

## A. Related Work

Several existing approaches model uncertainty using feature-space density but underperform without fine-tuning on OoD data. Our work identifies feature collapse and objective mismatch as possible reasons for this. Among these approaches, [43] uses Mahalanobis distances to quantify uncertainty by fitting a class-wise Gaussian distribution (with shared covariance matrices) on the feature space of a pre-trained ResNet encoder. The competitive results they report require input perturbations, ensembling GMM densities from multiple layers, and fine-tuning on OoD hold-out data. They do not discuss any constraints which the ResNet encoder should satisfy, and therefore, are vulnerable to feature collapse. In Fig. 1c, for example, the feature density of a LeNet and a VGG are unable to distinguish OoD from iD samples. [59] also propose a density-based estimation of aleatoric and epistemic uncertainty. Similar to [43], they do not constrain their pre-trained ResNet encoder. They do discuss feature collapse though, noting that they do not address this problem. Moreover, they do not consider the objective mismatch that arises (see Proposition 5.3 below) and use a single estimator for both epistemic and aleatoric uncertainty. Consequently, they report worse epistemic uncertainty: 74% AUROC on CIFAR-10 vs SVHN, which we show to considerably fall behind modern approaches for uncertainty estimation in deep learning in §4. Likewise, [46] compute an unnormalized density based on the softmax logits without taking into account the need for inductive biases to ensure smoothness and sensitivity of the feature space.

[69] use contrastive training on the feature extractor before estimating the feature-space density. Our method is orthogonal from this work as we restrict ourselves to the supervised setting and show that the inductive biases that result in bi-Lipschitzness [45; 65] are sufficient for the feature-space density to reliably capture epistemic uncertainty.

Lastly, our method improves upon [65] and [45] by alleviating the need for additional hyperparameters: DDU only needs minimal changes from the standard softmax setup to outperform DUQ and SNGP on uncertainty benchmarks, and our GMM parameters are optimised for the already trained model using the training set. In particular, DDU does not require training or fine-tuning with OoD data. Moreover, our insights in §5 explain why [45] found that a baseline that uses *the softmax entropy instead of the feature-space density* of a deterministic network with bi-Lipschitz constraint underperforms.

### A.1. Predictive Entropy and Confidence in Recent Works

Table 4 shows a selection of recently published papers which use entropy or confidence as OoD score. Only two papers examine using Mutual Information with Deep Ensembles as OoD score at all. None of the papers examines the



Figure 8. Samples from Ambiguous-MNIST.

possible confounding of aleatoric and epistemic uncertainty when using predictive entropy or confidence, or the consistency issues of softmax entropy (and softmax confidence), detailed in §5. This list is not exhaustive, of course.

## B. Ambiguous- and Dirty-MNIST

Each sample in Ambiguous-MNIST is constructed by decoding a linear combination of latent representations of 2 different MNIST digits from a pre-trained VAE [35]. Every decoded image is assigned several labels sampled from the softmax probabilities of an off-the-shelf MNIST neural network ensemble, with points filtered based on an ensemble’s MI (to remove ‘junk’ images) and then stratified class-wise based on their softmax entropy (some classes are inherently more ambiguous, so we “amplify” these; we stratify per-class to try to preserve a wide spread of possible entropy values, and avoid introducing additional ambiguity which will increase all points to have highest entropy). All off-the-shelf MNIST neural networks were then discarded and new models were trained to generate Fig 1 (and as can be seen, the ambiguous points we generate indeed have high entropy regardless of the model architecture used). We create 60K such training and 10K test images to construct Ambiguous-MNIST. Finally, the Dirty-MNIST dataset in this experiment contains MNIST and Ambiguous-MNIST samples in a 1:1 ratio (with 120K training and 20K test samples). In Fig. 8, we provide some samples from Ambiguous-MNIST.

## C. Algorithm

### C.1. Increasing sensitivity

Using residual connections to enforce sensitivity works well in practice when the layer is defined as  $x' = x + f(x)$ . However, there are several places in the network where additional spatial downsampling is done in  $f(\cdot)$  (through a strided convolution), and in order to compute the residual operation  $x$  needs to be downsampled as well. These downsampling

Table 4. *A sample of recently published papers and OoD metrics.* Many recently published papers only use Predictive Entropy or Predictive Confidence (for Deep Ensembles) or Softmax Confidence (for deterministic models) as OoD scores without addressing the possible confounding of aleatoric and epistemic uncertainty, that is ambiguous iD samples with OoD samples. Only two papers examine using Mutual Information with Deep Ensembles as OoD score at all.

