\doteq \min(1 - \bar{Q}_{00.z}(x), q_{1+1.z}(x) + q_{0+1.z}(x)) \\ L_x &\doteq \max(\bar{Q}_{11.z}(x), 1 - q_{0+0.z}(x) - q_{0+1.z}(x)) \end{aligned} \quad \text{E.43}$$

which is a minor modification—the only difference is that I append covariates—of the collection that is proposed in KÉDAGNI and MOURIFIÉ (2020). I then propose that

$$P^*(\{Y_1(x) = y_{1x}, Y_0(x) = y_{0x} : x \in \mathfrak{X}\}, T_z(x) = t|z, x) = p_{y_{1x}y_{0x}t}(x) \cdot \prod_{r \in \mathfrak{X}-x} \sum_t p_{y_{1r}y_{0r}t.\mathbf{e}_1^{\bar{z}}} \bar{z}^{\leftarrow r}(r) \quad \text{E.44}$$

such that the marginals of interest incorporate statistical independence of potential outcomes—specifically, those pertaining to response—for different values of the covariates.²⁶ By construction, Equations E.35 to E.42 satisfy Equation E.15. It remains for me to show that Equations E.35 to E.42 are compatible with the generalised Balke-Pearl model (i.e., that P^* is a proper distribution and that it satisfies

$$P^*(\{Y_1(x) = y_{1x}, Y_0(x) = y_{0x} : x \in \mathfrak{X}\}|z, x) = P^*(\{Y_1(x) = y_{1x}, Y_0(x) = y_{0x} : x \in \mathfrak{X}\}) \quad \text{E.45}$$

as is required) and generates the observable distribution.²⁷

I reiterate that Equations E.35 to E.42 amount to a minor modification of the collection that is proposed in KÉDAGNI and MOURIFIÉ (2020). I am, therefore, able to use several of the results that are stated therein: first, that Equations E.35 to E.42 are non-negative and sum to one (i.e., that the proposed marginal distributions are proper); second, that P^* more generally—being a product of proper marginal

²⁶ Something that the generalised Balke-Pearl model does not impose, but does not exclude either.

²⁷ I emphasise that I do not require that

$$P^*(\{Y_1(x) = y_{1x}, Y_0(x) = y_{0x} : x \in \mathfrak{X}\}, T_z(x)|z, x) = P^*(\{Y_1(x) = y_{1x}, Y_0(x) = y_{0x} : x \in \mathfrak{X}\}, T_z(x)) \quad \text{E.46}$$

here. I recall that RICHARDSON and ROBINS (2024, §Lemma 4) means that there exists a joint distribution satisfying a stronger condition than this provided that I can find marginals that align on $\{Y_1(x) = y_{1x}, Y_0(x) = y_{0x} : x \in \mathfrak{X}\}$; I claim that Equations E.35 to E.42 have this property.

distributions—is also proper; and third, that

$$p_{100|z}(x) + p_{101|z}(x) = A_x \tag{E.47}$$

$$p_{110|z}(x) + p_{111|z}(x) = L_x - A_x \tag{E.48}$$

$$p_{010|z}(x) + p_{011|z}(x) = U_x + A_x - L_x \tag{E.49}$$

$$p_{000|z}(x) + p_{001|z}(x) = 1 - U_x - A_x \tag{E.50}$$

which supports Equation E.45.²⁸ I remark that these properties can all be easily verified—the definition in Equation E.35 corresponds, for instance, to the maximum amount that can be subtracted from Equations E.36, E.37, E.40 and E.41 so that they remain non-negative, contingent upon an appropriate specification of the terms that are defined in Equation E.43. Moreover,

$$P^*(Y = j, T = 1|z, x) = P^*(Y_1(x) = j, T_z(x) = 1|z, x) \tag{E.51}$$

$$= p_{j11.z}(x) + p_{j01.z}(x) \tag{E.52}$$

$$= Q_{j1.z}(x) \tag{E.53}$$

$$= P(Y = j, T = 1|z, x) \tag{E.54}$$

and

$$P^*(Y = j, T = 0|z, x) = P^*(Y_0(x) = j, T_z(x) = 0|z, x) \tag{E.55}$$

$$= p_{1j0.z}(x) + p_{0j0.z}(x) \tag{E.56}$$

$$= Q_{j0.z}(x) \tag{E.57}$$

$$= P(Y = j, T = 0|z, x) \tag{E.58}$$

such that P^* and P are observationally equivalent. Given Equation E.22 and its primitive Equation E.21, however, this is also sufficient to determine that the proposed distribution is compatible with the restrictions that are imposed by the generalised Balke-Pearl model and that link the joint distribution over unobservable variates to the joint distribution over observable variates (i.e., P^* satisfies the inequalities in Equation E.21 with equality).

