

Simple Regularisation for Uncertainty-Aware Knowledge Distillation

Martin Ferianc¹ Miguel Rodrigues¹

Abstract

Considering uncertainty estimation of modern neural networks (NNs) is one of the most important steps towards deploying machine learning systems to meaningful real-world applications such as in medicine, finance or autonomous systems. At the moment, ensembles of different NNs constitute the state-of-the-art in both accuracy and uncertainty estimation in different tasks. However, ensembles of NNs are unpractical under real-world constraints, since their computation and memory consumption scale linearly with the size of the ensemble, which increase their latency and deployment cost. In this work, we examine a simple regularisation approach for distribution-free knowledge distillation of ensemble of machine learning models into a single NN. The aim of the regularisation is to preserve the diversity, accuracy and uncertainty estimation characteristics of the original ensemble without any intricacies, such as fine-tuning. We demonstrate the generality of the approach on combinations of toy data, SVHN/CIFAR-10, simple to complex NN architectures and different tasks.

1. Introduction

Neural networks (NNs) have enjoyed overwhelming interests in the recent past, due to their automatic feature learning abilities translating to super-human accuracy (LeCun et al., 2015). However, the deployment of the NNs in the real-world requires more than high accuracy. The NN-based system needs to be trustworthy, especially when presented with data that it has previously not observed. Trust can be built through estimating the uncertainty of the model and its estimation enables the users to gauge whether the model is wrong or lacks sufficient knowledge to solve the task at hand (Bhatt et al., 2021). Therefore, considering uncertainty estimation in modern NNs is increasingly important especially for safety-critical applications such as in healthcare

The code is at: https://github.com/martinferianc/hydra_plus. This is a work-in-progress.
¹Department of Electronic and Electrical Engineering, University College London, London, UK. Correspondence to: Martin Ferianc <martin.ferianc.19@ucl.ac.uk>.

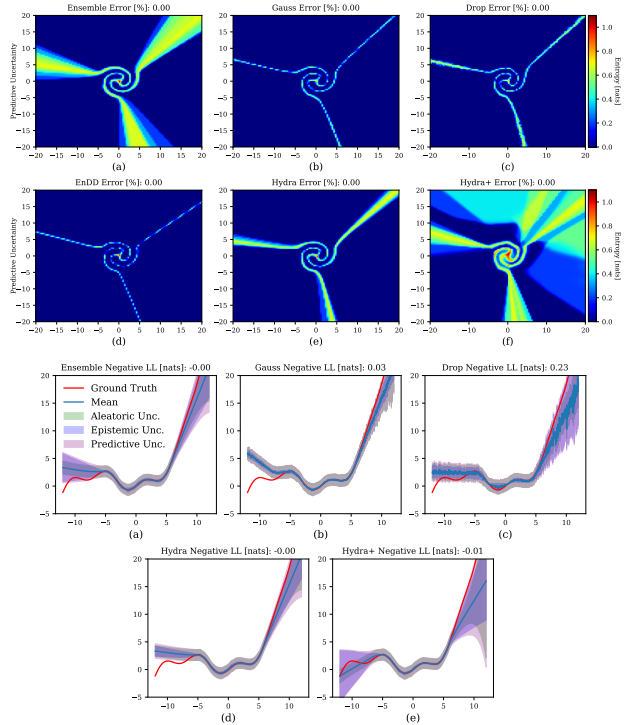


Figure 1: Toy classification/regression demonstrating KD of an *Ensemble* with $N = 20$ members into a student, while comparing different methods: *Gauss*, *Drop*, *EnDD*, *Hydra* and proposed *Hydra+* (a-f, a-e). The plots compare predictive, epistemic or aleatoric uncertainties, with one standard deviation for regression. The titles show the error/negative log-likelihood (LL) on test samples/curve. See Sections 2 or 4 for details. The decomposition of aleatoric and epistemic uncertainty in classification is in the Appendix A.

or self-driving (Abdar et al., 2021).

At the moment, ensembles of NNs (Zaidi et al., 2021) provide the best quantitative and qualitative results in terms of accuracy or uncertainty estimation. The ensemble can be created as simply as training different machine learning models with different seeds, meaning different initialisation of their parameters. The initialisation differentiation facilitates distinct optimisation trajectories of the ensemble members, ending in diverse local minima, which gives ensembles their representation capacity (Lakshminarayanan

et al., 2016). Then, it is possible to estimate the uncertainty of the complete ensemble through the disagreement of its individual members. In comparison to rigorous distribution-free uncertainty estimation methods (Angelopoulos & Bates, 2021), the user only needs to have access to the training data and train a machine learning model several times without any other assumptions to form an ensemble.

However, deploying ensembles of NNs in the real-world constitutes a challenge, since their resource consumption scales linearly with the size of the ensemble. *Knowledge distillation* (KD) (Hinton et al., 2015; Wang & Yoon, 2021) has been previously successfully utilised to compress the representation of the ensemble to a single NN, without any assumptions about the distilled model. Nonetheless, capturing the uncertainty of the ensemble without additional data or fine-tuning remains a challenge even with modern KD methods for capturing the ensemble’s uncertainty, shown in Figure 1.

In this work, we build on the *Hydra* KD idea proposed by (Tran et al., 2020) and we examine a simple regularisation to improve KD in capturing the uncertainty of the original ensemble, which we denote *Hydra+*. The regularisation is composed of two parts: *a*) modification of the loss function to capture *correctness*, *aggregated* and *individual* performance of the ensemble; *b*) minimisation of similarity between weights of different predictive heads of the multi-head student NN to promote *diversity*. These changes result in improved uncertainty estimation and calibration without requiring additional data, modelling assumptions or fine-tuning. We demonstrate the generality of the examined approach with respect to classification and regression on toy-data, SVHN, CIFAR-10 and simple feed-forward or convolutional, residual (He et al., 2016) architectures.

2. Preliminaries & Related Work

We now cover the preliminaries about ensembles, knowledge distillation and uncertainty decomposition. We also overview the related work.

2.1. Preliminaries

Ensembles Uncertainty estimation is increasingly gaining traction in the machine learning community in order to boost interpretability of the NNs deployed in the real-world (Bhatt et al., 2021). Despite different sophisticated attempts (Ovadia et al., 2019), ensembles maintain the state-of-the-art in uncertainty estimation, without requiring any particular assumptions about the data or the task (Lakshminarayanan et al., 2016; Wenzel et al., 2020; Zaidi et al., 2021). The process of building a baseline ensemble is simple - train a set of NNs on the same data, with the same architecture, but initialised with different random seeds (Lakshminarayanan

et al., 2016). Despite their generalisation and uncertainty estimation performance, ensembles are difficult to deploy in practice, since their compute and memory demands scale with complexity $\mathcal{O}(N)$ where N is the size of the ensemble.

Knowledge Distillation (KD) The deployment challenge of ensembles served as one of the inspirations for the idea of *Knowledge distillation* (KD) (Hinton et al., 2015), which aims to compress a large model, or a set of models to a single, smaller less demanding model. In general, KD achieves this through guiding the small model - the student to mimic the behaviour of the large ensemble - the teacher (Wang & Yoon, 2021). Concretely, the guidance between teacher and student is implemented through minimising the Kullback-Leibler (KL) (Kullback & Leibler, 1951) divergence between the likelihoods of the teacher and the student as: $KL(p(\hat{y}|\mathbf{x}, \theta_T)||p(\hat{y}|\mathbf{x}, \theta_S))$ where \hat{y} , \mathbf{x} are the output prediction and input and θ_T , θ_S are the parametrisations of the ensemble and the student respectively. In finer granularity, the guidance can be pointed towards different characteristics of the teacher, which is the investigation of this work. Namely, our focus is on the ability of the student to capture both the generalisation performance and aleatoric and epistemic uncertainty (Hüllermeier & Waegeman, 2021) of the ensemble as close as possible.