Title	Citation	Softmax Confidence	Predictive Confidence	Predictive Entropy	Mutual Information
A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks	[24]	✓	✗	✗	✗
Deep Anomaly Detection with Outlier Exposure	[25]	✓	✗	✗	✗
Enhancing The Reliability of Out-of-distribution Image Detection in Neural Networks	[44]	✓	✗	✗	✗
Training Confidence-calibrated Classifiers for Detecting Out-of-Distribution Samples	[42]	✓	✗	✗	✗
Learning Confidence for Out-of-Distribution Detection in Neural Networks	[8]	✓	✗	✗	✗
Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles	[40]	✗	✓	✓	✗
Predictive Uncertainty Estimation via Prior Networks	[50]	✗	✓	✓	✓
Ensemble Distribution Distillation	[51]	✗	✗	✓	✓
Generalized ODIN: Detecting Out-of-Distribution Image Without Learning From Out-of-Distribution Data	[29]	✓	✗	✗	✗
Being Bayesian, Even Just a Bit, Fixes Overconfidence in ReLU Networks	[38]	✗	✓	✗	✗

operations are crucial for managing memory consumption and generalisation. The way this is traditionally done in ResNets is by introducing an additional function  $g(\cdot)$  on the residual branch (obtaining  $x' = g(x) + f(x)$ ) which is a strided 1x1 convolution. In practice, the stride is set to 2 pixels, which leads to the output of  $g(\cdot)$  only being dependent on the top-left pixel of each 2x2 patch, which reduces sensitivity. We overcome this issue by making an architectural change that improves uncertainty quality without sacrificing accuracy. We use a strided average pooling operation instead of a 1x1 convolution in  $g(\cdot)$ . This makes the output of  $g(\cdot)$  dependent on all input pixels. Additionally, we use leaky ReLU activation functions, which are equivalent to ReLU activations when the input is larger than 0, but below 0 they compute  $p * x$  with  $p = 0.01$  in practice. These further improve sensitivity as all negative activations still propagate in the network.

## C.2. Algorithm Implementation

The algorithm is provided in Algorithm 1. Note that in order to compute thresholds for low and high density or entropy, we can simply use the training set containing iD data. We set all points having density lower than 99% quantile as OoD.

## C.3. Computational Complexity

Let  $N$  be the number of samples;  $D$ , the feature space dimensionality; and  $C$ , the number of classes with  $\approx N/C$  samples per class (balanced). For fitting the GMM via GDA: computing the covariance matrix per class requires  $(C(N/C)D^2) = (ND^2)$  complexity. Computing the inverse and determinant of the covariance matrices via the Cholesky decomposition requires  $(D^3)$  per class. Thus, the total com-

putational cost for GDA is  $(ND^2 + CD^3)$ . Evaluating density of a single point: distance from class means requires  $(CD)$ , and matrix vector multiplications requires  $(CD^2)$ . Hence, the total cost for evaluating density on a single point is  $(CD^2)$ .

## D. Experimental Details

### D.1. Dirty-MNIST

We train for 50 epochs using SGD with a momentum of 0.9 and an initial learning rate of 0.1. The learning rate drops by a factor of 10 at training epochs 25 and 40. Following SNGP [45], we apply online spectral normalisation with one step of a power iteration on the convolutional weights. For 1x1 convolutions, we use the exact algorithm, and for 3x3 convolutions, the approximate algorithm from [18]. The coefficient for SN is a hyper-parameter which we set to 3 using cross-validation.

### D.2. OoD Detection Training Setup

We train the softmax baselines on CIFAR-10/100 for 350 epochs using SGD as the optimiser with a momentum of 0.9, and an initial learning rate of 0.1. The learning rate drops by a factor of 10 at epochs 150 and 250. We train the 5-Ensemble baseline using this same training setup. The SNGP and DUQ models were trained using the setup of SNGP and hyper-parameters mentioned in their respective papers [45; 65]. For models trained on ImageNet, we train for 90 epochs with SGD optimizer, an initial learning rate of 0.1 and a weight decay of 1e-4. We use a learning rate warmup decay of 0.01 along with a step scheduler with step size of 30 and a step factor of 0.1.

---

**Algorithm 1** Deep Deterministic Uncertainty

---

1: **Definitions:**

- Regularized feature extractor  $f_\theta : x \rightarrow \mathbb{R}^d$
- Softmax output predictions:  $p(y|x)$
- GMM density:  $q(z) = \sum_y q(z|y = c) q(y = c)$
- Dataset  $(X, Y)$

2: **procedure** TRAIN

- ```
3:   train NN  $p(y|f_\theta(x))$  with  $(X, Y)$ 
4:   for each class  $c$  with samples  $\mathbf{x}_c \subset X$  do
5:      $\mu_c \leftarrow \frac{1}{|\mathbf{x}_c|} \sum_{\mathbf{x}_c} f_\theta(\mathbf{x}_c)$ 
6:      $\Sigma_c \leftarrow \frac{1}{|\mathbf{x}_c|-1} (f_\theta(\mathbf{x}_c) - \mu_c)(f_\theta(\mathbf{x}_c) - \mu_c)^T$ 
7:      $\pi_c \leftarrow \frac{\sum_{\mathbf{x}_c} 1}{|X|}$ 
8:   end for
9: end procedure
```

- ```
10: function DISENTANGLE_UNCERTAINTY(sample  $x$ )
11:   compute feature representation  $z = f_\theta(x)$ 
12:   compute density under GMM:  $q(z) = \sum_y q(z|y) q(y)$  with  $q(z|y) \sim \mathcal{N}(\mu_y; \Sigma_y)$ ,  $q(y) = \pi_y$ 
13:   compute softmax entropy:  $H_p[Y|x]$ 
14:   if low density  $q(z)$  then
15:     return  $(-q(z), )$ 
16:   else if high density  $q(z)$  then
17:     if high entropy  $H_p[Y|x]$  then
18:       return ambiguous iD
19:     else if low entropy  $H_p[Y|x]$  then
20:       return  $(0, H_p[Y|x])$ 
21:     end if
22:   end if
23: end function
```
- 