Having proved that Equation E.15 characterises $P(\{Y_1(x), Y_0(x) : x \in \mathfrak{X}\})$, I now turn my attention to the problem of combining elements of Equation E.15 to derive tight bounds on the Average Structural Functions. If a cover of the Average Structural Functions only uses observable quantities with the same

²⁸ I emphasise that Equation E.44 defines the left-hand side of Equation E.45 as the product of elements of Equations E.47 to E.50 and so is invariant to the exogenous observable variates.

value of covariates then Equations E.8 and E.9 are correct—RICHARDSON and ROBINS (2024, §Theorem 2) enumerates the vertices of the linear programme in this case. The only thing that I need to show then is that any bound that uses observable quantities with different values of covariates must be slack or is trivial. Clearly, any bound that uses observable quantities with different values of covariates cannot improve upon Equations E.8 and E.9 if it simply adds to Equations E.8 and E.9—all of the constituent types of the Average Structural Functions are already covered by Equations E.8 and E.9 (i.e., it is slack). As such, the only way that a bound that uses observable quantities with different values of covariates might improve upon Equations E.8 and E.9 is if there exists $P(Y_1(x) = y_1, Y_0(x) = y_0)$ that intersects with the Average Structural Function (or its complement) of interest and is not covered. For example, if interest is in $\bar{Y}_1(x)$ then the only possibilities for improvement are if (i.) $P(Y_1(x) = 1, Y_0(x) = 1)$ is covered but $P(Y_1(x) = 1, Y_0(x) = 0)$ is not, (ii.) $P(Y_1(x) = 1, Y_0(x) = 0)$ is covered but $P(Y_1(x) = 1, Y_0(x) = 1)$, or (iii.) if neither $P(Y_1(x) = 1, Y_0(x) = 1)$ nor $P(Y_1(x) = 1, Y_0(x) = 0)$ are covered. Here, by covered, I mean that the probability is bounded in the sense of Equation E.15 by a combination $Q_{11.z}(x) + Q_{y_0.z}(x)$. There exists, therefore, some $P(Y_1(x) = y_1, Y_0(x) = y_0, T_z(x) = t)$ that intersects with the Average Structural Function (or its complement) of interest and is wholly uncovered. For example, in (i.) $P(Y_1(x) = 1, Y_0(x) = 0, T_z(x) = 0)$ is wholly uncovered, in (ii.) $P(Y_1(x) = 1, Y_0(x) = 1, T_z(x) = 0)$ is wholly uncovered, and in (iii.) $P(Y_1(x) = 1, Y_0(x) = y, T_z(x) = t)$ is wholly uncovered, which subsumes the other two instances. Extending this argument to observable quantities with different values of covariates though, if a candidate does not cover the entirety of the support then there exists

$$P(\{Y_1(x^*) : x^* \in \mathfrak{X} - x\} = \mathbf{y}_1, \{Y_0(x^*) : x^* \in \mathfrak{X} - x\} = \mathbf{y}_0, \{T_z(x^*) : z \in \mathfrak{Z} \leftarrow x^*, x^* \in \mathfrak{X} - x\} = \mathbf{t}) \quad \text{E.59}$$

that is wholly uncovered. This implies, however, that

$$P(\{Y_1(x^*) : x^* \in \mathfrak{X}\} = \{\mathbf{y}_1, y_1\}, \{Y_0(x^*) : x^* \in \mathfrak{X}\} = \{\mathbf{y}_0, y_0\}, \{T_z(x^*) : z \in \mathfrak{Z} \leftarrow x^*, x^* \in \mathfrak{X}\} = \{\mathbf{t}, t\}) \quad \text{E.60}$$

(i.e., the intersection of the event in Equation E.59 with the event that is wholly uncovered in (i.), (ii.) or (iii.), say) is wholly uncovered, and so such a candidate does not constitute a proper cover of the Average Structural function. A proper cover involving observable quantities with different values of covariates must therefore cover the entirety of the support of heterogeneity, and so must be at least one (i.e., it is trivial); such a bound is clearly not an improvement upon Equations E.8 and E.9.

—| REFERENCES |—

BIBLIOGRAPHY

ABADIE, A., J. ANGRIST, and G. IMBENS (2002).

Instrumental variables estimates of the effect of subsidized training on the quantiles of trainee earnings.

In: *Econometrica*.

ABRAMOWITZ, M. and I. A. STEGUN (1964).

Handbook of mathematical functions with formulas, graphs, and mathematical tables.

US Government printing office.

ABREVVAYA, J., Y.-C. HSU, and R. P. LIELI (2015).

Estimating conditional average treatment effects.

In: *Journal of Business & Economic Statistics*.

ADJAHO, C. and T. CHRISTENSEN (2022).