Uncertainty Decomposition Unless the knowledge of the data generating process is known, the proposed models always contain a notion of uncertainty (Hüllermeier & Waegeman, 2021). Furthermore, this uncertainty can be decomposed into uncertainty relating to incorrect model assumptions: *epistemic* or noisy data: *aleatoric* uncertainty. The uncertainty’s decomposition enables practitioners to understand what the model does not know. For example, in classification, if D is the training dataset, \hat{y} are the softmax probabilities produced by a model on data \mathbf{x} , parametrised by θ as $p(\hat{y}|\mathbf{x}, \theta)$ and \mathbb{H} , \mathbb{E} and \mathbb{MII} are the entropy, expectation and mutual information operators, it corresponds to (Hüllermeier & Waegeman, 2021):

$$\underbrace{\mathbb{H}(\mathbb{E}_{p(\theta|D)}[p(\hat{y}|\mathbf{x}, \theta)])}_{\text{Predictive Unc.}} = \underbrace{\mathbb{E}_{p(\theta|D)}[\mathbb{H}(p(\hat{y}|\mathbf{x}, \theta))]}_{\text{Aleatoric Unc.}} + \underbrace{\mathbb{MII}[\hat{y}, \theta|\mathbf{x}, D]}_{\text{Epistemic Unc.}} \quad (1)$$

In practice the \mathbb{MII} term cannot be computed in a closed-form in NNs, but it can be simply computed by subtracting the aleatoric uncertainty from the predictive uncertainty. Additionally, the expectations are often approximated with empirical averages using N Monte Carlo samples from the learnt posterior distribution $p(\theta|D)$ for that given method. In ensembles, the samples would correspond to the N different learnt models.

If considering regression and Gaussian likelihood, then:

$$\underbrace{\sigma^2}_{\text{Predictive Unc.}} = \underbrace{\hat{\sigma}^2}_{\text{Aleatoric Unc.}} + \underbrace{\mathbb{V}_{p(\theta|D)}[p(\hat{y}|\mathbf{x}, \theta)]}_{\text{Epistemic Unc.}} \quad (2)$$

$\hat{\sigma}^2$ is the predicted aleatoric variance, \hat{y} is the predicted mean and \mathbb{V} is the variance operator. In this work, we concentrate specifically on the KD methods that focus on distilling aleatoric and epistemic uncertainty, along with generalisation performance, of the teacher, irrespective of the teacher training procedure or its architecture.

2.2. Related Work

In distillation of uncertainty, Malinin et al. (2019) propose *EnDD* by using a prior network (Malinin & Gales, 2018) as the student, however, their approach requires further fine-tuning on auxiliary data to fully capture the ensemble’s uncertainty and it works only for classification problems. Tran et al. (2020) proposed *Hydra*: a multi-headed model where each head is paired with a member of the ensemble, while reusing a shared core architecture. The heads aim to capture the diversity of the ensemble, while reusing the common features. However, their approach requires multiple-steps, fine-tuning and inflexibility in choosing the student architecture. Lindqvist et al. (2020); Shen et al. (2021) propose a simple distillation method for learning the conditional predictive distribution of the ensemble or a Bayesian NN, into a flexible parametric distribution modelled by the last layer of the NN. For comparison, we consider two such distributions, through adding dropout before the last layer (Gal & Ghahramani, 2016) or using the local-reparametrisation trick with a Gaussian mean-field prior (Kingma et al., 2015) and we denote them as *Drop* and *Gauss* respectively. The downside of parametrising a distribution through the last layer is that it is necessary to assume some prior distribution, which can be unintentionally misspecified (Fortuin, 2022). In summary, ensemble constitutes the baseline assumption-free performance which distribution-specified *EnDD*, *Drop* and *Gauss* and assumption-free: *Hydra* or *Hydra+* KD methods target to match. *Hydra+* builds on (Tran et al., 2020) and attempts to avoid multi-step training while maintaining the generalisation and improving the uncertainty estimation performance through modifying the training loss function and including a diversity-inducing term.

3. Method

We now define the loss decomposition along with the diversity-inducing regularisation applied during training.

3.1. Loss Decomposition

The goal of KD is to match the performance of the teacher ensemble with N members, parametrised by $\theta_T = \{\theta_T^n\}_{n=1}^N$ through a student parametrised by θ_S on data tuples (\mathbf{x}, \mathbf{y})

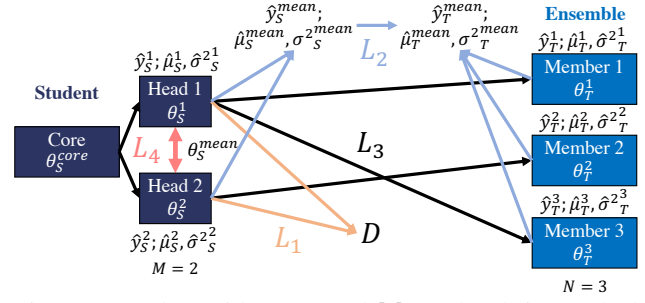


Figure 2: Student with a core and $M = 2$ heads is matched with ensemble consisting of $N = 3$ members. Training is performed with respect to the 4 loss components: L_1 correctness, L_2 aggregation, L_3 individuality and L_4 diversity.

initially coming from training set D , where \mathbf{x} and \mathbf{y} are the input and the desired output. If the end task is classification, e.g. categorisation of images into some classes, the output are the probabilities of the one-hot encoded labels $\mathbf{y} \in \mathbb{R}^K$ with K classes as $\text{Cat}(\mathbf{y}|\hat{\mathbf{y}}_T)$ with $\hat{\mathbf{y}}_T \in \mathbb{R}^{N \times K}$ or $\text{Cat}(\mathbf{y}|\hat{\mathbf{y}}_S)$ and $\hat{\mathbf{y}}_S \in \mathbb{R}^K$ for the teacher and the student respectively, provided that $\hat{\mathbf{y}}$ are obtained through the softmax activation at the output of both ensemble members and the student. If the task is regression, e.g. a prediction of a stock price, the output in both instances $y \in \mathbb{R}$ is modelled as a Gaussian with mean and aleatoric variance as $\mathcal{N}(y|\hat{\mu}_T, \hat{\sigma}_T^2)$ with $\hat{\mu}_T, \hat{\sigma}_T^2 \in \mathbb{R}^N$ or $\mathcal{N}(y|\hat{\mu}_S, \hat{\sigma}_S^2)$ with $\hat{\mu}_S, \hat{\sigma}_S^2 \in \mathbb{R}$ for the teacher and the student.

Inspired by *Hydra* (Tran et al., 2020), in this work we decompose the student network into two parts, the shared core parametrised by θ_S^{core} and M heads with the same structure such that $\theta_S = \{\theta_S^{\text{core}}, \{\theta_S^m\}_{m=1}^M\}$. The purpose of the core part of the network is to capture common features, while the M heads are supposed to capture the individual intricacies of the teacher ensemble members and enable the decomposition of the prediction into aleatoric and epistemic uncertainty, since the output is $\hat{\mathbf{y}}_S \in \mathbb{R}^{M \times K}$ and $\hat{\mu}_S, \hat{\sigma}_S^2 \in \mathbb{R}^M$. The relationship between the teacher and the student along with the used notation is visualised in Figure 2. Next, we introduce the 4 components L_1, L_2, L_3, L_4 of the proposed loss function L that aims to promote *correctness* in the student and capture *aggregated* and *individual* behaviour of the ensemble together with its *diversity*.

Correctness While proposing KD Hinton et al. (2015) noted that it is necessary to ensure correctness of the student by not only teaching it to mimic the behaviour of the teacher but also making it correct with respect to the training labels or values depending on classification or regression. We adopt this notion with respect to each student head’s m output $\hat{y}_S^m; \hat{\mu}_S^m, \hat{\sigma}_S^2$, where each head should be independently correct. This concept was previously considered in parameter-shared ensembles (Lee et al., 2015). For classi-

fication with the correct label y_k for class k it corresponds to minimising the mean of the cross-entropy across the M heads:

$$L_1^{class}(\hat{y}_S; y_k) = -\frac{1}{M} \sum_{m=1}^M y_k \log \hat{y}_S^m \quad (3)$$

For regression for target y and unit variance the Gaussian negative log-likelihood reduces to:

$$L_1^{reg}(\hat{\mu}_S^m, \hat{\sigma}_S^{2m}; y) = \frac{1}{2M} \sum_{m=1}^M \left(\frac{(\hat{\mu}_S^m - y)^2}{\hat{\sigma}_S^{2m}} + \log \hat{\sigma}_S^{2m} \right) + C \quad (4)$$

C is a constant term not affecting θ_S during optimisation.

