



Pitfalls in post hoc analyses of population receptive field data

Susanne Stoll^{a,*}, Elisa Infanti^a, Benjamin de Haas^b, D. Samuel Schwarzkopf^c

^a Experimental Psychology, University College London, 26 Bedford Way, London, WC1H 0AP, UK

^b Abteilung Allgemeine Psychologie, Justus-Liebig-Universität Gießen, Otto-Behagel-Str. 10F, 35394 Gießen, Germany

^c School of Optometry and Vision Science, The University of Auckland, Private Bag 92019, Auckland 1142, New Zealand

ARTICLE INFO

Keywords:

Regression towards the mean
Circularity
Double-dipping
Validation
Functional magnetic resonance imaging

ABSTRACT

Data binning involves grouping observations into bins and calculating bin-wise summary statistics. It can cope with overplotting and noise, making it a versatile tool for comparing many observations. However, data binning goes awry if the same observations are used for binning (selection) and contrasting (selective analysis). This creates circularity, biasing noise components and resulting in artifactual changes in the form of regression towards the mean. Importantly, these artifactual changes are a statistical necessity. Here, we use (null) simulations and empirical repeat data to expose this flaw in the scope of post hoc analyses of population receptive field data. In doing so, we reveal that the type of data analysis, data properties, and circular data cleaning are factors shaping the appearance of such artifactual changes. We furthermore highlight that circular data cleaning and circular sorting of change scores are selection practices that result in artifactual changes even without circular data binning. These pitfalls might have led to erroneous claims about changes in population receptive fields in previous work and can be mitigated by using independent data for selection purposes. Our evaluations highlight the urgency for us researchers to make the validation of analysis pipelines standard practice.

1. Introduction

Data binning refers to grouping observations into bins or subgroups and calculating bin-wise summary statistics, such as the arithmetic mean. It is often applied to large datasets in order to prevent overplotting and control noise. As such, data binning has become commonplace in population receptive field (pRF) modeling (Dumoulin and Knapen, 2018; Dumoulin and Wandell, 2008), where researchers are commonly interested in comparing visual field maps with thousands of observations between different (experimental) conditions. However, pRF modeling is only one out of several research areas where some form of differential data binning has been adopted (e.g., Gignac and Zajenkowski, 2020; Holmes, 2009; Kriegeskorte et al., 2009; Preacher et al., 2005; Shanks, 2017).

Although data binning can help us see an overall pattern in the face of an abundance of details, it goes awry if the same observations are used for *binning* (selection) and *contrasting* (selective analysis). This is because dipping into noise-tainted data (i.e., most data) more than once violates assumptions of independence, favoring some noise components over others and eventually biasing descriptive and inferential statistics (Kriegeskorte et al., 2009). As such, double-dipping in data binning prevents us from – amongst other things – controlling for *regression towards*

the mean (e.g., Galton, 1886; Gignac and Zajenkowski, 2020; Holmes, 2009; Makin and De Xivry, 2019; Shanks, 2017; Stigler, 1997).

Regression towards the mean is a statistical artifact occurring when two variables are imperfectly correlated (e.g., due to random noise¹). In this case, extreme observations for one variable will on average be less extreme for the other² (e.g., Campbell and Kenny, 1999; Cohen et al., 2003; Galton, 1886; Shanks, 2017; Stigler, 1997). The magnitude of regression towards the mean tends to be higher the lower the correlation between the variables (e.g., Campbell and Kenny, 1999, for systematic simulations, see Holmes 2009).

Double-dipping and/or regression towards the mean are of particular concern in what is known as *post hoc subgrouping* (Preacher et al., 2005), *post hoc data selection* (Shanks, 2017), and *extreme groups approach* (Preacher et al., 2005), all of which can be considered as subtypes of data binning. Post hoc subgrouping refers to collecting two measures, defining extreme subgroups post hoc using one measure (e.g., the lower and upper quantile), and then performing statistics on these measures for

¹ Note that random noise is only one factor weakening the correlation between two variables (for more details, see Shanks, 2017).

² To be precise, regression towards the mean refers to standard scores (z -scores; Campbell and Kenny, 1999; Kenny, 2005).

* Corresponding author.

E-mail address: stollsus@gmail.com (S. Stoll).

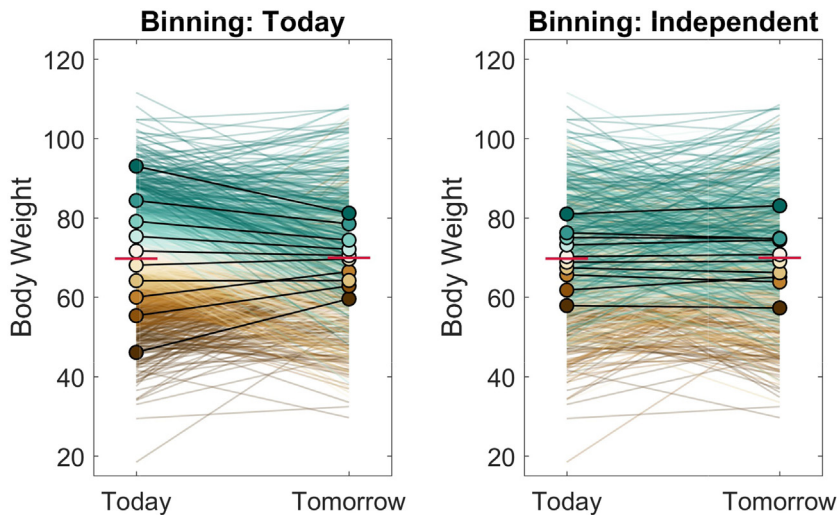


Fig. 1. Simulated post hoc binning analysis on fictive body weight data. Bin-wise fictive body weight data and means for Today and Tomorrow in the same group of adults and different data binning scenarios. Data for Today and Tomorrow were either binned according data for Today (1st column) or an Independent test occasion (2nd column). Fictive body weight data were simulated by sampling the body weight of 1000 adults from a Gaussian distribution ($M = 70$ kg; $SD = 10$ kg) and disturbing each adult's body weight with random Gaussian noise ($SD = 10$ kg), separately for each test occasion (Today, Tomorrow, and Independent). The red horizontal lines indicate the location of the overall mean for Today and Tomorrow. Dark brown colors correspond to lower and dark blue-green colors to higher decile bins. The endpoints of the colorful lines represent individual data points and the colorful dots with the black outline represent bin-wise means. Note that the graphs displayed here are referred to as Galton squeeze diagrams (Campbell and Kenny, 1999; Galton, 1886; Shanks, 2017).

the extreme subgroups (Preacher et al., 2005). Post hoc data selection is similar but involves only one extreme subgroup (Shanks, 2017). Both of these practices are different from the extreme groups approach, where extreme subgroups are selected a priori based on one measure; that is, without collecting the whole range of the other measure (Preacher et al., 2005). Here, we focus on a post hoc scenario where essentially all subgroups are considered, not just the extreme ones (see also Gignac and Zajenkowski, 2020; Holmes, 2009). We label this procedure including its subtypes *post hoc binning analysis*.

An intuitive way to think about the link between double-dipping, regression towards the mean, and post hoc binning are repeat data. Assume we measure body weight in a population of adults twice – Today and Tomorrow (see endpoints of colorful lines, Fig. 1; 1st column). Further assume that any weight we measure involves a *permanent* and a *transient* component (true value + random noise). When determining Today's and Tomorrow's overall mean weight, all things being equal, the permanent component persists and the transient component cancels out (see red horizontal lines, Fig. 1; 1st column). However, this is not the case when we select adults with extremely high measurements for Today (relative to the overall mean) and compare these measurements to Tomorrow's in the same adults by calculating the means (see lines and dots in dark green color, Fig. 1; 1st column). This is because we used Today's measurements twice: for selection (binning) and selective analysis (comparing bin-wise means). We therefore favored Today's noise components over Tomorrow's. Why is this? The noise components of our selection criterion are not independent of the noise components of Today's measurements. This renders the subgroup we selected Today on average heavier than it really is. This is not the case for Tomorrow's measurements. As a result, Tomorrow's measurements for this subgroup regress on average to Tomorrow's overall mean (see dots in dark green color, Fig. 1; 1st column; for a similar example see Stigler, 1997). This artifactual change in average weight might look like a real phenomenon, although – of course – it is not.

The analysis we just performed can be regarded as an instantiation of post hoc data selection involving one extreme subgroup. If we additionally select a subgroup of adults with extremely low measurements for Today (see lines and dots in dark brown color, Fig. 1; 1st column), regression towards the overall mean from below occurs for this subgroup. Such an approach would qualify as post hoc subgrouping involving two extreme subgroups. If we incorporate additional less extreme subgroups, we perform a full-blown post hoc binning analysis (see lines and dots in various colors, Fig. 1; 1st column), where the bin-wise means for Tomorrow's measurements regress towards the overall mean to various degrees. Importantly, this regression artifact is a statistical necessity not hinging upon body weight data. Once we use Independent data for

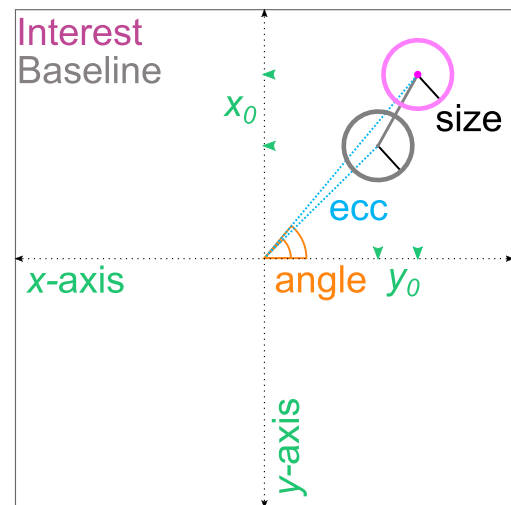


Fig. 2. Population receptive field estimates. The large black square outline represents a cutout of the visual field and the black dashed arrows a Cartesian coordinate system. The two circles represent a pRF that changes its position (gray solid line) in an Interest (magenta) compared to a Baseline (gray) condition. The pRF was modeled as a 2D Gaussian function. The center of the 2D Gaussian (midpoint of the gray and magenta circles) represents the position of the pRF. PRF position can be expressed in terms of x_0 and y_0 coordinates (green arrow heads) or eccentricity (blue dashed line) and polar angles (orange solid line). Eccentricity corresponds to the Euclidean distance between the center of gaze (origin) and the center of the 2D Gaussian. Polar angle corresponds to the counter-clockwise angle running from the positive x -axis to the eccentricity vector. The standard deviation of the Gaussian (1σ ; black solid line) represents pRF size. Both pRF position and size are typically expressed in degrees of visual angle. Polar angles are typically expressed in degrees. Ecc = Eccentricity. pRF = Population receptive field.

binning purposes (e.g., body weight measurements collected for the day after tomorrow), we break the circularity, and the regression artifact disappears (Fig. 1, 2nd column).