Externally Valid Treatment Choice.

In: *arXiv preprint arXiv:2205.05561*.

ALQUIER, P., J. RIDGWAY, and N. CHOPIN (2016).

On the properties of variational approximations of Gibbs posteriors.

In: *The Journal of Machine Learning Research*.

AMOS, D. E. (1974).

Computation of modified Bessel functions and their ratios.

In: *Mathematics of Computation*.

ANDREWS, D. W. and G. SOARES (2010).

Inference for parameters defined by moment inequalities using generalized moment selection.

In: *Econometrica*.

ANGRIST, J. D. and W. N. EVANS (1998).

Children and Their Parents' Labor Supply: Evidence from Exogenous Variation in Family Size.

In: *The American Economic Review*.

ARIAS, J., J. RUBIO-RAMÍREZ, and D. WAGGONER (2018).

- Inference Based on Structural Vector Autoregressions Identified with Sign and Zero Restrictions:
Theory and Applications.
In: *Econometrica*.
- ARTSTEIN, Z. (1983).
Distributions of random sets and random selections.
In: *Israel Journal of Mathematics*.
- ATHEY, S. and S. WAGER (2021).
Policy learning with observational data.
In: *Econometrica*.
- AUMANN, R. J. (1965).
Integrals of set-valued functions.
In: *Journal of Mathematical Analysis and Applications*.
- BALKE, A. and J. PEARL (1997).
Bounds on treatment effects from studies with imperfect compliance.
In: *Journal of the American Statistical Association*.
- BÉGIN, L., P. GERMAIN, F. LAVIOLETTE, and J.-F. ROY (2014).
PAC-Bayesian theory for transductive learning.
In: *Artificial Intelligence and Statistics*.
- (2016).
PAC-Bayesian bounds based on the Rényi divergence.
In: *Artificial Intelligence and Statistics*.
- BERESTEANU, A., I. MOLCHANOV, and F. MOLINARI (2012).
Partial identification using random set theory.
In: *Journal of Econometrics*.
- BEST, D. and N. I. FISHER (1979).
Efficient simulation of the von Mises distribution.
In: *Journal of the Royal Statistical Society: Series C (Applied Statistics)*.
- BEYGELZIMER, A. and J. LANGFORD (2009).
The offset tree for learning with partial labels.
In: *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*.
ACM.
- BHATTACHARYA, D. and P. DUPAS (2012).
Inferring welfare maximizing treatment assignment under budget constraints.

- In: *Journal of Econometrics*.
- BHATTACHARYYA, A. (1943).
On a measure of divergence between two statistical populations defined by their probability distributions.
In: *Bull. Calcutta Math. Soc.*
- BINGHAM, C. (1974).
An antipodally symmetric distribution on the sphere.
In: *The Annals of Statistics*.
- BISSIRI, P. G., C. C. HOLMES, and S. G. WALKER (2016).
A general framework for updating belief distributions.
In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*.
- BLOOM, H. S., L. L. ORR, S. H. BELL, et al. (1997).
The benefits and costs of JTPA Title II-A programs: Key findings from the National Job Training Partnership Act study.
In: *Journal of human resources*.
- BLUNDELL, R. W. and J. L. POWELL (2004).
Endogeneity in semiparametric binary response models.
In: *The Review of Economic Studies*.
- BRETAGNOLLE, J. and C. HUBER (1978).
Estimation des densités: risque minimax.
In: *Séminaire de probabilités de Strasbourg*.
- BUGNI, F. A. (2010).
Bootstrap inference in partially identified models defined by moment inequalities: Coverage of the identified set.
In: *Econometrica*.
- CATONI, O. (2007).
PAC-Bayesian Supervised Classification.
In: *Lecture Notes-Monograph Series. IMS*.
- CHAMBERLAIN, G. (2011).
Bayesian aspects of treatment choice.
In: *The Oxford Handbook of Bayesian Econometrics*.
- CHERNOZHUKOV, V., S. LEE, and A. M. ROSEN (2013).
Intersection Bounds: Estimation and Inference.
In: *Econometrica*.