Aggregation Next, [Hinton et al. \(2015\)](#) proposed the distillation itself where the output of the ensemble is averaged with respect to N to capture its overall behaviour. For classification this means $\hat{y}_T^{mean} = \frac{1}{N} \sum_{n=1}^N \hat{y}_T^n$. Similarly, we average the student output $\hat{y}_S^{mean} = \frac{1}{M} \sum_{m=1}^M \hat{y}_S^m$ where the loss is then the minimisation of the rearranged KL divergence between teacher and student outputs as:

$$L_2^{class}(\hat{y}_T^{mean}; \hat{y}_S^{mean}) = -\sum_{k=1}^K \hat{y}_T^{k,mean} \log \hat{y}_S^{k,mean} + C \quad (5)$$

the logits of the student and teacher can also be softened via division with temperature $T_{mean} \geq 1$ prior to softmax to give \hat{y}_T or \hat{y}_S .

For regression, we decide to aggregate the teacher and student prediction as means $\hat{\mu}_T^{mean} = \frac{1}{N} \sum_{n=1}^N \hat{\mu}_T^n$ and $\hat{\mu}_S^{mean} = \frac{1}{M} \sum_{m=1}^M \hat{\mu}_S^m$. However, for variance we not only take the mean of the aleatoric variance but we calculate the variance of both the teacher and student prediction to give $\sigma_T^{2mean} = \frac{1}{N} \sum_{n=1}^N \hat{\sigma}_T^{2n} + \mathbb{V}[\hat{\mu}_T]$ and $\sigma_S^{2mean} = \frac{1}{M} \sum_{m=1}^M \hat{\sigma}_S^{2m} + \mathbb{V}[\hat{\mu}_S]$ and thus we capture the complete predictive uncertainty. Then, the student output is compared to the teacher again via KL divergence between Gaussians:

$$L_2^{reg}(\hat{\mu}_T^{mean}, \sigma_T^{2mean}; \hat{\mu}_S^{mean}, \sigma_S^{2mean}) = \frac{1}{2} \left(\frac{\sigma_T^{2mean} + (\hat{\mu}_T^{mean} - \hat{\mu}_S^{mean})^2}{\sigma_S^{2mean}} + \log \sigma_S^{2mean} \right) + C \quad (6)$$

Individuality The primary loss with respect to which [Tran et al. \(2020\)](#) were training the *Hydra* in their second phase was matching the M individual heads to the N ensemble members, also conditioning that $M = N$. The motivation behind this loss was to urge each head to learn the representation of the individual ensemble member. If only L_3 is being used, all observable

individuality is lost ([Lee et al., 2015](#)). We relax the equality constraint on the number of heads M in order to explore algorithmic-hardware trade-offs from reducing M such that $2 \leq M \leq N$. If $N > M$, the remaining $N - M$ ensemble members are fairly divided between the M heads. Again for classification and outputs of the teacher ensemble \hat{y}_T and the student \hat{y}_S this KL divergence between the teacher and the student rearranges to:

$$L_3^{class}(\hat{y}_T; \hat{y}_S) = -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \hat{y}_T^{k,n} \log \hat{y}_S^{k,n \% M} + C \quad (7)$$

where $\%$ represents the modulo operator. The logits of the student and teacher can again be softened via division with temperature $T_{ind} \geq 1$ prior to softmax to give \hat{y}_T or \hat{y}_S .

Likewise for regression and outputs $\hat{\mu}_T, \hat{\sigma}_T^2; \hat{\mu}_S, \hat{\sigma}_S^2$ for teacher and student the KL divergence between the teacher and the student can be formulated as:

$$L_3^{reg}(\hat{\mu}_T, \hat{\sigma}_T^2; \hat{\mu}_S, \hat{\sigma}_S^2) = \frac{1}{N} \sum_{n=1}^N \frac{1}{2} \left(\frac{\hat{\sigma}_T^{2n} + (\hat{\mu}_T^n - \hat{\mu}_S^{n \% M})^2}{\hat{\sigma}_S^{2n \% M}} + \log \hat{\sigma}_S^{2n \% M} \right) + C \quad (8)$$

Diversity We empirically observed that it is not possible to induce diversity in one-shot training in the student by using L_3 alone. Therefore, we examine a differentiable diversity-inducing term calculated between the weights of the heads of the student $\{\theta_S^m\}_{m=1}^M$ at each layer-level $l = 1, \dots, L$, where $l = 1$ is the first weight-containing layer in the head and $l = L$ is the output layer. The core idea is to reduce the similarity between the weights and repulse them from each other at each level in order to obtain diverse responses to the same input processed through the shared core.

We define the mean head weight at an arbitrary level simply as $\theta_S^{mean} = \frac{1}{M} \sum_{m=1}^M \theta_S^m$. Given the abstract mean weight representation, we propose to minimise the similarity between the mean head weight and the individual head weights for each layer-level $l = 1, \dots, L$ of the head as:

$$L_4(\theta_S^{mean}; \{\theta_S^m\}_{m=1}^M) = \sum_{l=1}^L \sum_{m=1}^M \frac{1 + \cos(\theta_S^{l,mean}, \theta_S^{l,m})}{2} \quad (9)$$

We adopt the rescaled (0-1) cosine (cos) similarity as the main measure for pushing the weights apart from each other. The rationale behind using the cosine similarity is that it is a proper, differentiable distance metric which is irrespective of magnitude of the weights. By minimising this objective, each head's weights try to move away from the average of the heads' weights at that given level. Thus, the response

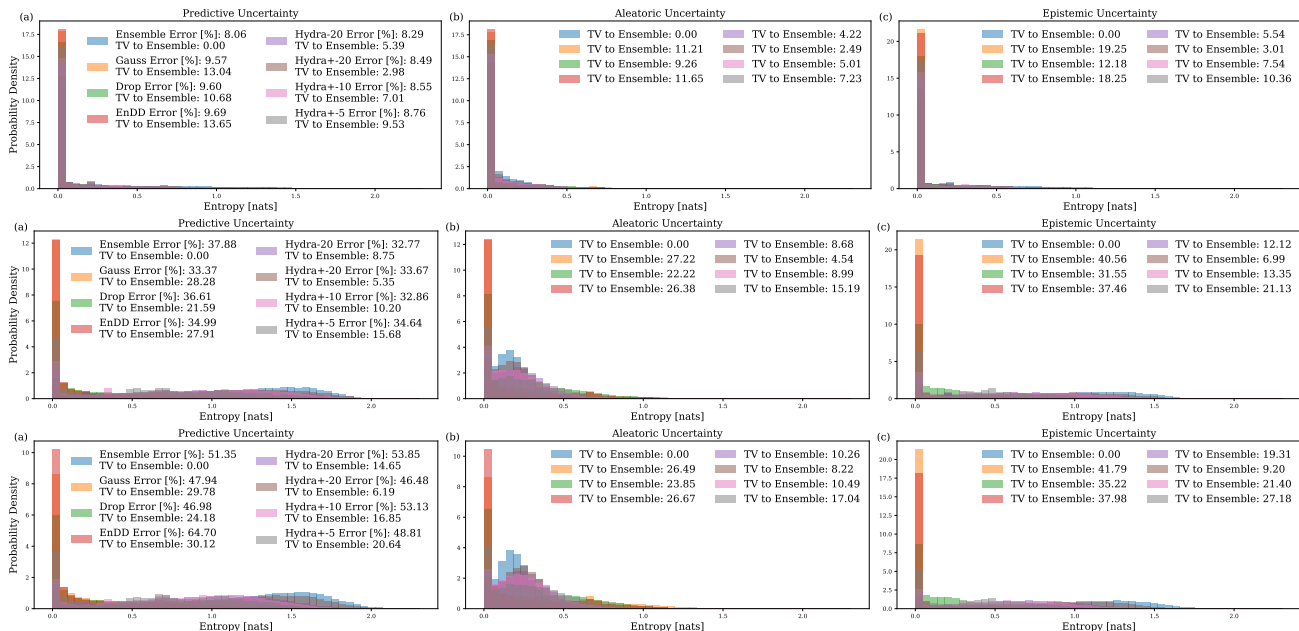


Figure 3: Uncertainty decomposition for augmentations applied to the SVHN test set and the compared methods for one seed/experiment. From the top, in the first row (a-c) is increasing the brightness by 30%, second row (a-c) is the 15° rotation and the third row (a-c) is 20% vertical shift. TV denotes total variation and Error denotes the error on the augmented test set.

of each head is induced to be different to all the other ones. Additionally, by comparing the head weights to their abstract mean, instead of a pair-wise comparison, we reduce the compute demand from $\mathcal{O}(M^2)$ to $\mathcal{O}(M)$. On the lowest level, the similarity is implemented between the weights of each separate node or filter, if considering the convolution operation. Related arguments for inducing diversity in NNs were reviewed in (Gong et al., 2019). However, to the best of our knowledge, there are no other related works to minimising the similarity between heads in multi-headed KD focused on capturing decomposable uncertainty estimation.