How does all of this relate to post hoc analyses involving pRF data? Imagine we conduct a retinotopic mapping experiment (Dumoulin and Wandell, 2008), where we estimate pRF position and pRF size for each voxel in the visual brain under a *Baseline* condition as well as a condition of *Interest* (see Fig. 2 for a single pRF). We can think of the Interest and Baseline conditions as repeat data (e.g., Benson et al., 2018; van Dijk et al., 2016; Senden et al., 2014), different attention conditions

(e.g., van Es et al., 2018; de Haas et al., 2014; 2020; Klein et al., 2014; Vo et al., 2017), mapping sequences (e.g., Binda et al., 2013; Infanti and Schwarzkopf, 2020), mapping stimuli (e.g., Alvarez et al., 2015; Binda et al., 2013; Le et al., 2017; Yildirim et al., 2018), magnetic field strengths (e.g., Morgan and Schwarzkopf, 2020), scotoma conditions (e.g., Barton and Brewer, 2015; Binda et al., 2013; Haak et al., 2012; Prabhakaran et al., 2020), and pRF modeling techniques (e.g., Carvalho et al., 2020) – to name but a few examples. Similarly, apart from visual scenarios, we can also interpret the Baseline and Interest condition as adaptation conditions (e.g., Tsouli et al., 2021), different finger movements (e.g., Schellekens et al., 2018), or uni- and multisensory conditions (see Holmes, 2009, for a discussion on non-pRF work).

As a pRF model, we adopt a 2D Gaussian, where pRF position represents the center of a pRF in visual space (the center of the Gaussian) and pRF size its spatial extent (the standard deviation of the Gaussian; see Fig. 2). We then fit this model to the voxel-wise brain responses we measured in the retinotopic mapping experiment (Dumoulin and Wandell, 2008). To compare pRF positions in the Interest and Baseline condition voxel-by-voxel, we bin the pRF positions from both conditions according to the pRF positions from the Baseline condition. Subsequently, we quantify for each voxel the position shift from the Baseline to the Interest condition (see Fig. 2 for a single pRF). Finally, we calculate the bin-wise mean shift. This is equivalent to calculating the bin-wise simple means for each condition and comparing them subsequently.

Either way, by adopting such a post hoc binning analysis, we essentially assume that binning voxels according to pRF positions from the Baseline condition and aggregating them subsequently for this condition ensures that bin-wise noise components are unbiased on average (see also Shanks, 2017). This, however, is not the case. The underlying reason is the same as for our body weight analysis further above: we dipped into the Baseline condition twice, namely to define bins (selection) and to estimate bin-wise means for further comparison (selective analysis). This circularity leads to a favoring of noise components, skewing the bin-wise means in the Baseline condition and eventually resulting in regression towards the overall mean for the bin-wise means of the Interest condition.

Here, we expose and explore this flaw in the scope of post hoc analyses of pRF data using (null) simulations and empirical repeat data from the Human Connectome Project (HCP; Benson et al., 2018; 2020). Unlike empirical data, simulations allowed us to separate true values from noise components. They also provided an excellent test bed for determining that the type of data analysis (change scores or simple scores, 1D or 2D binning, equidistant or decile binning), data properties (presence or absence of heteroskedasticity or a true effect) and additional circular selection practices (presence or absence of circular data cleaning) influence the appearance of the regression artifact. Moreover, they allowed us to pinpoint that circular data cleaning and circular sorting of change scores represent selection practices that yield artifactual changes even without circular data binning. Unlike empirical data from different experimental conditions, repeat data permitted us to assume a null effect between conditions, allowing for more straightforward conclusions about any systematic differences we might observe.

2. Methods

2.1. Post hoc binning using simulated data

For the post hoc binning analysis involving simulations, we used an empirical V1 visual field map of a single human participant as a basic data distribution. This map originated from a functional magnetic resonance imaging experiment (fMRI) aimed at mapping pRFs under different attention conditions using a drifting bar stimulus (2 sessions each with 4 runs per condition). One of these conditions was selected for simulation purposes. The maximal eccentricity of the mapping area subtended 8.5 degrees of visual angle (dva). We fit a 2D Gaussian function to preprocessed fMRI responses projected onto the cortical surface. For

each vertex (gray matter node on the cortical surface), we obtained 6 estimates: pRF position (x_0 and y_0 coordinates), pRF size (σ), pRF baseline (β_0), pRF amplitude (β_1), and goodness-of-fit (R^2). We first smoothed the resulting parameter maps and delineated V1 hemifield maps manually (for more details, see Supplementary methods, 1. Retinotopic mapping experiment). We then pooled the x_0 and y_0 coordinates across V1 hemifield maps and removed empty data points.

2.1.1. 1D post hoc binning analysis on eccentricity

To uncover the regression artifact, we first simulated a simplified contrast scenario with a null effect. To this end, we disturbed the x_0 and y_0 coordinates (Fig. 2) 200 times with random Gaussian noise ($SD = 2$ dva). We repeated this to generate a *Baseline*, *Interest*, and *Independent* condition. We then converted the x_0 and y_0 coordinates to eccentricity values (Fig. 2), as is often done in the pRF literature (see Fig. s1 for interpretational difficulties with eccentricity when it comes to position shifts). This resulted in a gamma-like eccentricity distribution. Lastly, we binned the eccentricity values in the Baseline and Interest condition according to the eccentricity values of any of the 3 conditions using deciles and calculated the bin-wise means.³ A schematic workflow of this simulated 1D post hoc binning analysis can be found in Fig. 3. Bin-wise eccentricity means were visualized as a color-coded scatter plot along with individual observations per bin and marginal histograms (bin width = 0.5 dva) reflecting the simulated distributions.

Building upon the simulated null effect, we performed the 1D post hoc binning analysis on 4 more simulation cases: a null effect with condition cross-thresholding based on the Baseline condition, a null effect with condition cross-thresholding based on both the Baseline and Interest condition, a null effect with eccentricity-dependent noise, and a true effect. We use the term ‘condition cross-thresholding’ to refer to the pairwise or list-wise deletion of data points across experimental conditions (see below). The selected simulation cases reflect analysis practices and data properties we consider characteristic of pRF studies. For all simulation cases, the Independent condition consisted of a second draw (re-sample) of the Baseline condition. Moreover, to ensure reproducibility and comparability, all simulation cases were based on the same seed for random number generation. However, our conclusions do not depend on the choice of seed for random number generation.

For the simulation cases involving condition cross-thresholding, we removed simulated observations falling outside a certain eccentricity range (≥ 0 and ≤ 6 dva) in the Baseline or Baseline and Interest condition from all conditions (i.e., Baseline, Interest, and Independent). For the simulation case involving eccentricity-dependent noise, we used a small standard deviation ($SD = 0.25$ dva) of random Gaussian noise to disturb empirical observations with smaller eccentricities (≥ 0 and < 3 dva) and a larger standard deviation ($SD = 2$ dva) to disturb empirical observations with larger eccentricities (≥ 3 dva). For the simulation case involving a true effect, we induced a radial increase in eccentricity of 2 dva in the Interest condition.

Apart from simple bin-wise means, we performed the 1D post hoc binning analysis also on change scores. The change scores were obtained by subtracting individual simulated observations or means in the Baseline condition from those in the Interest condition. Both simple means and mean change scores have been used for post hoc binning in previous pRF studies (e.g., Barton and Brewer, 2015; Binda et al., 2013; Carvalho et al., 2020; Haak et al., 2012; de Haas et al., 2014; 2020; Prabhakaran et al., 2020; Tsouli et al., 2021; Yildirim et al., 2018). Similarly, we re-

³ Note that when evaluating data distributions with unequal means, variances, or non-linearity, z -standardization might be necessary to detect regression towards or away from the mean (Campbell and Kenny, 1999; Shanks, 2017). In particular, z -standardization makes data distributions directly comparable. As such, bin-wise means should regress to wherever they intersect the identity line. Here, we always display data in native space, as this is typically done in the pRF literature. However, we use crosshairs to indicate the location of the mean and thus provide a visual guideline.