- CHERNOZHUKOV, V. and H. HONG (2003a).
 An MCMC approach to classical estimation.
 In: *Journal of Econometrics*.
- (2003b).
 An MCMC approach to classical estimation.
 In: *Journal of Econometrics*.
- CHESHER, A. and A. M. ROSEN (2013).
 What do instrumental variable models deliver with discrete dependent variables?
 In: *The American Economic Review*.
- (2017).
 Generalized instrumental variable models.
 In: *Econometrica*.
- (2020).
 Generalized instrumental variable models, methods, and applications.
 In: *Handbook of Econometrics*.
 Elsevier.
- CSABA, D. and B. SZOKE (2020).
 Learning with misspecified models.
 Tech. rep.
 mimeo.
- D’ADAMO, R. (2021).
 Policy Learning Under Ambiguity.
 In: *arXiv preprint arXiv:2111.10904*.
- DAWID, A. P. (2008).
 Beware of the DAG!
 In: *Proceedings of the 2008 International Conference on Causality: Objectives and Assessment*.
- DEHEJIA, R. H. (2005).
 Program evaluation as a decision problem.
 In: *Journal of Econometrics*.
- DERBEKO, P., R. EL-YANIV, and R. MEIR (2004).
 Explicit learning curves for transduction and application to clustering and compression algorithms.
 In: *Journal of Artificial Intelligence Research*.
- DHILLON, I. S. and S. SRA (2003).
 Modeling data using directional distributions.

- Tech. rep.
University of Texas at Austin, Department of Computer Sciences.
- DIETHE, T. (2015).
A Note on the Kullback-Leibler Divergence for the von Mises-Fisher distribution.
In: *arXiv preprint arXiv:1502.07104*.
- DWORK, C., M. HARDT, T. PITASSI, O. REINGOLD, and R. ZEMEL (2012).
Fairness through awareness.
In: *Proceedings of the 3rd innovations in theoretical computer science conference*.
- FISHER, R. A. (1953).
Dispersion on a sphere.
In: *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*.
- GALICHON, A. and M. HENRY (2011).
Set identification in models with multiple equilibria.
In: *The Review of Economic Studies*.
- GERMAIN, P., A. LACASSE, F. LAVIOLETTE, and M. MARCHAND (2009).
PAC-Bayesian learning of linear classifiers.
In: *Proceedings of the 26th Annual International Conference on Machine Learning*.
ACM.
- GIACOMINI, R. and T. KITAGAWA (2021a).
Robust Bayesian Inference for Set-identified Models.
In: *Econometrica*.
- GIACOMINI, R. and T. KITAGAWA (2021b).
Robust Bayesian inference for set-identified models.
In: *Econometrica*.
- GRONAU, R. (1977).
Leisure, home production, and work—the theory of the allocation of time revisited.
In: *Journal of political economy*.
- GUEDJ, B. (2019).
A primer on PAC-Bayesian learning.
In: *arXiv preprint arXiv:1901.05353*.
- GUNSILIUS, F. (2019).
A path-sampling method to partially identify causal effects in instrumental variable models.
In: *arXiv preprint arXiv:1910.09502*.
- (2020).

- A path-sampling method to partially identify causal effects in instrumental variable models.
 In: *arXiv preprint arXiv:1910.09502v2*.
- HAAVELMO, T. (1944).
 The probability approach in econometrics.
 In: *Econometrica*.
- HAN, S. (2022).
 Optimal dynamic treatment regimes and partial welfare ordering.
 In: *unpublished manuscript*.
- HECKMAN, J. J., H. ICHIMURA, and P. E. TODD (1997).
 Matching as an econometric evaluation estimator: Evidence from evaluating a job training programme.
 In: *The review of economic studies*.
- HECKMAN, J. J., S. URZUA, and E. J. VYTLACIL (2006).
 Understanding instrumental variables in models with essential heterogeneity.
 In: *The review of economics and statistics*.
- HILLEN, T., K. J. PAINTER, A. C. SWAN, and A. D. MURTHA (2017).
 Moments of von Mises and Fisher distributions and applications.
 In: *Mathematical Biosciences & Engineering*.
- HIRANO, K. and J. R. PORTER (2009).
 Asymptotics for statistical treatment rules.
 In: *Econometrica*.
- HORNIK, K. and B. GRÜN (2013).
 On conjugate families and Jeffreys priors for von Mises–Fisher distributions.
 In: *Journal of statistical planning and inference*.
- IACOVOU, M. (2001).
 Fertility and female labour supply.
 Institute for Social and Economic Research, University of Essex.
- IMBENS, G. W. and J. D. ANGRIST (1994).
 Identification and estimation of Local Average Treatment Effects.
 In: *Econometrica*.
- IMBENS, G. W. and D. B. RUBIN (1997).
 Estimating Outcome Distributions for Compliers in Instrumental Variables Models.
 In: *The Review of Economic Studies*.
- ISHIHARA, T. and T. KITAGAWA (2021).
 Evidence Aggregation for Treatment Choice.