3.2. Complete Loss

Finally, Eqs. 3 & 4, 5 & 6, 7 & 8 and 9 can be merged together, with a slight abuse of notation, to give a differentiable optimisation objective L :

$$L = (1 - \alpha)L_1 + \alpha((1 - \beta)L_2 + \beta L_3) + \lambda L_4 \quad (10)$$

which is optimised with respect to θ_S , the teacher responses and the training dataset D via gradient descent in one training session for both classification and regression. The $0 \leq \alpha, \beta \leq 1$ and $\lambda \geq 0$ are hyperparameters dividing the focus between the loss components. λ can be kept constant from the start of training or increased linearly from and to a certain iteration. Note that, if $T_{ind}, T_{mean} > 1$, the respective loss components are being rescaled through T_{ind}^2 or T_{mean}^2 as discussed in (Hinton et al., 2015).

4. Experiments

In this Section we evaluate the outlined methodology on different experiments. Moreover, we present ablations to the hyperparameter choices and we discuss known limitations that could be targeted in the future work. We varied the dataset and architecture choices to change the complexity of the experiments from toy-data to CIFAR-10 and from feed-forward NNs to ResNet-18 in classification and regression to thoroughly test all the discussed methods.

Hyperparameter Settings First, we discuss shared hyperparameters across all methods for a fair comparison. Dropout rate was set to 0.5 when used in *Drop* and γ for regularisation of KL-divergence in *Gauss* was set to $1/\text{training set size}$ with zero mean, unit variance mean-field prior. No data augmentation was used in training except normalisation. No fine-tuning with or without extra data after training was allowed for any methods. All training had to be performed in one session. For *Hydra*, *Drop* and *Gauss*, $\beta = 1.0$, for *Hydra* $\alpha = 1.0$ and for *EnDD* $\beta = 0$. Further settings for all experiments are described in the Appendix A.

Metrics Second, we discuss the metrics. For classification, we were observing the classification error and expected calibration error (ECE) (Guo et al., 2017) with 10 uniformly-spaced bins. For regression, we were measuring the negative log-likelihood (NLL) for a Gaussian output. Specifically, for image-based experiments, we visually and quantitatively compare the decomposed un-

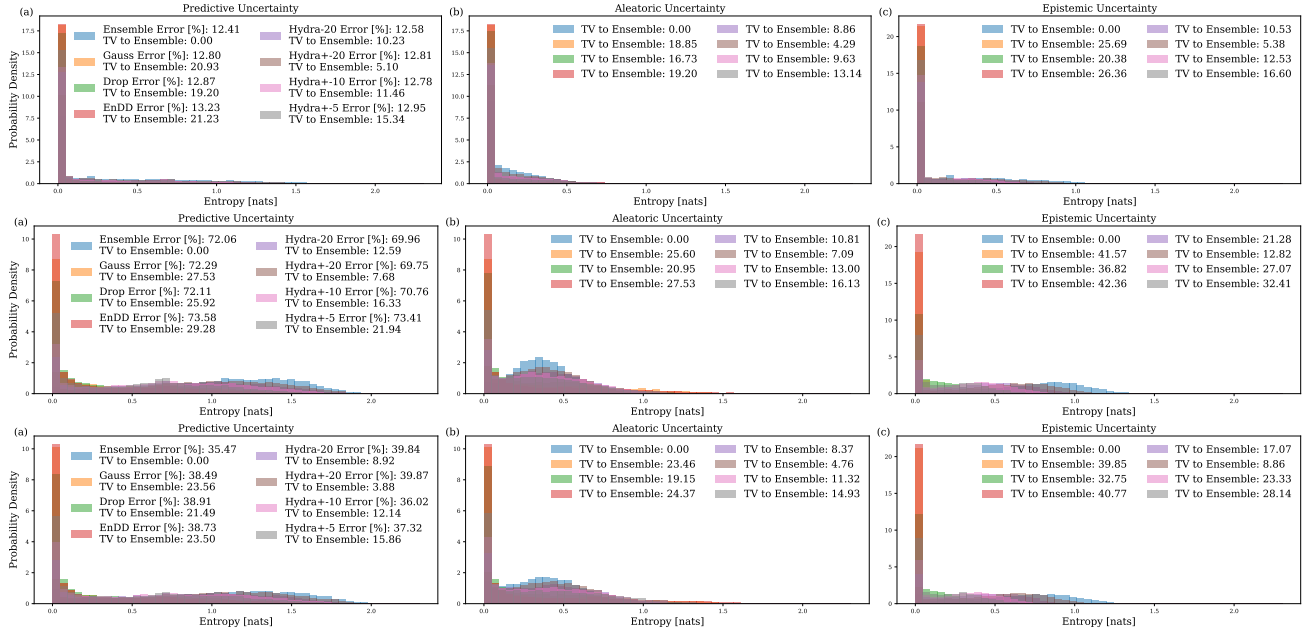


Figure 4: Uncertainty decomposition for augmentations applied to the CIFAR-10 test set and the compared methods for one seed/experiment. From the top, in the first row (a-c) is increasing the brightness by 30%, second row (a-c) is the 45° rotation and the third row (a-c) is 30% vertical shift. TV denotes total variation and Error denotes the error on the augmented test set.

certainty for all test set predictions according to Eq. 1 in normalised histograms and their total variation (TV): $TV(H_a, H_b) = \sum_{i=1}^B |H(i)_b - H(i)_a|$, where H_a, H_b are histograms with B bins. We used $B = 50$ bins. From the hardware perspective, we were comparing the number of parameters (#Params) or required floating-point operations (#FLOPS).

4.1. Toy Experiments

Classification In the toy classification experiment, seen in the first part of Figure 1, we constructed simple feed-forward NNs with 3 hidden layers and ReLU activations, where the last 2 layers served as the heads for *Hydra* and *Hydra+*. As it can be seen, without observing any additional data or fine-tuning, *Hydra+* was able to be at least as uncertain as the teacher ensemble, while having only 228560/618060, 37% of the parameters and 230860/626060 of FLOPS when compared to an ensemble when $N, M = 20$. As seen in Figure 9, the method can also near the quality of uncertainty decomposition of the original ensemble.

Regression Furthermore, we demonstrate the applicability of the examined method on regression as seen in the second part of Figure 1. We constructed a feed-forward NN with 2 hidden layers and ReLU activations where again the last 2 layers served as the heads for *Hydra* and *Hydra+*. In the unobserved regions, *Hydra+* was able to capture the uncertainty of the ensemble, while being able to generalise better than the ensemble, seen in the region where the in-

put $x < -6$. Through $N, M = 20$ for *Hydra+* or *Hydra*, the number of parameters or FLOPS were reduced from 106040 or 109040 for the ensemble to 55690 or 56790 for the student, denoting roughly a 48% decrease in the required computations and memory consumption for parameters.

4.2. Image-based Experiments

For the image-based experiments, we compared their performance under uncertainty on the test data and its augmentations through changes of brightness, rotations or vertical shifts. Additional details, along with experiments, are in the Appendix A. All experiments were end-to-end repeated with 3 different random seeds for robustness.