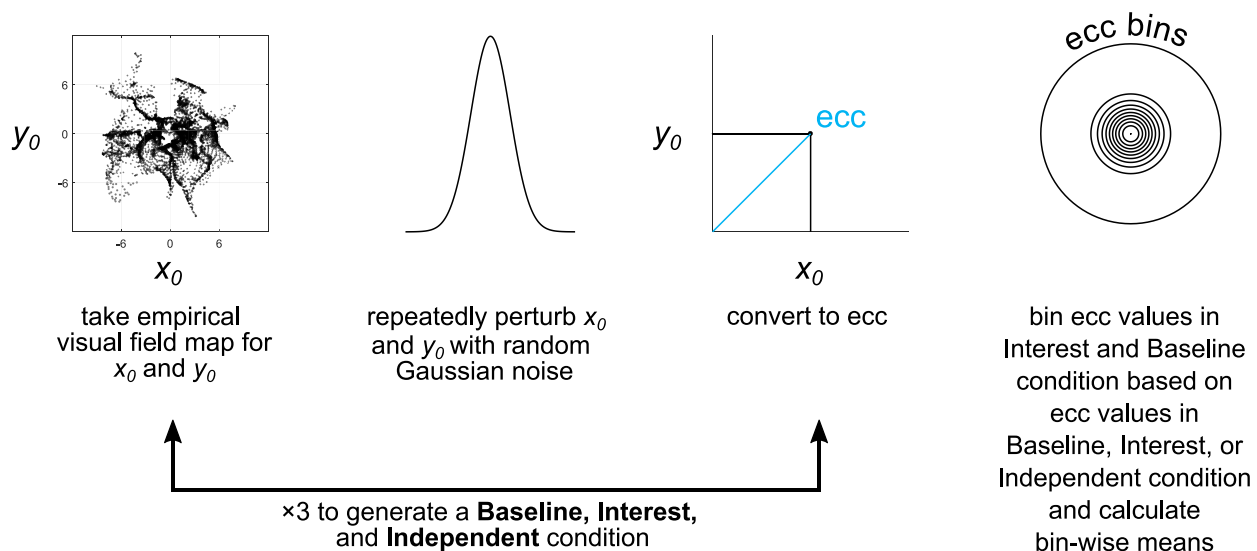


Fig. 3. Schematic workflow of 1D post hoc binning analysis on simulated eccentricity data | Null effect. Ecc = Eccentricity.

peated the binning analysis using equidistant instead of decile binning. To this end, we used a constant bin width of 1.75 dva and an overall binning range of 0 to 19.25 dva eccentricity. Unlike equidistant binning, decile binning ensures a roughly equal number of data points in each bin, which facilitates the interpretation of results. However, we consider equidistant binning as the most common binning type in the pRF literature. For both the change score analysis and equidistant binning, we used the simulation case involving a null effect as a data basis.

2.1.2. 2D post hoc binning analysis on x_0 and y_0

Apart from the 1D binning analysis on eccentricity, we also conducted a 2D binning analysis on the simulated x_0 and y_0 values. To this end, we converted the x_0 and y_0 values to polar coordinates; that is, polar angle and eccentricity (Fig. 2). We then binned the x_0 and y_0 values in the Baseline or Interest condition according to their polar coordinates in the Baseline, Interest, or Independent condition using equidistant bins and calculated the bin-wise x_0 and y_0 means for each condition. The condition-wise means were visualized as vector graphs along with marginal histograms (bin width = 0.5 dva) illustrating the simulated distributions. Vector graphs have been used in prior pRF work (e.g., van Es et al., 2018; Klein et al., 2014; Vo et al., 2017). The 2D binning analysis was performed for all aforementioned simulation cases. The polar angle bins ranged from 0° to 360° with a constant bin width of 45°. The eccentricity bins ranged from 0 to 22 dva (for the simulation case involving a true effect) or from 0 to 20 dva (for all other simulation cases) with a constant bin width of 2 dva.

2.2. Post hoc binning using empirical repeat data

For the post hoc binning analysis on repeat data, we used publicly available pRF estimates from the HCP 7 T Retinotopy Dataset (Benson et al., 2018; 2020). These estimates stem from a split-half analysis where a 2D isotropic Gaussian was fit to fMRI time series from the first and second half of 6 pRF mapping runs. For each half, 6 estimates were obtained for each grayordinate (vertex; <https://wiki.humanconnectome.org/display/WBPublic/Workbench+Glossary>); that is, pRF polar angle, pRF eccentricity, pRF size, pRF gain, percentage of R^2 , and mean signal intensity. The maximal eccentricity of the mapping area subtended 8 dva. For further details, see Benson et al. (2018).

Following Benson et al. (2018), we analyzed complexes of visual areas across hemispheres for the 25th and 75th percentile participants of

the R^2 distribution using delineations from Wang et al.'s (2015) atlas. Benson et al. (2018) generated the R^2 distribution by calculating the median R^2 for each participant across grayordinates from both cortical hemispheres within all areas of Wang et al.'s (2015) atlas. For our purposes, we focused on the posterior complex (V1-V3) and the dorsal complex (V3A/B and IPS05), as those came with a larger number of available data points (which was, amongst other things, necessary to perform the 2D post hoc binning analysis and generate vector graphs).

To obtain x_0 and y_0 values, polar angle and eccentricity estimates were converted to Cartesian coordinates. The eccentricity, x_0 , and y_0 values of the first half were used as a Baseline condition and those of the second half as an Interest condition. Similar to the simulation-based analyses, binning was either based on the Interest or Baseline condition and bin-wise means were calculated. Moreover, binning was either performed without or with condition cross-thresholding. As for the latter case, we removed observations outside a certain eccentricity range (≥ 0 and ≤ 8 dva) or below a certain R^2 cut-off ($\leq 2.2\%$) in the Baseline or Baseline and Interest condition from both conditions. The R^2 cut-off was adopted from Benson et al. (2018).

We then performed a 1D binning analysis on eccentricity and a 2D binning analysis on x_0 and y_0 as we did for the simulated data. However, here, the eccentricity bins for the 2D analysis ranged from 0 to 18 dva with a constant bin width of 2 dva. All binning analyses and visualizations (including those on simulated data) were implemented in Matlab 2016b (9.1; <https://uk.mathworks.com/>) using custom code (Data and code availability). The color scheme used for color-coding was an adapted version of the BrBG palette from ColorBrewer (2.0; Brewer et al., 2021) retrieved via R (3.5.3; R Core Team, 2018) and the package RColorBrewer (1.1-2; Neuwirth, 2014).

3. Results and discussion

3.1. The many faces of regression towards the mean and other problems

To expose the regression artifact, we repeatedly perturbed the x_0 and y_0 values of an empirical visual field map with random Gaussian noise to generate a Baseline and Interest condition. We then converted the x_0 and y_0 values to eccentricity. Subsequently, we binned the eccentricity values of either condition according to eccentricity values in the Baseline condition using deciles and calculated bin-wise means. The bin-wise means from both conditions were plotted against one another along with individual observations per bin and marginal histograms re-

flecting the simulated distributions⁴ (Fig. 4, 1st column). Since there was no true difference between conditions, the bin-wise means should lie on the identity line. Contrary to this prediction, they systematically diverged from the identity line. Strikingly, when using the Interest instead of the Baseline condition for binning, this systematic pattern of divergence flipped (Fig. 4, 2nd column). This bidirectionality is a typical sign of regression towards the mean (Campbell and Kenny, 1999; Shanks, 2017) and due to circularity. This leads to asymmetric bins (see bin-wise ranges of observations for the Baseline and Interest condition, Fig. 4, 1st and 2nd columns) and on average biases bin-wise noise components for the condition that was used for contrasting and binning (henceforth *circular* condition). On the contrary, for the other condition (henceforth *non-circular* condition), this is not the case.

The skew in average noise renders the bin-wise eccentricity means of the circular condition more extreme, especially for lower and higher decile bins. As a result, the bin-wise eccentricity means for the non-circular condition regress – by statistical necessity – to the overall mean⁵ for this condition (red crosshair); that is, they are less extreme. This becomes clear when looking at the different ranges of bin-wise means for the circular and non-circular conditions (Fig. 4, 1st and 2nd columns). If the Interest condition is then contrasted to the Baseline condition, a mean increase in eccentricity for lower deciles and a mean decrease for higher deciles or vice versa occurs, depending on whether the data are binned on the Baseline or Interest condition (Fig. 4, 1st and 2nd columns). This artifact arises because we did not always use independent conditions for binning and contrasting; that is, conditions with independent noise components.

Apart from simple means (e.g., Binda et al., 2013; Carvalho et al., 2020; Haak et al., 2012; Yildirim et al., 2018), post hoc binning analyses have also been performed on change scores in previous pRF studies (e.g., Barton and Brewer, 2015; de Haas et al., 2014; 2020; Prabhakaran et al., 2020; Tsouli et al., 2021). Here, the difference between the Interest and Baseline condition is typically plotted against the binning (i.e., circular) condition (Fig. 5, A., 1st and 2nd columns). Consequently, the bin-wise means now regress to the overall mean of the change score distribution (see also Gignac and Zajenkowski, 2020; Holmes, 2009) and bin-wise noise components are neither unbiased for the change scores nor the binning conditions. This is because the noise components of the change scores are not independent of those in either binning condition. What is more, scatter plots of change scores disguise important aspects readily available with scatter plots of simple scores. Specifically, they prevent us from directly appreciating the larger bin-wise range of eccentricity means for the circular as compared to the non-circular condition (see explanations further above and compare Fig. 5, A., and Fig. 4, 1st and 2nd columns). This makes it difficult to spot the source of the problem graphically when only looking at a single plot. On the other hand, since both the *x*- and *y*-axis feature the Baseline or Interest condition and either of these conditions are used for data binning, the act of double-dipping becomes much more obvious.

Critically, scattering change scores against one of the conditions involved in change score calculation also results in a biased visualization of individual change scores. This is because the noise components of the variables on the *x*- and *y*-axis are not independent, rendering this sorting procedure circular. When plotting individual change scores against the Baseline condition, this results in a downwards sloping data cloud, suggesting an effect although there is none (Fig. 5, A., 1st column). Why

⁴ Note that apart from the visualizations provided here, it might be beneficial to additionally look at Galton squeeze diagrams to detect regression towards or away from the mean (see Fig. 1, Campbell and Kenny, 1999; Galton, 1886; Shanks, 2017).

⁵ Note that for skewed distributions (such as the gamma-like distribution here), the regression effect might be actually towards the mode and away from the mean of the overall distribution (Schwarz and Reike, 2018). If the location of the overall mode and mean are sufficiently close, our visualizations would be unable to distinguish these two cases.

does this happen? Owing to noise, the change scores are more likely to be positive for lower Baseline eccentricities and negative for higher Baseline eccentricities (Fig. 5, A., 1st column). When plotting individual change scores against the Interest condition, the reverse is true (Fig. 5, A., 2nd column). This means visualizing or analyzing the data using such a circular sorting procedure is misleading irrespective of circular data binning (for more details on circular data sorting, see Holmes, 2009; Kriegeskorte et al., 2009).