- In: *arXiv preprint arXiv:2108.06473*.
- KAIDO, H., F. MOLINARI, and J. STOYE (2019).
Confidence intervals for projections of partially identified parameters.
In: *Econometrica*.
- KARLIN, S. and H. RUBIN (1956).
The theory of decision procedures for distributions with monotone likelihood ratio.
In: *The Annals of Mathematical Statistics*.
- KASY, M. (2016).
Partial identification, distributional preferences, and the welfare ranking of policies.
In: *Review of Economics and Statistics*.
- KÉDAGNI, D. and I. MOURIFIÉ (2020).
Generalized instrumental inequalities: testing the instrumental variable independence assumption.
In: *Biometrika*.
- KENT, J. T., A. M. GANEIBER, and K. V. MARDIA (2013).
A new method to simulate the Bingham and related distributions in directional data analysis with applications.
In: *arXiv preprint arXiv:1310.8110*.
- KENT, J. T. (1982).
The Fisher-Bingham Distribution on the Sphere.
In: *Journal of the Royal Statistical Society. Series B (Methodological)*.
- AL-KHAJA, A. (2016).
Essays on female empowerment and women’s status.
PhD thesis. University Colege London.
- KIDO, D. (2022).
Distributionally Robust Policy Learning with Wasserstein Distance.
In: *arXiv preprint arXiv:2205.04637*.
- KITAGAWA, T. (2021).
The identification region of the potential outcome distributions under instrument independence.
In: *Journal of Econometrics*.
- KITAGAWA, T., S. LEE, and C. QIU (2022a).
Treatment Choice with Nonlinear Regret.
In: *arXiv preprint arXiv:2205.08586*.
- KITAGAWA, T., H. LOPEZ, and J. ROWLEY (2022b).
Stochastic Treatment Choice with Empirical Welfare Updating.

- In: *arXiv preprint arXiv:2211.01537*.
- KITAGAWA, T., S. SAKAGUCHI, and A. TETENOV (2021).
 Constrained classification and policy learning.
 In: *arXiv preprint arXiv:2106.12886*.
- KITAGAWA, T. and A. TETENOV (2018a).
 Who should be treated? empirical welfare maximization methods for treatment choice.
 In: *Econometrica*.
- (2018b).
 Who should be treated? empirical welfare maximization methods for treatment choice.
 In: *Econometrica*.
- (2021).
 Equality-Minded Treatment Choice.
 In: *Journal of Business Economics and Statistics*.
- KLINE, B. and E. TAMER (2016).
 Bayesian inference in a class of partially identified models.
 In: *Quantitative Economics*.
- KOCK, A. B., D. PREINERSTORFER, and B. VELIYEV (2022).
 Treatment recommendation with distributional targets.
 In: *Journal of Econometrics*.
- KURZ, G. and U. D. HANEBECK (2015).
 Stochastic sampling of the hyperspherical von Mises–Fisher distribution without rejection methods.
 In: *2015 Sensor Data Fusion: Trends, Solutions, Applications (SDF)*.
 IEEE.
- LATTIMORE, T. and C. SZEPESVÁRI (2020).
 Bandit algorithms.
 Cambridge University Press.
- LIN, L., V. RAO, and D. DUNSON (2017).
 Bayesian nonparametric inference on the Stiefel manifold.
 In: *Statistica Sinica*.
- LIU, Y. (2022).
 Policy Learning under Endogeneity Using Instrumental Variables.
 In: *arXiv preprint arXiv:2206.09883*.
- MANSKI, C. F. (1975).
 Maximum Score Estimation of the Stochastic Utility Model of Choice.

- In: *Journal of Econometrics*.
- (1985).
Semiparametric Analysis of Discrete Response: Asymptotic Properties of the Maximum Score Estimator.
In: *Journal of Econometrics*.
 - (2000a).
Identification problems and decisions under ambiguity: empirical analysis of treatment response and normative analysis of treatment choice.
In: *Journal of Econometrics*.
 - (2000b).
Using Studies of Treatment Response to Inform Treatment Choice in Heterogeneous Populations.
 - (2003).
Partial identification of probability distributions.
Springer Science & Business Media.
 - (2004a).
Statistical treatment rules for heterogeneous populations.
In: *Econometrica*.
 - (2004b).
Statistical treatment rules for heterogeneous populations.
In: *Econometrica*.
 - (2005).
Social choice with partial knowledge of treatment response.
Princeton University Press.
 - (2007a).
Identification for prediction and decision.
Harvard University Press.
 - (2007b).
Minimax-regret treatment choice with missing outcome data.
In: *Journal of Econometrics*.
 - (2009).
The 2009 Lawrence R. Klein Lecture: Diversified treatment under ambiguity.
In: *International Economic Review*.
 - (2022).
Identification and Statistical Decision Theory.