SVHN The SVHN results are presented in Figure 3 and Table 1. For SVHN experiments, we adapted the LeNet architecture where the last fully-connected layers served as the heads for *Hydra* or *Hydra+*. Examining results in Table 1, it can be seen that mainly *Hydra-20* and *Hydra+-20* with $M = 20$ heads were able to come close to the performance of the ensemble, seen in error or calibration, while having approximately 13% of the FLOPS and 56% of the parameters of the ensemble. Primarily, as seen in Figure 3 *Hydra+-20* was able to significantly reduce the TV to the ensemble in all uncertainty types and in some instances, rows 2 and 3, improve on the error, also improving the overall calibration of the model. However, as seen in the Table 1 or Figure 3, reducing the number of heads to 10 or 5 has a detrimental effect on the overall performance.

Method	Error [%]	ECE [%]	#FLOPS [M]	#Params [M]
Ensemble-20	7.20±0.05	3.81±0.11	117.75	3.04
Gauss	8.53±0.21	6.08±0.13	13.48	0.13
Drop	8.70±0.14	5.17±0.08	15.07	0.12
EnDD	8.91±0.08	6.45±0.11	13.48	0.12
Hydra-20	7.49±0.06	3.10±0.10	15.07	1.71
Hydra+20	7.56±0.06	3.08±0.03	15.07	1.71
Hydra+10	7.59±0.04	3.14±0.01	14.23	0.88
Hydra+5	7.67±0.08	3.69±0.08	13.81	0.46

Table 1: Comparison on the test dataset of SVHN. The number after method name denotes N or M .

Method	Error [%]	ECE [%]	#FLOPS [G]	#Params [M]
Ensemble-20	10.92±0.10	4.62±0.20	2.81	55.95
Gauss	11.70±0.10	9.53±0.09	0.80	4.55
Drop	11.61±0.19	8.84±0.14	0.80	4.55
EnDD	11.80±0.29	9.70±0.34	0.80	4.55
Hydra-20	11.44±0.28	4.83±0.14	1.04	19.23
Hydra+20	11.54±0.16	3.74±0.10	1.04	19.23
Hydra+10	11.37±0.05	4.98±0.07	0.92	11.51
Hydra+5	11.29±0.07	6.47±0.11	0.85	7.64

Table 2: Comparison on the test dataset of CIFAR-10. The number after method name denotes N or M .

From the runtime perspective, including L_4 regularisation increased the training cycle by $\sim 43\%$ without any particular hardware optimisations if compared to *Hydra*.

CIFAR-10 The CIFAR-10 results are presented in Figure 4 and Table 2. For CIFAR-10 experiments we adapted the ResNet-18 where the last 2 blocks served as the heads for *Hydra* or *Hydra+*. As it can be seen from the Table 2, given the complexity of the task, no method was able to achieve lower error than the ensemble. However *Hydra+* with $M = 20$ was able to be better calibrated than the ensemble with 37% of the FLOPS and 34% of the parameters. Interestingly, as the number of heads was decreased the error did not significantly deteriorate, but due to the reduced capacity, the representation power was smaller, primarily seen in Figure 4. In Figure 4 it can be seen that the examined loss in Eq. 10 and regularisation were able to significantly decrease TV between the teacher and the student for predictive, aleatoric and epistemic uncertainty, especially when compared *Hydra-20* to *Hydra+-20*. In a more complex model, compared to the SVHN experiment, including L_4 regularisation prolonged the training by $\sim 6\%$ when compared to *Hydra*, given that the most of the parameters are in convolutions instead of fully-connected layers as in the SVHN experiment.

4.3. Ablations

Changing λ, β Additionally, we wanted to demonstrate the effects of changing λ and β hyperparameters for the student models. For a clear visualisation of the responses, the changes are illustrated on the toy problems as seen in

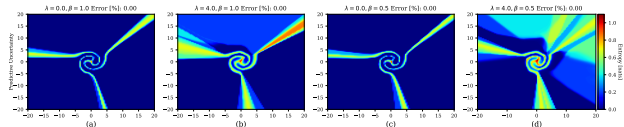


Figure 5: From the left to right, (a) disabling L_4 , $\beta = 1.0$, (b) enabling L_4 , $\beta = 1.0$, (c) disabling L_4 , $\beta = 0.5$, (d) enabling L_4 , $\beta = 0.5$ for the toy classification problem.

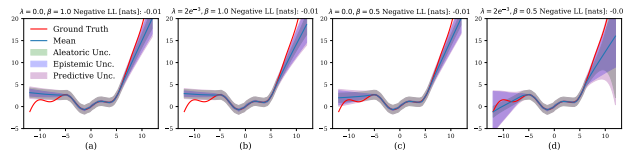


Figure 6: From the left to right, (a) disabling L_4 , $\beta = 1.0$, (b) enabling L_4 , $\beta = 1.0$, (c) enabling L_4 , $\beta = 0.5$, (d) enabling L_4 , $\beta = 0.5$ for the toy regression problem.

Figures 5 & 6. In general, we observe that the main portion of the observed epistemic uncertainty is caused by including the L_4 term in the loss function through λ , as seen in subplots (b) in both Figures 5 & 6. Conversely, β being less than 1 and thus including L_2 in the optimisation, guides the overall fit of the model and rectifies its predictive uncertainty, most notably the aleatoric component, which could otherwise result in uncalibrated predictions, as seen if comparing Figures 5 & 6 (a,c) in the regions close to the origin, where training data was actually observed. Interestingly, we observed that L_3 and L_4 are closely related and empirically we were unable to obtain as good results as seen in Figures 5 & 6 (d) without enabling both L_3 and L_4 during training.

Changing M Next, we discuss changing the number of heads in *Hydra+* through decreasing $M \leq N$. For the toy datasets, this is visualised in Figures 7 & 8 and for CIFAR-10 and SVHN experiments in Tables 1 & 2 or Figures 3 & 4. Decreasing M corresponds to decreasing the representation ability of the student, through practically decreasing its number of parameters and pushing a head to learn from multiple ensemble members. As a result, the smaller, less representative students' uncertainty representation power deteriorates. Nevertheless, their accuracy, does not necessarily need to follow the same trend, where a head learns to generalise better if focusing on more than one teacher. This is best seen in Table 2 for the CIFAR-10 experiment. Additionally, having to run fewer heads on hardware reduces the required computation and memory storage as seen in Tables 1 & 2 when comparing FLOPS or the number of parameters.

4.4. Discussion, Limitations & Future Work

Discussion [Mehrtens et al. \(2022\)](#), considering ([Benjamin et al., 2018](#)) came to the conclusion that parameter distance is no good measure for functional differences between ensemble members. Nevertheless, [Mehrtens et al. \(2022\)](#) sug-

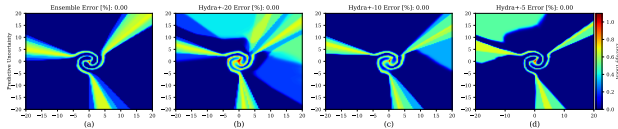


Figure 7: From the left to right, (a) Ensemble, (b) $M = 20$, (c) $M = 10$, (d) $M = 5$ for the toy classification problem.

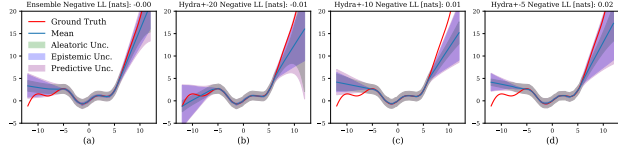


Figure 8: From the left to right, (a) Ensemble, (b) $M = 20$, (c) $M = 10$, (d) $M = 5$ for the toy regression problem.

gest that in certain instances (Wenzel et al., 2020), where parameters are not easily separable from each other, increasing orthogonality in parameter space could induce functional diversity. We believe that this particular setting, with KD and a shared core and intra-dependent training of the heads is another such case. We did try the distance mentioned (Mehrtens et al., 2022), adapted to multiple heads instead of ensemble members, or different cosine-based distance measures (Sohangir & Wang, 2017), but empirically they did not outperform the proposed measure quantitatively or qualitatively. We used $N = 20$ in the ensemble, such that it had strong representation power. We empirically observed that if N is small e.g. 5, irrespective of the hyperparameters, the students are unable to capture barely any uncertainty representation. Moreover, we noticed that if the N ensemble members are randomly shuffled for Eqs. 7 & 8, the heads are only able to retain general representation without any individuality, avoiding the capture of any uncertainty estimation of the teacher. We also empirically observed that *Hydra+* performs better when L_4 was not applied to batch normalisation weights. We empirically observed that it is beneficial to increase λ as M or N decrease. In certain instances, e.g. the CIFAR-10 experiment, it was better to apply L_4 later in the training and increase its influence as the training progresses. Similarly to Tran et al. (2020), we empirically observed that increasing the individual head’s size, the depth or width, improves their representation power, but it also increases the overall number of required parameters and FLOPS.