The fact that circular sorting of change scores and circular data binning are separate issues can be further appreciated by imagining what happens when we plot the individual change scores against the Baseline condition, but bin on the Interest condition (instead of the Baseline condition as before). In this case, the individual change scores are sorted in a way (downwards sloping; just like in Fig. 5, A., 1st column) that is opposite to the trend implied by the bin-wise means (upwards sloping).

How the regression artifact induced by circular data binning manifests can change when data are thresholded across conditions; that is, deleted in a pair- or list-wise fashion (Fig. 5, B. and C., 1st and 2nd columns). In fact, in the event of condition cross-thresholding, noise components are reshaped and might thus not necessarily be unbiased on average even for the non-circular condition (Fig. 5, B., 2nd column as well as Fig. 5, C., 1st and 2nd columns). Condition cross-thresholding is common practice in the pRF literature where data are cleaned across conditions according to eccentricity, goodness-of-fit (R^2), pRF size, missing data or other criteria from one or multiple conditions.

Here, we cross-thresholded the eccentricity values in the Interest and Baseline condition using the eccentricity values from the Baseline condition (Fig. 5, B., 1st and 2nd columns) or both the Baseline and Interest condition (Fig. 5, C., 1st and 2nd columns). This cross-thresholding procedure is circular whenever the noise components of the data used for cross-thresholding are not independent of the noise components of the data involved in contrasting. This is evidently true even without circular data binning. As such, the reason why the noise components in our cross-thresholding scenarios are sometimes biased even for the non-circular condition⁶ (Fig. 5, B., 2nd column as well as Fig. 5, C., 1st and 2nd columns) is because we introduced another layer of circularity.

The fact that circular cross-thresholding and circular data binning are somewhat distinct but also highly similar issues can, for instance, be appreciated when comparing the overall instead of the bin-wise means. Without circular cross-thresholding, the overall mean in both the Baseline and Interest condition amounts to 4.66 dva (Fig. 4, 1st and 2nd columns). With circular cross-thresholding based on the Baseline condition, the overall mean in the Baseline condition amounts to 3.40 dva, whereas it amounts to 3.97 dva in the Interest condition (Fig. 5, B., 1st and 2nd columns). Here, the introduced bias for the Baseline condition can be appreciated by directly comparing the overall means in the Baseline and Interest condition. With circular cross-thresholding based on both the Baseline and Interest condition, the overall means in the Baseline and Interest condition amount to 3.24 dva and 3.25 dva, respectively (Fig. 5, C., 1st and 2nd columns). Here, the introduced bias for the Baseline and Interest condition can be appreciated by comparing the overall means in these conditions to the overall mean of an Independent condition (retest of the Baseline condition) that was cross-thresholded based on both the Baseline and Interest condition. This overall mean amounts to 3.66 dva. We will return to the usefulness of such an Independent condition further below (3.2). In any case, circular cross-thresholding biases the overall means as compared to when no such circular cross-thresholding is performed.

⁶ For reasons of clarity and simplicity, we use the term ‘circular condition’ or ‘non-circular condition’ exclusively when referring to circular data binning. However, other circular selection procedures, such as circular data sorting or cleaning, might of course render a condition circular above and beyond circular data binning.

Simulated null effect

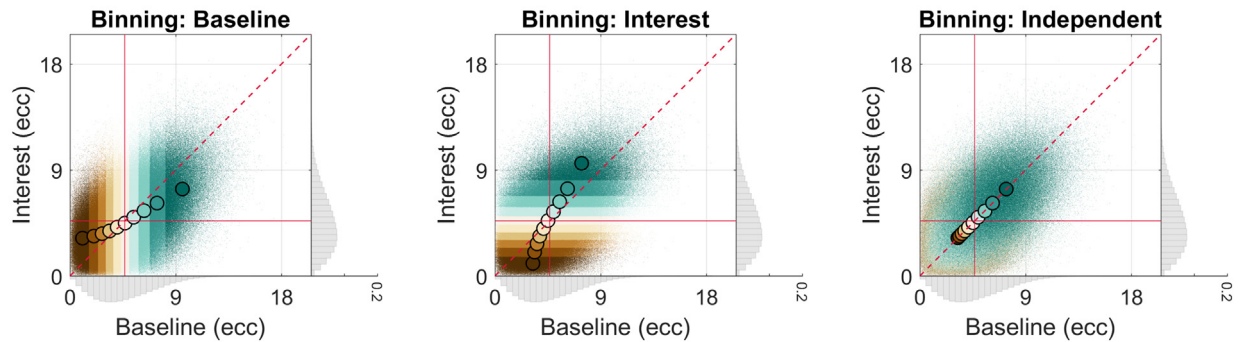


Fig. 4. Simulated 1D post hoc binning analysis on eccentricity | Null effect. Bin-wise eccentricity values and means in the Interest and Baseline condition for a simulated null effect and different data binning scenarios. The eccentricity values in the Baseline and Interest condition were either binned according to eccentricity values in the Baseline (1st column), Interest (2nd column), or an Independent condition (equivalent to repeat data of the Baseline condition; 3rd column). The gray marginal histograms (bin width = 0.5 dva; y-axis: relative frequency) show the simulated eccentricity distributions for each condition, obtained by repeatedly disturbing the x_0 and y_0 values of an empirical visual field map with random Gaussian noise ($SD = 2$ dva) and subsequently converting them to eccentricity values. Note that the range of the marginal y-axis is the same for all histograms. The red crosshair indicates the location of the overall mean for the Interest and Baseline condition. The red dashed line corresponds to the identity line. Dark brown colors correspond to lower and dark blue-green colors to higher decile bins. The smaller colorful dots represent individual data points and the larger colorful dots with the black outline bin-wise means. The maximal eccentricity of the stimulated visual field area subtended 8.5 dva. Dva = Degrees of visual angle. Ecc = Eccentricity.

Importantly, however, only circular cross-thresholding based on the Baseline condition results in artificial differences between the overall means. Why is this? Given that the level of noise in the Interest and Baseline condition was equivalent (2.1 Post hoc binning using simulated data), circular cross-thresholding based on both the Baseline and Interest condition on average skewed the noise components for these conditions similarly, resulting in biased overall means, but a valid difference of around 0 between them. However, as for empirical data, the assumption of equivalent noise levels can probably only be safely made for repeat data (and even then, this needs to be justifiable). In any case, conceptually, circular cross-thresholding without data binning can be regarded as a single bin or region-of-interest analysis (Kriegeskorte et al., 2009), essentially constituting another subtype of a post hoc binning analysis.

The appearance of the regression artifact arising from circular data binning can furthermore change when the level of noise depends on eccentricity – a property better known as *heteroskedasticity* (Fig. 6, A., 1st and 2nd columns; see also Holmes, 2009). In fact, the case of eccentricity-dependent noise shows that the artifact can include some clear regression away from the mean – a phenomenon referred to as *egression*⁷ (Fig. 6, A., 1st and 2nd columns; see e.g., Campbell and Kenny, 1999; Schwarz and Reike, 2018). Eccentricity-dependent noise might arise from fitting errors that differ across visual space. This could be due to partial stimulation of pRFs (especially near the edge of the stimulated mapping area), higher variability in pRF position estimates for wider pRFs as well as fluctuations in the signal-to-noise ratio of brain responses from the central to the peripheral visual field or as a result of manipulating attention.

The regression artifact due to circular data binning also manifested when simulating a true effect (Fig. 6, B., 1st and 2nd columns). The same was true for equidistant binning (Fig. 6, C., 1st and 2nd columns), which is frequently applied in the pRF literature. However, unlike decile binning (which we used further above), equidistant binning resulted in a lower number of observations for higher equidistant bins (due to the gamma-like eccentricity distribution; Fig. 6, C., 1st and 2nd columns). Consequently, for higher equidistant bins, the skew in average noise for the circular condition was generally larger here (compare Fig. 6, C.,

and Fig. 4, 1st and 2nd columns). Similarly, for higher equidistant bins, noise components were not always unskewed on average for the non-circular condition (see Fig. 6, C., 1st and 2nd columns, where the pattern of bin-wise means is not entirely mirror-symmetric). This is because for random noise to be unskewed on average, the number of observations needs to be sufficiently large.

Critically, both true effects and equidistant binning can substantially modify the appearance of the regression artifact. Along with circular condition cross-thresholding and eccentricity-dependent noise, this teaches us an important lesson: the regression artifact can take pretty much *any* form.⁸

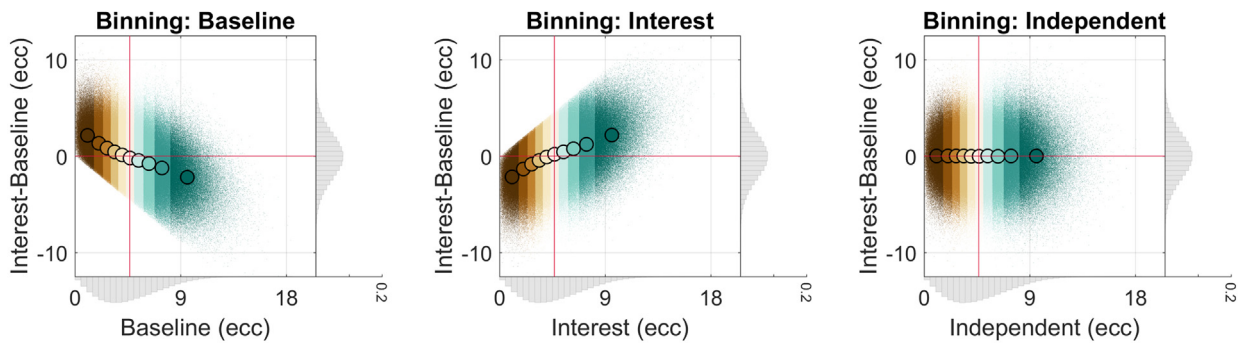
For all presented simulation cases (null effect, null effect with cross-thresholding or eccentricity-dependent noise, and true effect), the regression artifact likewise manifested for another kind of binning analysis, namely, when binning the x_0 and y_0 values according to both eccentricity and polar angle (i.e., 2D segments) and computing shift vectors (Fig. 2 as well as Fig. 7 and Fig. S2-S5, 1st row). Here, the bin-wise means regressed towards and away from the overall means of the x_0 and y_0 distribution. The calculation of shift vectors is not uncommon in pRF studies (e.g., van Es et al., 2018; Klein et al., 2014; Vo et al., 2017).