- In: *arXiv preprint arXiv:2204.11318*.
- MANSKI, C. F. and A. TETENOV (2007).
Admissible treatment rules for a risk-averse planner with experimental data on an innovation.
In: *Journal of Statistical Planning and Inference*.
- MANSKI, C. F. (1989).
Anatomy of the Selection Problem.
In: *The Journal of Human Resources*.
- MARDIA, K. V. (1975).
Statistics of Directional Data.
In: *Journal of the Royal Statistical Society. Series B (Methodological)*.
- MARDIA, K. V. and S. EL-ATOUM (1976).
Bayesian inference for the von Mises-Fisher distribution.
In: *Biometrika*.
- MARDIA, K. V., G. HUGHES, C. C. TAYLOR, and H. SINGH (2008).
A multivariate von Mises distribution with applications to bioinformatics.
In: *Canadian Journal of Statistics*.
- MARDIA, K. V. and P. E. JUPP (2009).
Directional statistics.
John Wiley & Sons.
- MAURER, A. (2004).
A note on the PAC Bayesian theorem.
In: *arXiv preprint cs/0411099*.
- MBAKOP, E. and M. TABORD-MEEHAN (2021).
Model selection for treatment choice: Penalized welfare maximization.
In: *Econometrica*.
- MCALLESTER, D. A. (1999).
Some PAC-Bayesian Theorems.
In: *Machine Learning*.
- (2003).
PAC-Bayesian stochastic model selection.
In: *Machine Learning*.
- MOLCHANOV, I. (2005).
Theory of Random Sets.
Springer verlag.

- MOON, H. R. and F. SCHORFHEIDE (2012).
Bayesian and frequentist inference in partially identified models.
In: *Econometrica*.
- MOURIFIÉ, I. (2015).
Sharp bounds on treatment effects in a binary triangular system.
In: *Journal of Econometrics*.
- NIE, X., E. BRUNSKILL, and S. WAGER (2021).
Learning When-to-Treat Policies.
In: *Journal of the American Statistical Association*.
- NIST (2021).
Digital Library of Mathematical Functions, Release 1.1.1 of 2021-03-15.
F. W. J. Olver, A. B. Olde Daalhuis, D. W. Lozier, B. I. Schneider, R. F. Boisvert, C. W. Clark, B. R. Miller, B. V. Saunders, H. S. Cohl, and M. A. McClain, eds.
- OECD (2007).
Mothers in Paid Employment.
In.
- PEARL, J. (2009).
Causality.
Cambridge University Press.
- PELLATT, D. F. (2022).
PAC-Bayesian Treatment Allocation Under Budget Constraints.
In: *arXiv preprint arXiv:2212.09007*.
- PENTINA, A. and C. H. LAMPERT (2015).
Lifelong learning with non-iid tasks.
In: *Advances in Neural Information Processing Systems*.
- PINSKER, M. S. (1964).
Information and information stability of random variables and processes.
Holden-Day.
- RACINE, J. S., L. SU, and A. ULLAH (Feb. 2014).
The Oxford Handbook of Applied Nonparametric and Semiparametric Econometrics and Statistics.
Oxford University Press.
- RAI, Y. (2019).
Statistical Inference for Treatment Assignment Policies.
In: *unpublished manuscript*.

- RICHARDSON, T. S. and J. M. ROBINS (2014).
 ACE bounds; SEMs with equilibrium conditions.
 In: *Statistical Science*.
- (2024).
 Assumptions and bounds in the instrumental variable model.
 In: *arXiv preprint arXiv:2401.13758*.
- ROSENBAUM, P. R. and D. B. RUBIN (1983).
 The central role of the propensity score in observational studies for causal effects.
 In: *Biometrika*.
- RUBIN, D. B. (1974).
 Estimating causal effects of treatments in randomized and nonrandomized studies.
 In: *Journal of Educational Psychology*.
- RUGGLES, S., K. GENADEK, R. GOEKEN, J. GROVER, and M. SOBEK (2015).
 Integrated Public Use Microdata Series: Version 6.0 [dataset].
 Minneapolis, MN: University of Minnesota.
- SABELFELD, K. K. (2018).
 Application of the von Mises–Fisher distribution to Random Walk on Spheres method for solving high-dimensional diffusion–advection–reaction equations.
 In: *Statistics & Probability Letters*.
- SAKAGUCHI, S. (2019).
 Estimating Optimal Dynamic Treatment Assignment Rules under Intertemporal Budget Constraints.
 Tech. rep.
 cemmap working paper.
- SASAKI, Y. and T. URA (2020).
 Welfare analysis via marginal treatment effects.
 In: *arXiv preprint arXiv:2012.07624*.
- SCHLAG, K. H. (2006).
 ELEVEN-Tests needed for a Recommendation.
 In: *Economics Working Papers*.
- SEEGER, M. (2002).
 Pac-bayesian generalisation error bounds for gaussian process classification.
 In: *Journal of machine learning research*.
- SEI, T., H. SHIBATA, A. TAKEMURA, K. OHARA, and N. TAKAYAMA (2013).
 Properties and applications of Fisher distribution on the rotation group.