In observing the performance of *Gauss*, *Drop* or *EnDD* we see a trend that if the prior distribution is potentially misspecified the student network is unable to capture the uncertainty estimation of the teacher. This gives preference, to rather simpler, however more compute expensive, from the parameter of FLOPS standpoint, distribution-free methods such as *Hydra* which does not require any particular assumptions, but only hyperparameter tuning for the provided loss function.

Limitations In our experimentation we observed several corner cases. To begin with, at the moment it is necessary to manually define M . The architecture of the student was manually selected and it had to be hand-tuned on the validation dataset, along with other hyperparameters. Moreover, the examined regularisation focuses primarily on improving the estimation of the epistemic uncertainty, without the ability to control the aleatoric uncertainty other than relying on the performance of teacher. Last but not least, other loss functions could be examined which would focus on capturing the uncertainty estimation more prominently in comparison to KL divergence.

Future Work As seen for more challenging tasks, such as SVHN or CIFAR-10 classification, it is difficult to match the performance of the ensemble. We believe that to provide further adaptation through diversity in both the function space (Mehrtens et al., 2022) and parameter space, the student architecture could be optimised in addition to the parameters. Our current efforts are in applying neural architecture search (Chen et al., 2020) to find the architecture of the required M heads and core automatically. We are also experimenting in automatising the search for M , such that we can further reduce the FLOPS count or the number of parameters.

5. Conclusion

In this work-in-progress we examined a regularisation for knowledge distillation, to capture both the generalisation performance as well as uncertainty estimation of the teacher ensemble without any fine-tuning or extra data. We demonstrated the explored methodology on toy and real-world data and different architectures to show its versatility. In comparison to the underlying ensembles, the discussed regularisation was able to approach near or improve upon its quality of calibration and uncertainty estimation. In the future work, we aim to improve the student’s performance through automatic adaptation of the student architecture to the task and the teacher network.

Acknowledgements

Martin Ferianc was sponsored through a scholarship from the Institute of Communications and Connected Systems at UCL and through the PhD Enrichment scheme at The Alan Turing Institute. Lastly, we thank DFUQ’22 reviewers for feedback and encouragement.

References

Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi, A., Acharya, U. R., et al. A review of uncertainty

- quantification in deep learning: Techniques, applications and challenges. *Information Fusion*, 76:243–297, 2021.
- Angelopoulos, A. N. and Bates, S. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511*, 2021.
- Benjamin, A. S., Rolnick, D., and Kording, K. Measuring and regularizing networks in function space. *arXiv preprint arXiv:1805.08289*, 2018.
- Bhatt, U., Antorán, J., Zhang, Y., Liao, Q. V., Sattigeri, P., Fogliato, R., Melançon, G., Krishnan, R., Stanley, J., Tickoo, O., et al. Uncertainty as a form of transparency: Measuring, communicating, and using uncertainty. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 401–413, 2021.
- Chen, X., Wang, R., Cheng, M., Tang, X., and Hsieh, C.-J. Drnas: Dirichlet neural architecture search. *arXiv preprint arXiv:2006.10355*, 2020.
- Fortuin, V. Priors in bayesian deep learning: A review. *International Statistical Review*, 2022.
- Gal, Y. and Ghahramani, Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pp. 1050–1059. PMLR, 2016.
- Gong, Z., Zhong, P., and Hu, W. Diversity in machine learning. *IEEE Access*, 7:64323–64350, 2019. doi: 10.1109/access.2019.2917620. URL <https://doi.org/10.1109%2Faccess.2019.2917620>.
- Guo, C., Pleiss, G., Sun, Y., and Weinberger, K. Q. On calibration of modern neural networks. In *International Conference on Machine Learning*, pp. 1321–1330. PMLR, 2017.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Hinton, G., Vinyals, O., and Dean, J. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Hüllermeier, E. and Waegeman, W. Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. *Machine Learning*, 110(3):457–506, 2021.
- Kingma, D. P., Salimans, T., and Welling, M. Variational dropout and the local reparameterization trick. *Advances in neural information processing systems*, 28, 2015.
- Kullback, S. and Leibler, R. A. On information and sufficiency. *The annals of mathematical statistics*, 22(1): 79–86, 1951.
- Lakshminarayanan, B., Pritzel, A., and Blundell, C. Simple and scalable predictive uncertainty estimation using deep ensembles. *arXiv preprint arXiv:1612.01474*, 2016.
- LeCun, Y., Bengio, Y., and Hinton, G. Deep learning. *nature*, 521(7553):436–444, 2015.
- Lee, S., Purushwalkam, S., Cogswell, M., Crandall, D., and Batra, D. Why m heads are better than one: Training a diverse ensemble of deep networks. *arXiv preprint arXiv:1511.06314*, 2015.
- Lindqvist, J., Olmin, A., Lindsten, F., and Svensson, L. A general framework for ensemble distribution distillation. In *2020 IEEE 30th International Workshop on Machine Learning for Signal Processing (MLSP)*, pp. 1–6. IEEE, 2020.
- Malinin, A. and Gales, M. Predictive uncertainty estimation via prior networks. *Advances in neural information processing systems*, 31, 2018.
- Malinin, A., Mlodozieniec, B., and Gales, M. Ensemble distribution distillation. *arXiv preprint arXiv:1905.00076*, 2019.
- Mehrtens, H. A., González, C., and Mukhopadhyay, A. Improving robustness and calibration in ensembles with diversity regularization. *arXiv preprint arXiv:2201.10908*, 2022.
- Ovadia, Y., Fertig, E., Ren, J., Nado, Z., Sculley, D., Nowozin, S., Dillon, J., Lakshminarayanan, B., and Snoek, J. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. *Advances in neural information processing systems*, 32, 2019.
- Shen, Y., Zhang, Z., Sabuncu, M. R., and Sun, L. Real-time uncertainty estimation in computer vision via uncertainty-aware distribution distillation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 707–716, 2021.
- Sohangir, S. and Wang, D. Improved sqrt-cosine similarity measurement. *Journal of Big Data*, 4(1):1–13, 2017.
- Tran, L., Veeling, B. S., Roth, K., Swiatkowski, J., Dillon, J. V., Snoek, J., Mandt, S., Salimans, T., Nowozin, S., and Jenatton, R. Hydra: Preserving ensemble diversity for model distillation. *arXiv preprint arXiv:2001.04694*, 2020.
- Wang, L. and Yoon, K.-J. Knowledge distillation and student-teacher learning for visual intelligence: A review and new outlooks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.

Wenzel, F., Snoek, J., Tran, D., and Jenatton, R. Hyperparameter ensembles for robustness and uncertainty quantification. *Advances in Neural Information Processing Systems*, 33:6514–6527, 2020.

Zaidi, S., Zela, A., Elsken, T., Holmes, C. C., Hutter, F., and Teh, Y. Neural ensemble search for uncertainty estimation and dataset shift. *Advances in Neural Information Processing Systems*, 34, 2021.