Notably, for empirical repeat data from the HCP (Benson et al., 2018; 2020), both kinds of binning analyses produced patterns consistent with the regression artifact (Fig. S6-S13). This establishes its practical relevance. Moreover, some of us recently retracted an article on attention-induced differences in pRF position and size in V1-V3 (de Haas et al., 2014) because an in-house reanalysis suggested that circular post hoc binning along with circular condition cross-thresholding and heteroskedasticity yielded artificial results in the form of egression from the mean (de Haas et al., 2020). In this case, the apparent significant effect was an increase in eccentricity and pRF size in the Interest vs Baseline condition (expressed as change scores) for eccentricity bins (based on the Baseline condition) in the middle of the tested range. Importantly, the inferential statistical analysis in this study (de Haas et al., 2014; 2020) was based on unbinned data, and thus the overall means.

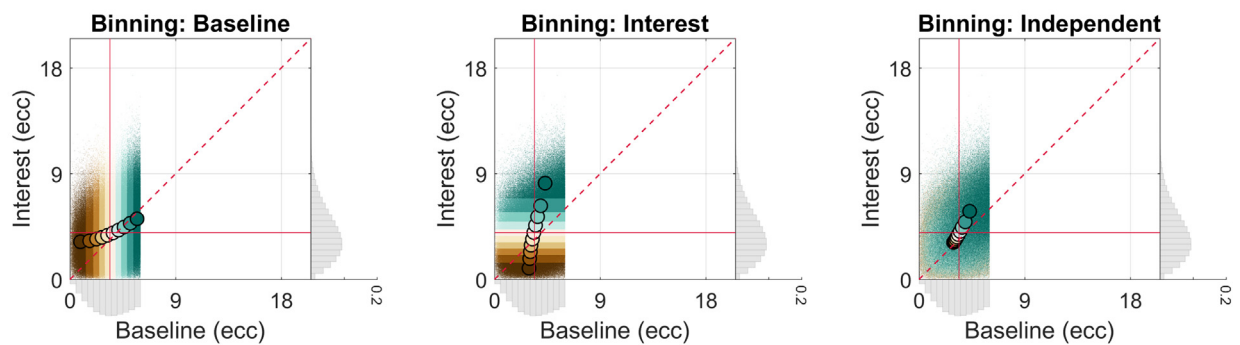
⁷ Note that the regression was presumably towards the nearest modes of the simulated bimodal distribution (see marginal histograms in Fig. 6, A., 1st and 2nd columns; Schwarz and Reike, 2018).

⁸ Note that floor/ceiling effects (due to physiological and methodological constraints on the minimum and maximum observable value) and/or the calculation of absolute (raw) vs proportional (%) differences are further factors influencing the appearance of the regression artifact (de Haas et al., 2014; 2020; Holmes, 2009).

A. Simulated null effect - Change score



B. Simulated null effect - Cross-thresholding (Baseline)



C. Simulated null effect - Cross-thresholding (Baseline and Interest)

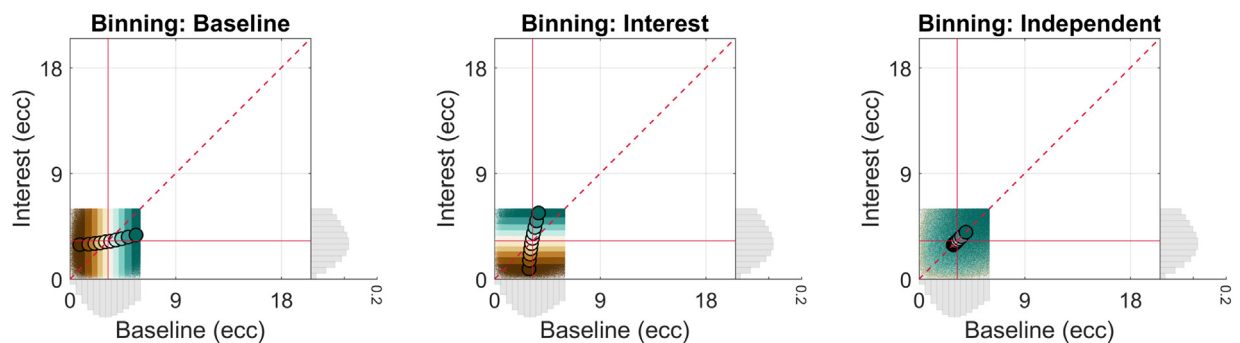


Fig. 5. Simulated 1D post hoc binning analysis on eccentricity | Null effect – Change score and cross-thresholding. A. The same as in Fig. 4, although here, the change score (Interest vs Baseline) is plotted against the respective binning condition. B. The same as in Fig. 4, although here, condition cross-thresholding was applied, i.e., simulated observations falling outside a certain eccentricity range (≥ 0 and ≤ 6 dva) in the Baseline condition were removed from all conditions. C. The same as in B., although here, condition cross-thresholding was based on both the Baseline and Interest condition. (Condition) cross-thresholding = The pair-wise or list-wise deletion of observations across conditions.

As such, the apparent significant effect was likely driven by or inflated due to circular cross-thresholding.

The example of de Haas et al. (2014, 2020) illustrates that data visualizations and associated inferential statistical analyses do not necessarily suffer from the same pitfalls. It is also possible that only one but not the other produces artifactual changes. This potential divergence adds another layer of complexity to the issues we discussed here.

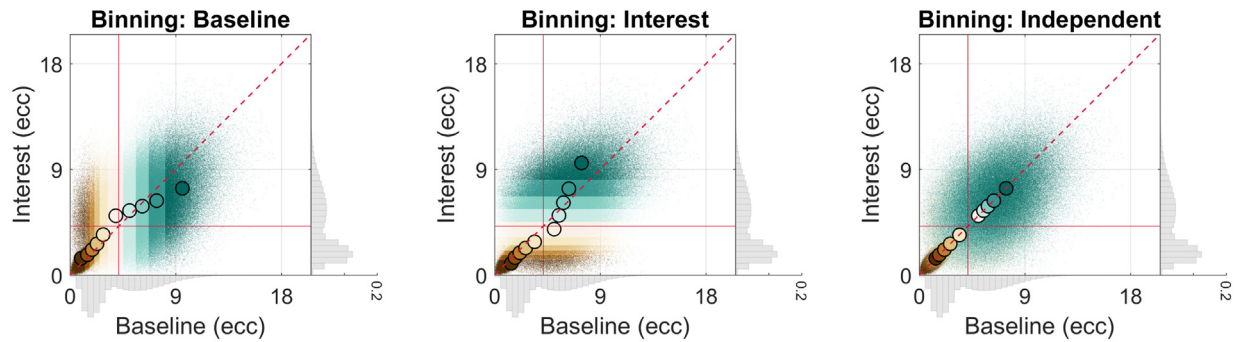
Taken together, the heterogeneity in manifestation we exposed here makes it hard to spot the regression artifact by visual inspection alone and highlights its dependency on the type of analysis, additional circular selection practices as well as exact distributional properties of the data at hand (see Campbell and Kenny, 1999; Holmes, 2009; Schwarz and Reike, 2018, for similar points). Importantly, circular data binning is

only but one pitfall resulting in artifactual changes. Other pitfalls, such as circular sorting of change scores and circular cross-thresholding are equally problematic.

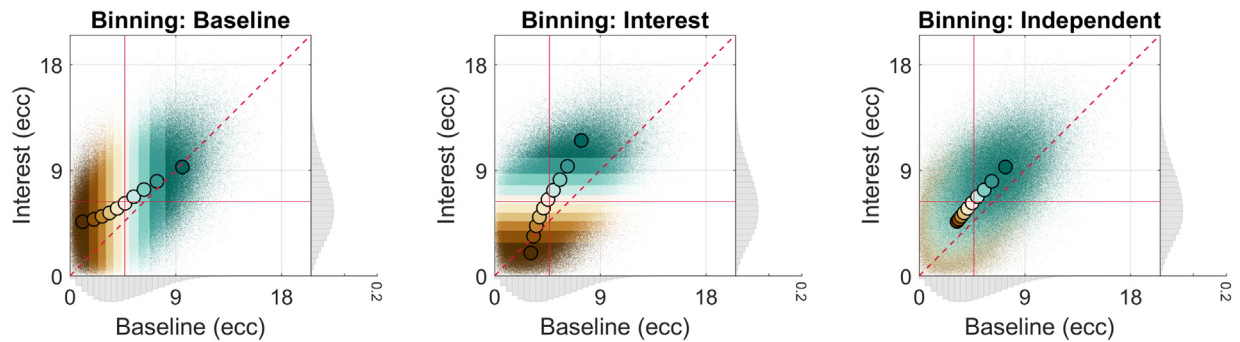
3.2. Potential mitigation strategies

How can we omit double-dipping and control for regression towards the mean? We could, for instance, use an Independent condition for binning (such as repeat data or odd or even runs for the Baseline condition; Fig. 4 and Fig. 5–6, A.-C., 3rd column as well as Fig. 7 and Fig. S2-S5, 2nd row) or an anatomical criterion (Kriegeskorte et al., 2009), such as cortical distance or anatomical atlases (Benson et al., 2014; 2012).