- In: *Journal of Multivariate Analysis*.
- SHAIKH, A. M. and E. J. VYTLACIL (2011).
 Partial identification in triangular systems of equations with binary dependent variables.
 In: *Econometrica*.
- SHAWE-TAYLOR, J. and R. C. WILLIAMSON (1997).
 A PAC analysis of a Bayesian estimator.
 In: *Annual Workshop on Computational Learning Theory: Proceedings of the tenth annual conference on Computational learning theory*.
- SPLAWA-NEYMAN, J. (1923).
 On the application of probability theory to agricultural experiments. Essay on principles.]
 In: *Roczniki Nauk Rolniczych*.
- STOYE, J. (2009).
 Minimax regret treatment choice with finite samples.
 In: *Journal of Econometrics*.
- (2012).
 Minimax regret treatment choice with covariates or with limited validity of experiments.
 In: *Journal of Econometrics*.
- STROTZ, R. H. and H. O. WOLD (1960).
 Recursive vs. nonrecursive systems: An attempt at synthesis (part i of a triptych on causal chain systems).
 In: *Econometrica: Journal of the Econometric Society*.
- SUN, L. (2021).
 Empirical welfare maximization with constraints.
 In: *arXiv preprint*.
- SWAMINATHAN, A. and T. JOACHIMS (2015).
 Counterfactual risk minimization: Learning from logged bandit feedback.
 In: *International Conference on Machine Learning*.
- TETENOV, A. (2012).
 Statistical treatment choice based on asymmetric minimax regret criteria.
 In: *Journal of Econometrics*.
- THOMPSON, W. R. (1933).
 On the likelihood that one unknown probability exceeds another in view of the evidence of two samples.
 In: *Biometrika*.
- TSAGRIS, M., G. ATHINEOU, C. ADAM, et al. (2022).

- Directional: A Collection of Functions for Directional Data Analysis.
R package version 5.5.
- UHLIG, H. (2005).
What are the Effects of Monetary Policy on Output? Results from an Agnostic Identification Procedure.
In: *Journal of Monetary Economics*.
- VALIANT, L. G. (1984).
A theory of the learnable.
In: *Proceedings of the sixteenth annual ACM symposium on Theory of computing*.
ACM.
- VAN ERVEN, T. and P. HARREMOËS (2014).
Rényi divergence and Kullback-Leibler divergence.
In: *IEEE Transactions on Information Theory*.
- VAPNIK, V. N. (1998).
Statistical Learning Theory.
Wiley.
- VIVIANO, D. (2021).
Policy targeting under network interference.
In: *arXiv preprint*.
- VON MISES, R. (1918).
Über die 'Ganzzahligkeit' der Atomgewichte und verwandte Fragen.
In: *Phys. Z.*
- VYTLACIL, E. J. (2002).
Independence, monotonicity, and latent index models: An equivalence result.
In: *Econometrica*.
- WAINWRIGHT, M. J. and M. I. JORDAN (2008).
Graphical Models, Exponential Families, and Variational Inference.
In: *Foundations and Trends in Machine Learning*.
- WATSON, G. S. (1984).
The theory of concentrated Langevin distributions.
In: *Journal of multivariate analysis*.
- WOOD, A. T. (1994).
Simulation of the von Mises Fisher distribution.
In: *Communications in statistics-simulation and computation*.

- WOOLDRIDGE, J. M. (2015).
Control function methods in applied econometrics.
In: *Journal of Human Resources*.
- YATA, K. (2021).
Optimal Decision Rules Under Partial Identification.
In: *arXiv preprint arXiv:2111.04926*.
- ZADROZNY, B. (2003).
Policy mining: Learning decision policies from fixed sets of data.
PhD thesis. University of California, San Diego.
- ZHANG, B., A. A. TSIATIS, M. DAVIDIAN, M. ZHANG, and E. LABER (2012).
Estimating optimal treatment regimes from a classification perspective.
In: *Stat*.
- ZHAO, Y., D. ZENG, A. J. RUSH, and M. R. KOSOROK (2012).
Estimating individualized treatment rules using outcome weighted learning.
In: *Journal of the American Statistical Association*.

—| STATEMENT |—

DIVISION OF LABOUR

This statement—regarding the division of labour for multi-authored work—supplements the preceding declaration.

—| DECLARATION 1 |—

THE f -DIVERGENCE OF A VON MISES-FISHER DISTRIBUTION FROM SOME REFERENCE DISTRIBUTIONS

Chapter C extracted from a paper co-authored with Toru Kitagawa.

The idea for the paper was arrived at jointly by Toru and I. We recognised that there was an absence of results concerning the statistical divergence of spherical distributions—in particular, von Mises-Fisher distributions—having already identified a use for such results. I conducted preliminary research into the topic and we agreed that it was feasible to proceed with the project.