A. Appendix

A.1. Further Experimental Settings

All experiments were performed with respect to 200 epochs, batch size 256, Adam optimiser with default β_1, β_2 hyperparameters, with cosine decreasing learning rate schedule and gradient clipping coefficient set to 5.0. 10% of any training set was set for validation to hand-tune the hyperparameters. PyTorch 1.11, CUDA 11.1 and Nvidia GeForce RTX 2080 SUPER GPU were used for implementation. The default PyTorch initialisation was used for the weights. Note that, we reimplemented the weight decay to be computed explicitly and then added to the loss and we did not use the default weight decay option for optimisers in PyTorch. The ensemble was always trained simply with respect to a cross-entropy loss or Gaussian negative log-likelihood. The M in context of *Gauss*, *Drop* or *EnDD* meant the number of Monte Carlo samples for training or evaluation. We have combined the predictions of the student or the teacher with respect to all N or M samples and averaged, after the softmax activation for classification and for regression, we combined the variance through Eq. 2.

For *Hydra* or *Hydra+* we initialised the architecture already with M heads to observe the performance after a singular training session. We wanted to benchmark the methods for the simplicity of implementation of a single training session and to avoid interference from multiple steps and choices of hyperparameters that would be needed if the original *training with multi-head growth* was to be considered. This is also reflected in the choices of other hyperparameters, where we did not want to disadvantage any other method, such that we can provide a fair comparison and outlook.

The image corruption experiments consisted of [20, 30, 50]% increases in brightness intensity, [15, 45, 75] degree rotations, [20, 30, 50]% vertical shifts for both datasets. We encourage the reader to see our code, how to generate additional results for corruptions not shown in the paper.

A.2. Toy Classification

The toy classification experiment was performed with respect to a feed-forward network with [2, 100, 100, 100, 100, 3] input, hidden nodes and output, ReLU activations, initial

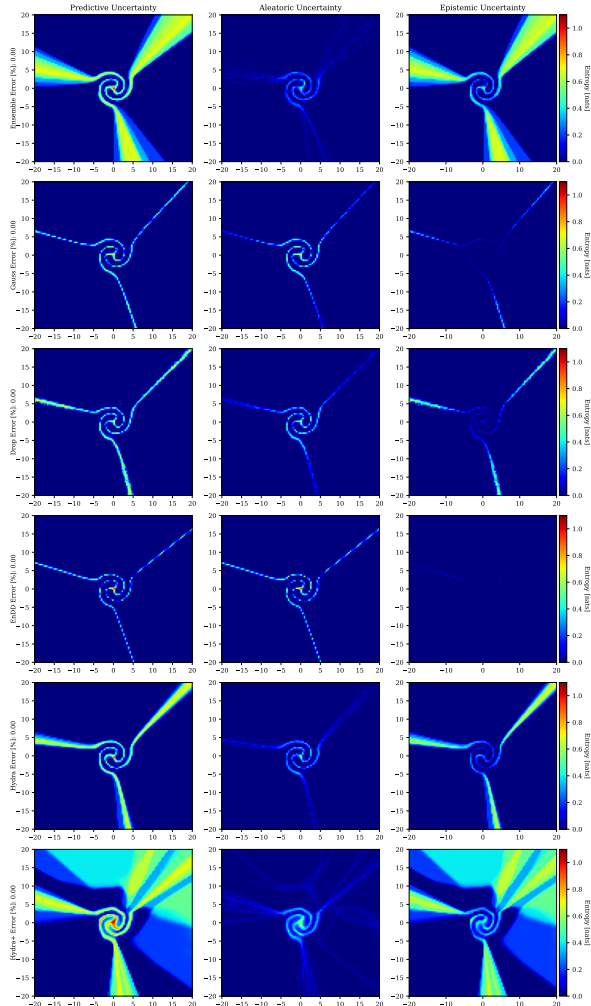


Figure 9: Decomposition of the uncertainty into aleatoric and epistemic parts for the toy classification spiral dataset and the compared methods from Figure 1.

learning rate 0.01, weight decay applied to the teacher training was set to $1e^{-4}$. The training, validation and test sets contained 240, 30, 30 points split equally across 3 classes sampled from 3 spirals originating at coordinates 0,0. For the KD of the student the weight decay was set to $1e^{-8}$, $\alpha = 0.9$, $T_{ind} = 3.0$, $T_{mean} = 1.0$ and β for *Hydra+* was set to 0.5. For 20, 10 and 5 heads of *Hydra+*, the λ was set to 4.0, 7.0 and 9.0 respectively.

Figure 9 shows further decomposition of uncertainty for the compared methods according to Eq. 1. Figure 10 shows further decomposition of uncertainty for the varying number of heads. Figure 11 shows further decomposition of uncertainty for the varying hyperparameters as discussed in the ablations in Section 4.3.

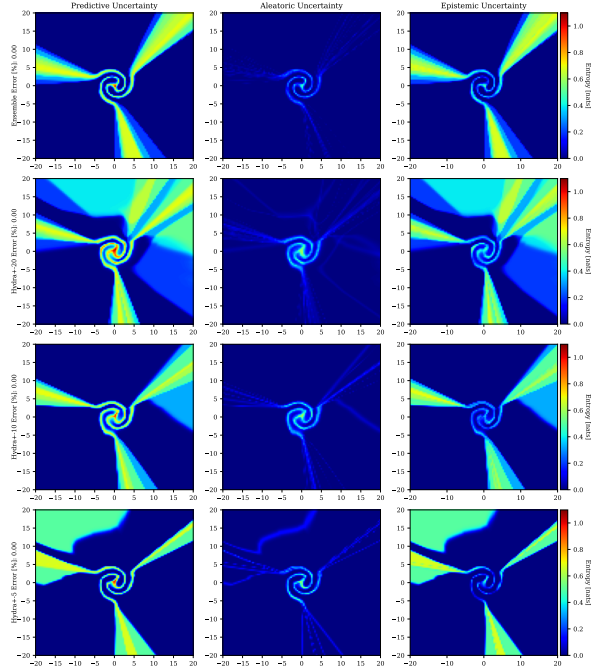


Figure 10: Decomposition of the uncertainty into aleatoric and epistemic parts for the toy classification spiral dataset and varying number of heads from Figure 7.

A.3. Toy Regression

The toy regression experiment was performed with respect to a feed-forward network with [2, 50, 50, 50, 2] input, hidden nodes and output, ReLU activations, initial learning rate 0.05, weight decay applied to the teacher training was set to $1e^{-5}$. The training, validation and test sets contained 240, 30, 30 points sampled from $f(x) = \sin(x) - 0.1x + 0.1x^2 + 0.01x^3$ from range -6 to 6 with additive Gaussian noise $\mathcal{N}(0, 1)$ added to the training points. For the KD of the student the weight decay was set to $1e^{-8}$, $\alpha = 0.9$ and β for *Hydra+* was set to 0.5. For 20, 10 and 5 heads of *Hydra+*, the λ were linearly increased from 0.0, beginning on the 50th to $2e^{-3}$, 0.02 or 0.6 culminating at the 150th epoch respectively.

A.4. SVHN

The SVHN classification experiment was performed with respect to a LeNet-like architecture with 2 convolutions, followed by ReLU and maxpooling with 2 fully-connected layers with [3, 32, 64, 1024, 128, 10] input, hidden and output channels for the teacher. The student had a similar architecture, however, with [3, 128, 32, 512, 128, 10] input, hidden and output split. The learning rate was initialised to 0.001, weight decay applied to the teacher training was set to $1e^{-4}$. The training, validation and test sets contained 65932, 7325 and 26032 samples split across 10 classes. For the KD of the student the weight decay was

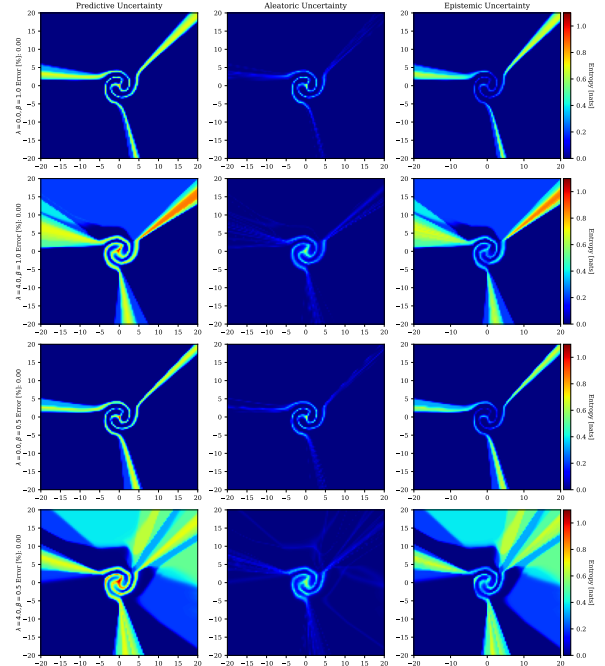


Figure 11: Decomposition of the uncertainty into aleatoric and epistemic parts for the toy classification spiral dataset and ablations from Figure 5.

set to $1e^{-8}$, $\alpha = 0.95$, $T_{ind} = 5.0$, $T_{mean} = 1.0$ and β for *Hydra+* was set to 0.9. For 20, 10 and 5 heads of *Hydra+*, the λ were set to $5e^{-4}$, $1e^{-3}$, $5e^{-3}$ respectively.