A. Simulated null effect - Eccentricity-dependent noise



B. Simulated true effect - Radial shift



C. Simulated null effect - Equidistant binning

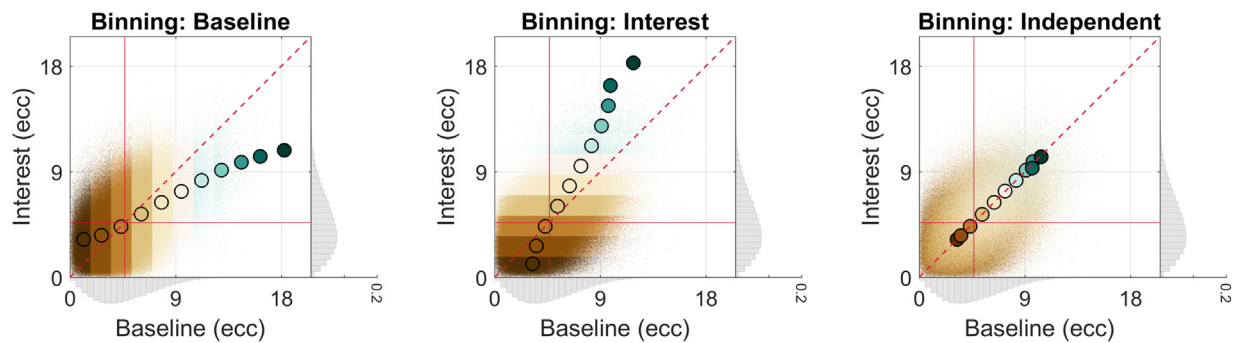


Fig. 6. Simulated 1D post hoc binning analysis on eccentricity | Null or true effect – Eccentricity-dependent noise, radial shift, and equidistant binning. **A.** The same as in Fig. 4, although here, original observations having smaller eccentricities (≥ 0 and < 3 dva) were disturbed by random Gaussian noise with a smaller standard deviation ($SD = 0.25$ dva) and those having larger eccentricities (≥ 3 dva) by random Gaussian noise with a larger standard deviation ($SD = 2$ dva). **B.** The same as in Fig. 4, although here, we simulated a true effect; that is, a radial increase in eccentricity of 2 dva in the Interest as compared to the Baseline condition. **C.** The same as in Fig. 4, although here, equidistant binning was used. The equidistant bins ranged from an eccentricity of 0 dva to an eccentricity of 19.25 dva with a constant bin-width of 1.75 dva. Please note the different number of bins here relative to the other figure panels (11 vs 10, respectively).

This way, noise components should be unbiased on average in both the Baseline and Interest condition.

Unbiased bin-wise noise components are of course less likely for sparsely populated bins (Fig. 6, C., 3rd column as well as Fig. 7 and Fig. S2-S5, 2nd row), which can be captured by quantifying uncertainty. Critically, however, for scatter plots of change scores, bin-wise noise components are not unbiased for the Independent binning condition (Fig. 5, A., 3rd column). The reason for this is the same as before: non-independence of noise components. Thus, only the bin-wise change scores can be readily interpreted here. Moreover, given that cross-thresholding reshapes noise components, they might not be unbiased when binning on an Independent condition (Fig. 5, B. and C.,

3rd column as well as Fig. 5, C., 2nd row). The same can evidently also happen with an anatomical criterion if the Baseline and/or the Interest condition are subjected to cross-thresholding. Consequently, unless cross-thresholding can be omitted or demonstrated to be unbiased (see below for further considerations), binning on an Independent condition might not be a safe option.

Of note, for the discussed cross-thresholding case where circular cross-thresholding was performed based on both the Interest and Baseline condition, binning on the Independent condition ensured that the bin-wise noise components for the Interest and Baseline condition are similarly biased (Fig. 5, C., 3rd column). As mentioned earlier, this is because cross-thresholding of this sort biases the noise components in the

Simulated null effect

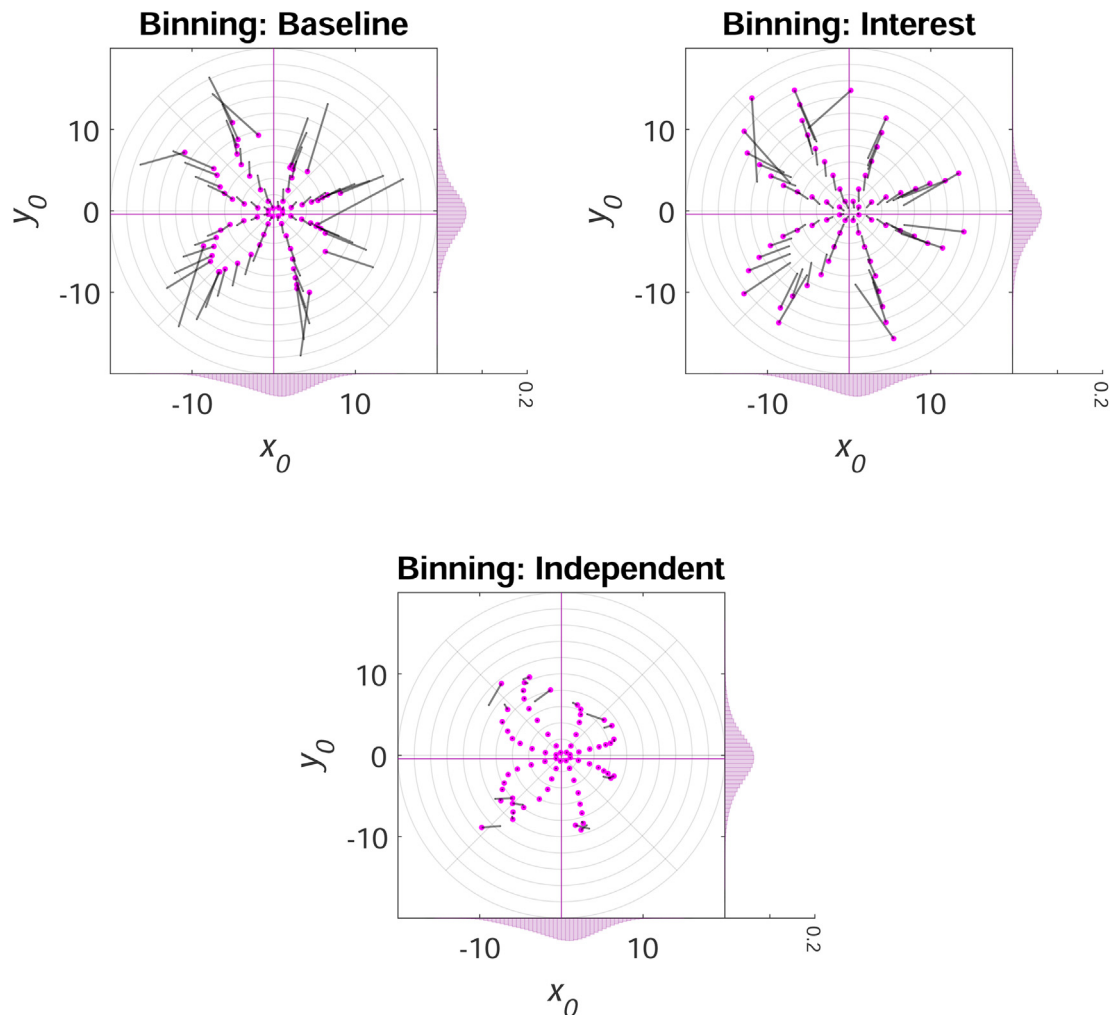


Fig. 7. Simulated 2D post hoc binning analysis on x_0 and y_0 | Null effect. Bin-wise x_0 and y_0 means in the Interest and Baseline condition for a simulated null effect and different data binning scenarios. The x_0 and y_0 values in the Baseline and Interest condition were either binned according to eccentricity and polar angle values in the Baseline (1st column, 1st row), Interest (2nd column, 1st row), or an Independent condition (equivalent to repeat data of the Baseline condition; 2nd row). The marginal histograms (bin width = 0.5 dva; y-axis: relative frequency) show the simulated x_0 and y_0 distributions for each condition, obtained by repeatedly disturbing the x_0 and y_0 values of an empirical visual field map with random Gaussian noise ($SD=2$ dva). Magenta histograms correspond to the Interest condition and gray histograms to the Baseline condition. Note that the range of the marginal y-axis is the same for all histograms. The large magenta dots (arrow tip) correspond to the means in the Interest condition and the endpoint of the gray line (arrow knock) to the means in the Baseline condition. The gray line itself (arrow shaft) depicts the shift from the Baseline to the Interest condition. The magenta crosshair indicates the location of the overall x_0 and y_0 means for the Interest condition and the gray crosshair the location of the overall means for the Baseline condition. Note that if there is no systematic difference between the Baseline and Interest condition, the histograms and crosshairs coincide (as is the case here). The light gray polar grid demarcates the bin segments. Polar angle bins ranged from 0° to 360° with a constant bin width of 45° and eccentricity bins from 0 to 20 dva with a constant bin width of 2 dva. The maximal eccentricity of the stimulated visual field area subtended 8.5 dva. Dva = Degrees of visual angle.

Baseline and Interest condition similarly (3.1.) and binning on an Independent condition introduces no further biases. Moreover, given that the noise components of both the Interest and Baseline condition were independent of those in the Independent condition, cross-thresholding did not bias the noise components in the Independent condition. As such, although the simple bin-wise means in the Baseline and Interest condition are biased, the difference between those amounts to around 0 (Fig. 5, C., 3rd column).

Apart from binning on an Independent condition, we could use analyses without binning that control for circularity and regression artifacts or effects could be evaluated against appropriate null distributions that take into account all statistical dependencies (e.g., Holmes, 2009; Kriegeskorte et al., 2009). For instance, errors-in-variables models (e.g.,

Deming regression) might be an option. Such models account for the noise in both the Baseline and Interest condition as well as for the fact that we often have no clear separation between independent and dependent variables in post hoc analyses of pRF data. However, as with any statistical approach, the underlying assumptions need to be checked carefully.

Just like circular data binning, circular sorting of change scores can be counteracted by plotting individual change scores against an Independent condition (Fig. 5, A., 3rd column). Similarly, one way to deal with circular cross-thresholding might be to cross-threshold all data according to an Independent condition/the Independent binning condition. However, condition-specific systematic errors, such as artifacts and outliers, might survive such independent data cleaning. As such, the us-

age of robust estimators might be advisable. Future research is necessary to evaluate this point more comprehensively.

A combination of the discussed approaches might prove most fruitful. Regardless of the specific mitigation strategy, we believe that in light of the many layers of complexity in our analysis pipelines, we need to make it common practice to perform sanity checks using (null) simulations and empirical repeat data. This is because such sanity checks provide a means for us researchers to ensure the validity of our analysis procedures.