Work on the project was divided as follows. I was responsible for obtaining the main theoretical results, and for drafting and editing the resulting paper; Toru was responsible for reviewing the resulting paper, and for providing helpful discussion and perspective throughout the project.

Given this allocation of responsibility, I find it appropriate to describe only Toru’s influence on the project and its output. First, Toru suggested the use of the moment-generating function—with its well-known form for exponential families—to derive the moments of a von Mises-Fisher distribution.¹ This approach was simpler than the approach that I had taken up until that point, which was to perform integration by substitution involving products of trigonometric functions. This simplicity extended to the resulting expressions of the moments. Second, Toru shaped the introduction by linking the subject-matter to a common economic model—the canonical binary choice model—and then to several other economic problems—optimisation of maximum score or empirical welfare criteria, and underidentifying linear

¹ The moments of the distribution—specifically, the first moment—are necessary to obtain an closed-form expression for the Kullback-Leibler divergence. We originally intended that the paper would also present the moments of a von Mises-Fisher distribution as a result in itself, but Toru identified that some of these results were available elsewhere.

simultaneous equation models. Third, Toru read-through the paper and suggested minor changes.

As a reflection of this attribution of work and output, I include the paper in its entirety herein.

—| DECLARATION 2 |—

STOCHASTIC TREATMENT CHOICE WITH EMPIRICAL WELFARE UPDATING

Chapter D extracted from a paper co-authored with Toru Kitagawa and Hugo Lopez.

The idea for the paper was arrived at jointly by Toru and Hugo. I was assigned to the project to provide research assistance and later upgraded to a co-author in recognition of my contribution to it. At the time I was assigned to the project, some theoretical results had been obtained and organised into an initial draft. I attribute these results to both Toru and Hugo since I am unaware of how work on the project was divided prior to my involvement (although I presume that Hugo prepared the initial draft).

The initial draft included an introduction, a framework, and several of the main theoretical results (first theorem and its proof, second theorem and its proof, third theorem and its proof, and first lemma and its proof). I was directed to apply these theoretical results to data and to rewrite the initial draft, which was incomplete and a work in progress. My influence on the paper can be seen throughout—due to my preparation of it—but my contribution to its substance is apparent in the proof of the second theorem, the third theorem and its proof, and parts that were not included in the initial draft.

Work on the project subsequent to my involvement was divided as follows. I was responsible for verifying the existing main theoretical results and for correcting any issues with these, for applying the aforementioned results to data, for rewriting the initial draft, and for editing the resulting paper; Toru was responsible for correcting any issues with the existing main theoretical results, and for providing helpful discussion and perspective throughout the project; Hugo was responsible for correcting any issues with the existing main theoretical results, and for editing the resulting paper.

Given this allocation of responsibility, I find it appropriate to describe the issues that I identified with the initial draft before discussing how these issues were corrected. First, the proof of the second theorem was incorrect and somewhat incomplete—although the theorem itself was correct. Second, the third theorem was stated for a von Mises—rather than a von Mises-Fisher—distribution and was, therefore, quite different.² Third, the proof of the third theorem was also incomplete in those parts that remained applicable.

² Although the idea underpinning the third theorem remains the same—find an expression for the constant that appears in the first lemma in terms of the sample size such that the right-hand side of the inequality converges—the substance of this calculation is different.

The issue with the second theorem was due to the non-convexity of the objective function, which made establishing the veracity of the derivative as a solution for the case of a continuous prior difficult. Toru, Hugo and I all worked on this issue for several months, holding regular meetings to discuss our progress, with approaches ranging from contraction mapping to functional derivatives all rejected. The eventual proof is the culmination of our collective effort—although I note that the key insight that the derivative is monotone in the Lagrange multiplier is due to Toru.

The issue with the third theorem was that a von Mises distribution was—as I understand it—initially selected for two reasons—one valid, and one based upon an erroneous assumption. I persuaded Toru and Hugo to instead work with von Mises-Fisher distributions and provided the theoretical results that were necessary to do this. Toru, Hugo and I worked together to adapt the initial proof of the third theorem to incorporate these results. The eventual proof is the culmination of our collective effort.

The issue with those parts of the proof of the third theorem that remained applicable were that they were rather vague as to the rate of convergence. I worked independently to derive an expression for the constant and in doing so provided a concise expression of this rate.

As a reflection of this attribution of work and output, I include a partial reproduction of the paper herein—I omit several proofs (proof of the first theorem and proof of the first lemma) but include the associated theorems for the reader's comprehension. My influence on the substance of the paper is most apparent in the third theorem and its proof, and in the empirical application and its associated appendix.