Figure 12 shows further decomposition of uncertainty for the compared methods according to Eq. 1 for more severe augmentations.

A.5. CIFAR-10

The CIFAR-10 classification experiment was performed with respect to a ResNet-18 architecture with 8 residual blocks and output channel sizes [32, 64, 128, 256]. The student had a similar architecture, however, with [96, 128, 256, 128] output channel sizes. The learning rate was initialised to 0.01, weight decay applied to the teacher training was set to $1e^{-4}$. The training, validation and test sets contained 45000, 5000 and 10000 samples split across 10 classes. For the KD of the student the weight decay was set to $1e^{-8}$, $\alpha = 0.95$, $T_{ind} = 8.0$, $T_{mean} = 1.0$ and β for *Hydra+* was set to 0.5. For 20, 10 and 5 heads of *Hydra+*, the λ were linearly increased from 0.0, beginning on the 20th to $2e^{-3}$, 0.05 or 0.1 culminating at the 150th epoch respectively.

Figure 13 shows further decomposition of uncertainty for the compared methods according to Eq. 1 for more severe augmentations.

Simple Regularisation for Uncertainty-Aware Knowledge Distillation

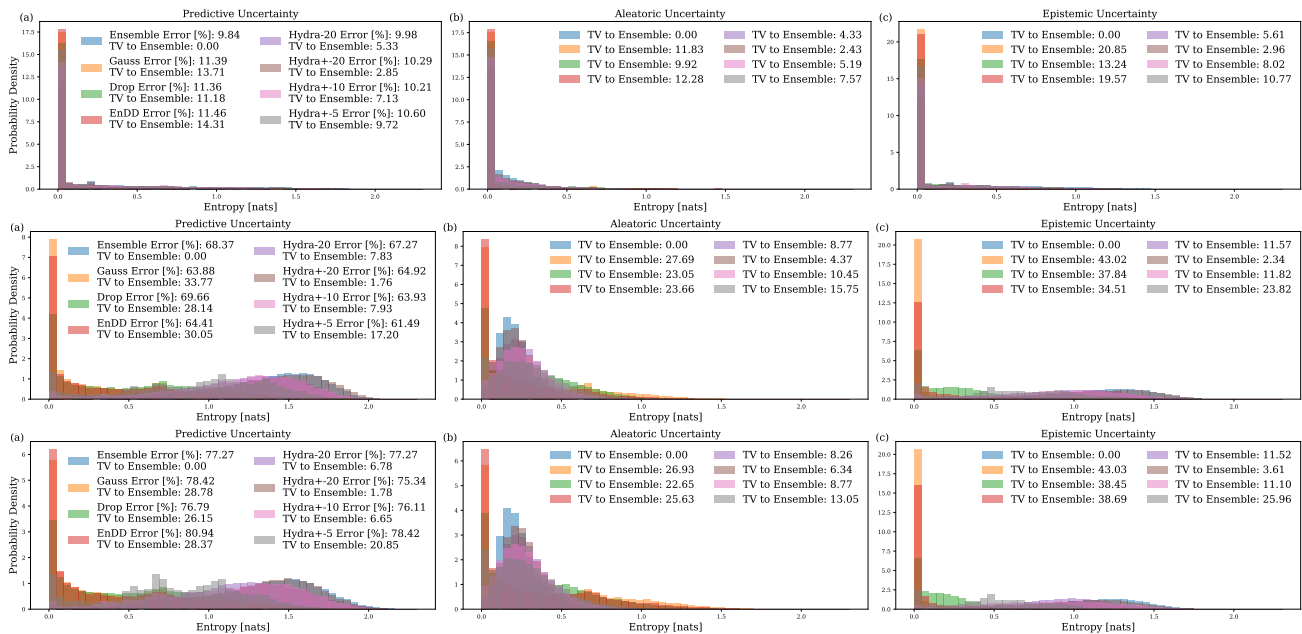


Figure 12: Uncertainty decomposition for augmentations applied to the SVHN test set and the compared methods for one seed/experiment. From the top, in the first row (a-c) is increasing the brightness by 50%, second row (a-c) is the 30° rotation and the third row (a-c) is 30% vertical shift. TV denotes total variation and Error denotes the error on the augmented test set.

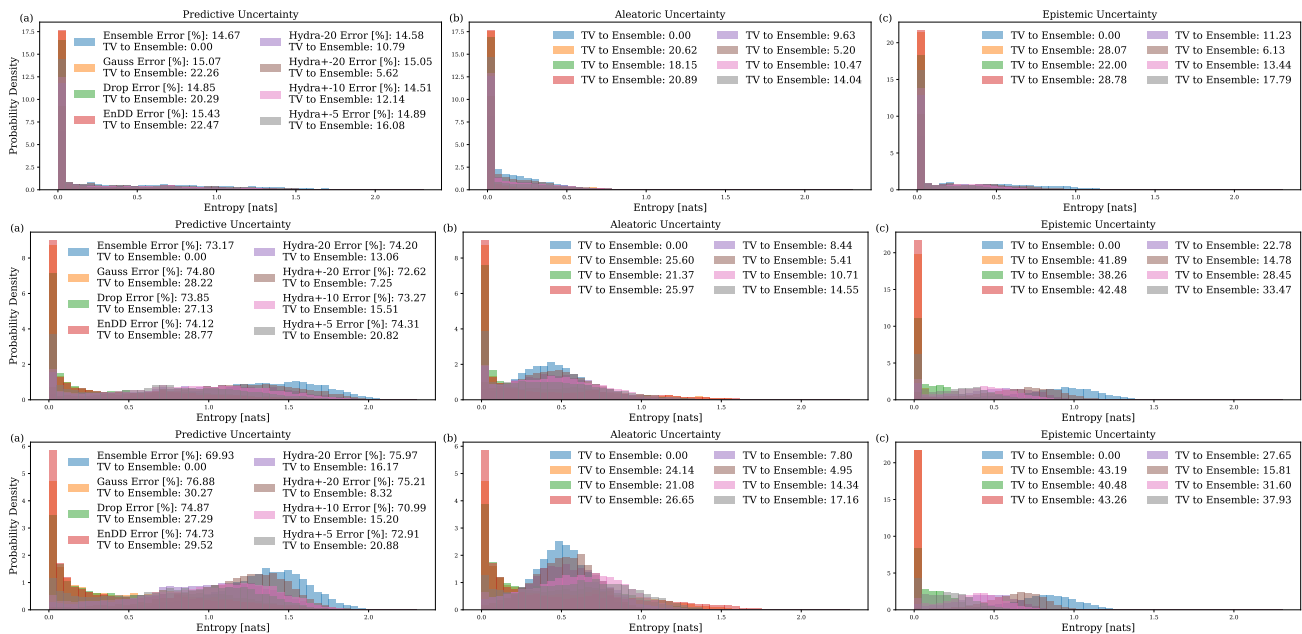


Figure 13: Uncertainty decomposition for augmentations applied to the CIFAR-10 test set and the compared methods for one seed/experiment. From the top, in the first row (a-c) is increasing the brightness by 50%, second row (a-c) is the 75° rotation and the third row (a-c) is 50% vertical shift. TV denotes total variation and Error denotes the error on the augmented test set.