3.3. The bigger picture

Circular post hoc binning analyses come in many flavors (e.g., centroids, shift vectors, eccentricity differences, x_0 and y_0 differences, and 1D or 2D bins) and cannot be assumed to be restricted to pRF position estimates. For instance, partial stimulation of pRFs likely results in heteroskedasticity and positively correlated errors for pRF size and position. This would, for instance, bias bin-wise pRF size vs pRF position or pRF size vs pRF size comparisons where binning is based on non-independent eccentricity values. Likewise, fitting errors due to partial stimulation should be more pronounced whenever pRF size is larger, leading to stronger artifactual effects (for simulations using different levels of noise, see Holmes, 2009). The same is to be expected based on a higher variability in pRF position estimates for wider pRFs. These factors might potentially explain why changes in pRF position and/or size have been reported to be tendentially larger in higher-level areas where pRFs are wider (e.g., Barton and Brewer, 2015; van Es et al., 2018; de Haas et al., 2014; 2020; Klein et al., 2014).

Moreover, the distribution of errors likely depends on the toolbox that was used for fitting (Lerma-Usabiaga et al., 2020), making it hard to generalize across studies. And lastly, delineations of visual areas in post hoc binning analyses should ideally also be based upon independent criteria as this is where selection starts. Importantly, the intricacies we just discussed do not only apply to circular data binning, but also circular sorting of change scores and circular condition cross-thresholding.

The application of circular data binning, circular sorting of change scores, and/or circular cross-thresholding in the pRF literature might have led to spurious claims about changes in pRFs (see de Haas et al., 2014; 2020, for an example). Consequently, we encourage researchers who used such procedures to check for the severity of biases in their analyses by running adequate simulations and reanalyzing the original data wherever possible. Likewise, we urge them to take into account the issues discussed here when conducting future studies, reviewing manuscripts, and when teaching and mentoring.

3.4. Limitations

Our simulations were designed to encapsulate a given issue succinctly and cannot be interpreted as reflecting the exact properties of empirical pRF data. For this, we would need to have a good understanding of the underlying noise components. Similarly, the level of random Gaussian noise we adopted for most simulations ($SD = 2$ dva) might be more reminiscent of higher than lower visual areas (although this depends on many factors, such as mapping stimulus and magnetic field strength). For the present purposes, it appeared important to settle on a level allowing for clear exposition. Moreover, as alluded to further above (1. Introduction), unless there is a perfect correlation between two variables (and thus no random noise), double-dipping and regression towards or away from the mean likely pose issues to post hoc analyses.

To fully parallel our simulations, the analyses of the HCP data would have benefited from binning on an Independent condition; that is, a second set of repeat data. pRF estimates for such an Independent condition are currently not publicly available (Benson et al., 2018; 2020), leaving this sanity check for future research. Moreover, unlike our simulations,

the condition cross-thresholding applied to the HCP data not only involved pRF position, but also goodness-of-fit (2.1 Post hoc binning using simulated data and 2.2 Post hoc binning using empirical repeat data). This is because such multivariate data cleaning is frequently applied in pRF studies. It is challenging to simulate these more complex scenarios and thus best addressed in a separate article.

Some post hoc binning analyses in the pRF literature are conducted in a hemifield-specific fashion, whereas others mirror observations across hemifields or quadrants. Our analyses do not capture these specificities. However, there is no reason to believe that they would alleviate the expression of the regression artifact. The primary component that might change when applying such procedures is the location of the overall mean and the shape of the data distribution and thus how exactly the artifact manifests (for preliminary analyses, see Stoll et al., 2022). Of course, if data points are not mirrored based on an Independent condition but, for instance, the Baseline condition, data mirroring in combination with circular data binning and/or circular cross-thresholding might favor noise components in multiple ways. Importantly, circular data mirroring is also problematic for analyses that do not involve any circular data binning and/or circular cross-thresholding, as are other procedures, such as circular data weighting (Kriegeskorte et al., 2009).

4. Conclusions

Without doubt, circularity and regression towards the mean are thorny and omnipresent problems that can manifest subtly and diversely (e.g., Ball et al., 2020; Barnett et al., 2005; Campbell and Kenny, 1999; Eriksson and Häggström, 2014; Gignac and Zajenkowski, 2020; Holmes, 2009; Kilner, 2013; Kriegeskorte et al., 2009; Preacher et al., 2005; Shanks, 2017; Stigler, 1997; Vul et al., 2009). As such, we need to ensure that the validation of analysis procedures becomes part and parcel of the scientific process.

Data and code availability

Preprocessed data, custom code, and figures are available at <https://doi.org/10.17605/OSF.IO/WJADP>.

Declaration of competing interest

The authors declare no conflict of interest.

Acknowledgements

This research was supported by a European Research Council Starting Grant to DSS (WMOSPOTWU, 310829). BdH was supported by a European Research Council Starting Grant (INDIVISUAL, 852885) and the Deutsche Forschungsgemeinschaft (222641018SFB/TRR 135 TP C9). We thank three peer reviewers for providing constructive feedback.

Supplementary material

Supplementary material associated with this article can be found in the online version at doi:10.1016/j.neuroimage.2022.119557.

References

- Alvarez, I., de Haas, B., Clark, C.A., Rees, G., Schwarzkopf, D.S., 2015. Comparing different stimulus configurations for population receptive field mapping in human fMRI. *Front. Hum. Neurosci.* 9, 96. doi:10.3389/fnhum.2015.00096.
- Ball, T.M., Squeglia, L.M., Tapert, S.F., Paulus, M.P., 2020. Double dipping in machine learning: Problems and solutions. *Biol. Psychiatry Cogn. Neurosci. Neuroimaging* 5 (3), 261–263. doi:10.1016/j.bpsc.2019.09.003.
- Barnett, A.G., van der Pols, J.C., Dobson, A.J., 2005. Regression to the mean: What it is and how to deal with it. *Int. J. Epidemiol.* 34 (1), 215–220. doi:10.1093/ije/dyh299.
- Barton, B., Brewer, A.A., 2015. fMRI of the rod scotoma elucidates cortical rod pathways and implications for lesion measurements. *Proc. Natl. Acad. Sci. U.S.A.* 112 (16), 5201–5206. doi:10.1073/pnas.1423673112.

- Benson, N.C., Butt, O.H., Brainard, D.H., Aguirre, G.K., 2014. Correction of distortion in flattened representations of the cortical surface allows prediction of V1-V3 functional organization from anatomy. *PLoS Comput. Biol.* 10 (3), e1003538. doi:10.1371/journal.pcbi.1003538.
- Benson, N.C., Butt, O.H., Datta, R., Radoeva, P.D., Brainard, D.H., Aguirre, G.K., 2012. The retinotopic organization of striate cortex is well predicted by surface topology. *Curr. Biol.* 22 (21), 2081–2085. doi:10.1016/j.cub.2012.09.014.
- Benson, N.C., Jamison, K.W., Arcaro, M.J., Vu, A.T., Glasser, M.F., Coalson, T.S., Van Essen, D.C., Yacoub, E., Ugurbil, K., Winawer, J., Kay, K., 2018. The Human Connectome Project 7 Tesla retinotopy dataset: Description and population receptive field analysis. *J. Vis.* 18 (13), 1–22. doi:10.1167/18.13.23.
- Benson, N.C., Jamison, K.W., Arcaro, M.J., Vu, A.T., Glasser, M.F., Coalson, T.S., Van Essen, D.C., Yacoub, E., Ugurbil, K., Winawer, J., Kay, K., 2020. The HCP 7T Retinotopy Dataset. <https://doi.org/10.17605/OSF.IO/BW9EC>.
- Binda, P., Thomas, J.M., Boynton, G.M., Fine, I., 2013. Minimizing biases in estimating the reorganization of human visual areas with bold retinotopic mapping. *J. Vis.* 13 (7), 1–16. doi:10.1167/13.7.13.
- Brewer, C.A., Harrower, M., University, T.P.S., 2021. ColorBrewer [Web tool]. <https://colorbrewer2.org>.
- Campbell, D.T., Kenny, D.A., 1999. *A primer on regression artifacts*. Guilford Press, New York, NY.
- Carvalho, J., Invernizzi, A., Ahmadi, K., Hoffmann, M.B., Renken, R.J., Cornelissen, F.W., 2020. Micro-probing enables fine-grained mapping of neuronal populations using fmri. *Neuroimage* 209, 116423. doi:10.1016/j.neuroimage.2019.116423.
- Cohen, J., Cohen, P., West, S.G., Aiken, L.S., 2003. *Applied multiple regression/correlation analysis for the behavioral sciences*. Lawrence Erlbaum Associates, Mahwah, NJ.
- van Dijk, J.A., de Haas, B., Moutsiana, C., Schwarzkopf, D.S., 2016. Intersession reliability of population receptive field estimates. *Neuroimage* 143, 293–303. doi:10.1016/j.neuroimage.2016.09.013.
- Dumoulin, S.O., Knapen, T., 2018. How visual cortical organization is altered by ophthalmologic and neurologic disorders. *Annu. Rev. Vis. Sci.* 4 (1), 357–379. doi:10.1146/annurev-vision-091517-033948.
- Dumoulin, S.O., Wandell, B.A., 2008. Population receptive field estimates in human visual cortex. *Neuroimage* 39 (2), 647–660. doi:10.1016/j.neuroimage.2007.09.034.
- Eriksson, K., Häggström, O., 2014. Lord's paradox in a continuous setting and a regression artifact in numerical cognition research. *PLoS ONE* 9 (4), e95949. doi:10.1371/journal.pone.0095949.
- van Es, D.M., Theeuwes, J., Knapen, T., 2018. Spatial sampling in human visual cortex is modulated by both spatial and feature-based attention. *Elife* 7, e36928. doi:10.7554/eLife.36928.
- Galton, F., 1886. Regression towards mediocrity in hereditary stature. *J. Anthropol. Inst. Gt. Britain Irel.* 1, 246–263. <http://www.jstor.org/stable/2841583>.
- Gignac, G.E., Zajenkowski, M., 2020. The Dunning-Kruger effect is (mostly) a statistical artefact: Valid approaches to testing the hypothesis with individual differences data. *Intelligence* 80, 101449. doi:10.1016/j.intell.2020.101449.
- Haak, K.V., Cornelissen, F.W., Morland, A.B., 2012. Population receptive field dynamics in human visual cortex. *PLoS ONE* 7 (5), e37686. doi:10.1371/journal.pone.0037686.
- de Haas, B., Schwarzkopf, D.S., Anderson, E.J., Rees, G., 2014. Perceptual load affects spatial tuning of neuronal populations in human early visual cortex. *Curr. Biol.* 24 (2), R66–R67. doi:10.1016/j.cub.2013.11.061. [Retracted, 2020, *Curr. Biol.* 30, 4814].
- de Haas, B., Schwarzkopf, D.S., Anderson, E.J., Rees, G., 2020. Retraction notice to: Perceptual load affects spatial tuning of neuronal populations in human early visual cortex. *Curr. Biol.* 30 (23), 4814. doi:10.1016/j.cub.2020.11.015.
- Holmes, N.P., 2009. The principle of inverse effectiveness in multisensory integration: Some statistical considerations. *Brain Topogr.* 21 (3–4), 168–176. doi:10.1007/s10548-009-0097-2.
- Infanti, E., Schwarzkopf, D.S., 2020. Mapping sequences can bias population receptive field estimates. *Neuroimage* 211, 116636. doi:10.1016/j.neuroimage.2020.116636.
- Kay, K.N., Winawer, J., Mezer, A., Wandell, B.A., 2013. Compressive spatial summation in human visual cortex. *J. Neurophysiol.* 110 (2), 481–494. doi:10.1152/jn.00105.2013.
- Kenny, D.A., 2005. *Regression artifacts*. In: B. S. Everitt & D. C. Howell (Eds.), *Encycl. Stat. Behav. Sci.*. Chichester, UK, pp. 1723–1725.
- Kilner, J.M., 2013. Bias in a common EEG and MEG statistical analysis and how to avoid it. *Clin. Neurophysiol.* 124 (10), 2062–2063. doi:10.1016/j.clinph.2013.03.024.
- Klein, B.P., Harvey, B.M., Dumoulin, S.O., 2014. Attraction of position preference by spatial attention throughout human visual cortex. *Neuron* 84 (1), 227–237. doi:10.1016/j.neuron.2014.08.047.
- Kriegeskorte, N., Simmons, W.K., Bellgowan, P.S., Baker, C.I., 2009. Circular analysis in systems neuroscience: The dangers of double dipping. *Nat. Neurosci.* 12 (5), 535–540. doi:10.1038/nn.2303.
- Le, R., Witthoft, N., Ben-Shachar, M., Wandell, B., 2017. The field of view available to the ventral occipito-temporal reading circuitry. *J. Vis.* 17 (4), 1–19. doi:10.1167/17.4.6.
- Lerma-Usabiaga, G., Benson, N., Winawer, J., Wandell, B.A., 2020. A validation framework for neuroimaging software: The case of population receptive fields. *PLoS Comput. Biol.* 16 (6), e1007924. doi:10.1371/journal.pcbi.1007924.
- Makin, T.R., De Xivry, J.J.O., 2019. Ten common statistical mistakes to watch out for when writing or reviewing a manuscript. *Elife* 8, e48175. doi:10.7554/eLife.48175.
- Morgan, C., Schwarzkopf, D.S., 2020. Comparison of human population receptive field estimates between scanners and the effect of temporal filtering. *F1000Research* 8, 1681. doi:10.12688/f1000research.20496.2.
- Neuwirth, E., 2014. RColorBrewer: ColorBrewer palettes [Computer software]. <https://cran.r-project.org/package=RColorBrewer>.
- Prabhakaran, G.T., Carvalho, J., Invernizzi, A., Kanowski, M., Renken, R.J., Cornelissen, F.W., Hoffmann, M.B., 2020. Foveal pRF properties in the visual cortex depend on the extent of stimulated visual field. *Neuroimage* 222, 117250. doi:10.1016/j.neuroimage.2020.117250.
- Preacher, K.J., MacCallum, R.C., Rucker, D.D., Nicewander, W.A., 2005. Use of the extreme groups approach: A critical reexamination and new recommendations. *Psychol. Methods* 10 (2), 178–192. doi:10.1037/1082-989X.10.2.178.
- R Core Team, 2018. R: A language and environment for statistical computing [Computer software].
- Schellekens, W., Petridou, N., Ramsey, N.F., 2018. Detailed somatotopy in primary motor and somatosensory cortex revealed by Gaussian population receptive fields. *Neuroimage* 179, 337–347. doi:10.1016/j.neuroimage.2018.06.062.
- Schwarz, W., Reike, D., 2018. Regression away from the mean: Theory and examples. *Br. J. Math. Stat. Psychol.* 71 (1), 186–203. doi:10.1111/bmsp.12106.
- Senden, M., Reithler, J., Gijzen, S., Goebel, R., 2014. Evaluating population receptive field estimation frameworks in terms of robustness and reproducibility. *PLoS ONE* 9 (12), e114054. doi:10.1371/journal.pone.0114054.
- Shanks, D.R., 2017. Regressive research: The pitfalls of post hoc data selection in the study of unconscious mental processes. *Psychon. Bull. Rev.* 24 (3), 752–775. doi:10.3758/s13423-016-1170-y.
- Stigler, S.M., 1997. Regression towards the mean, historically considered. *Stat. Methods Med. Res.* 6 (2), 103–114. doi:10.1177/096228029700600202.
- Stoll, S., Infanti, E., Schwarzkopf, D.S., 2022. The impact of multifocal attention on population receptive fields in human visual cortex - A tale of unexpected complexities <https://doi.org/10.17605/OSF.IO/G4HD2>.
- Tsouli, A., Cai, Y., van Ackooij, M., Hofstetter, S., Harvey, B.M., te Pas, S.F., van der Smagt, M.J., Dumoulin, S.O., 2021. Adaptation to visual numerosity changes neural numerosity selectivity. *Neuroimage* 229, 117794. doi:10.1016/j.neuroimage.2021.117794.
- Vo, V.A., Sprague, T.C., Serences, J.T., 2017. Spatial tuning shifts increase the discriminability and fidelity of population codes in visual cortex. *J. Neurosci.* 37 (12), 3386–3401. doi:10.1523/JNEUROSCI.3484-16.2017.
- Vul, E., Harris, C., Winkielman, P., Pashler, H., 2009. Puzzlingly high correlations in fMRI studies of emotion, personality, and social cognition. *Perspect. Psychol. Sci.* 4 (3), 319–324. doi:10.1111/j.1745-6924.2009.01132.x.
- Wang, L., Mruczek, R.E., Arcaro, M.J., Kastner, S., 2015. Probabilistic maps of visual topography in human cortex. *Cereb. Cortex* 25 (10), 3911–3931. doi:10.1093/cercor/bhu277.
- Yildirim, F., Carvalho, J., Cornelissen, F.W., 2018. A second-order orientation-contrast stimulus for population-receptive-field-based retinotopic mapping. *Neuroimage* 164, 183–193. doi:10.1016/j.neuroimage.2017.06.073.

Further reading

- Amano, K., Wandell, B.A., Dumoulin, S.O., 2009. Visual field maps, population receptive field sizes, and visual field coverage in the human MT+ complex. *J. Neurophysiol.* 102 (5), 2704–2718. doi:10.1152/jn.00102.2009.
- Brainard, D.H., 1997. The psychophysics toolbox. *Spat. Vis.* 10 (4), 433–436. doi:10.1163/156856897X00357.
- Breuer, F.A., Blaimer, M., Heidemann, R.M., Mueller, M.F., Griswold, M.A., Jakob, P.M., 2005. Controlled aliasing in parallel imaging results in higher acceleration (CAIPIRINHA) for multi-slice imaging. *Magn. Reson. Med.* 53 (3), 684–691. doi:10.1002/mrm.20401.
- Dale, A.M., Fischl, B., Sereno, M.I., 1999. Cortical surface-based analysis: I. Segmentation and surface reconstruction. *Neuroimage* 9 (2), 179–194. doi:10.1006/nimg.1998.0395.
- Engel, S.A., Glover, G.H., Wandell, B.A., 1997. Retinotopic organization in human visual cortex and the spatial precision of functional MRI. *Cereb. Cortex* 7 (2), 181–192. doi:10.1093/cercor/7.2.181.
- Fischl, B., Sereno, M.I., Dale, A.M., 1999. Cortical surface-based analysis: II. Inflation, flattening, and a surface-based coordinate system. *Neuroimage* 9 (2), 195–207. doi:10.1006/nimg.1998.0396.
- Kleiner, M., Brainard, D.H., Pelli, D.G., Broussard, C., Wolf, T., Niehorster, D., 2007. What's new in psychtoolbox-3? [ECVP abstract supplement]. *Perception* 36. <http://journals.sagepub.com/doi/10.1177/03010066070360S101>
- Lagarias, J.C., Reeds, J.A., Wright, M.H., Wright, P.E., 1998. Convergence properties of the Nelder-Mead simplex method in low dimensions. *SIAM J. Optim.* 9 (1), 112–147. doi:10.1137/S1052623496303470.
- Nelder, J.A., Mead, R., 1965. A simplex method for function minimization. *Comput. J.* 7 (4), 308–313. doi:10.1093/comjnl/7.4.308.
- Pelli, D.G., 1997. The VideoToolbox software for visual psychophysics: transforming numbers into movies. *Spat. Vis.* 10 (4), 437–442. doi:10.1163/156856897X00366.
- Sereno, M.I., Dale, A.M., Reppas, J.B., Kwong, K.K., Belliveau, J.W., Brady, T.J., Rosen, B.R., Tootell, R.B., 1995. Borders of multiple visual areas in humans revealed by functional magnetic resonance imaging. *Science* 268 (5212), 889–893. doi:10.1126/science.7754376.
- Wandell, B.A., Dumoulin, S.O., Brewer, A.A., 2007. Visual field maps in human cortex. *Neuron* 56 (2), 366–383. doi:10.1016/j.neuron.2007.10.012.