### D.3. Semantic Segmentation Training Setup

In Figure 9, we plot the L2 distance between feature space means of different classes for a pair of randomly chosen distant pixels on the Pascal VOC 2012 val set. We observe that feature space means between pairs of different classes are more distant compared to the same class irrespective of the location of the pixel for the class. This leads us to construct a Gaussian mean and covariance matrix per class as opposed to one mean and one covariance matrix per class per pixel, thereby greatly reducing the computational load of fitting a GMM in semantic segmentation. Similar to classification, we treat each pixel in the training set as a separate sample and fit a single Gaussian mean and covariance matrix per class.

For the semantic segmentation experiment, we use a DeepLab-v3+ [3] model with a ResNet-101 backbone as the architecture of choice. We train each of the models on Pascal VOC for 50 epochs using SGD as the optimizer, with

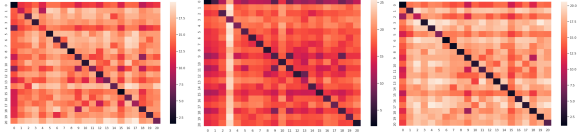


Figure 9. L2 distances between the feature space means of different classes for a pair of distant pixels on the Pascal VOC 2012 val set: (left) Pixels (10, 255) and (500, 255), (middle) Pixels (234, 349) and (36, 22) and (right) Pixels (300, 500) and (400, 255).

a momentum of 0.9 and a weight decay of  $5e-4$ . We set the initial learning rate to 0.007 with a polynomial decay during the course of training. Finally, we trained with a batch size of 32 parallelized over 4 GPUs.

### D.4. Compute Resources

Each model (ResNet-18, Wide-ResNet-28-10, ResNet-50, ResNet-110, DenseNet-121 or VGG-16) used for the large scale active learning, CIFAR-10 vs SVHN/CIFAR-100/Tiny-ImageNet/CIFAR-10-C and CIFAR-100 vs SVHN/Tiny-ImageNet tasks was trained on a single Nvidia Quadro RTX 6000 GPU. Each model (LeNet, VGG-16 and ResNet-18) used to get the results in Fig. 1 and Tab. 11 was trained on a single Nvidia GeForce RTX 2060 GPU. Each model (ResNet-50, Wide-ResNet-50-2, VGG-16) trained on ImageNet was trained using 8 Nvidia Quadro RTX 6000 GPUs.

### E. Additional Results

In this section, we provide details of additional results on the OoD detection task using CIFAR-10 vs SVHN/CIFAR-100/Tiny-ImageNet/CIFAR-10-C and CIFAR-100 vs SVHN/Tiny-ImageNet for ResNet-50, ResNet-110 and DenseNet-121 architectures. We present results on ResNet-50, ResNet-110 and DenseNet-121 for CIFAR-10 vs SVHN/CIFAR-100/Tiny-ImageNet and CIFAR-100 vs SVHN/Tiny-ImageNet in Tab. 5, Tab. 6 and Tab. 7 respectively. We also present results on individual corruption types for CIFAR-10-C for Wide-ResNet-28-10, ResNet-50, ResNet-110 and DenseNet-121 in Fig. 10, Fig. 11, Fig. 12 and Fig. 13 respectively.

Finally, we provide results for various ablations on DDU. As mentioned in §3, DDU consists of a deterministic softmax model trained with appropriate inductive biases. It uses softmax entropy to quantify aleatoric uncertainty and feature-space density to quantify epistemic uncertainty. In the ablation, we try to experimentally evaluate the following scenarios:

1. **Effect of inductive biases (sensitivity + smoothness):**

We want to see the effect of removing the proposed inductive biases (i.e. no sensitivity and smoothness constraints) on the OoD detection performance of a model. To do this, we train a VGG-16 with and without spectral nor-

Table 5. *OoD detection performance of different baselines using a ResNet-50 architecture with the CIFAR-10 vs SVHN/CIFAR-100/Tiny-ImageNet and CIFAR-100 vs SVHN/Tiny-ImageNet dataset pairs averaged over 25 runs.* Note: SN stands for Spectral Normalisation, JP stands for Jacobian Penalty. We highlight the best deterministic and best method overall in bold for each metric.

Train Dataset	Method	Penalty	Aleatoric Uncertainty	Epistemic Uncertainty	Test Accuracy (↑)	Test ECE (↓)	AUROC SVHN (↑)	AUROC CIFAR-100 (↑)	AUROC Tiny-ImageNet (↑)
CIFAR-10	Softmax	-		Softmax Entropy			93.80 ± 0.41	88.91 ± 0.07	88.32 ± 0.07
	Energy-based [46]	-	Softmax Entropy	Softmax Density	<b>95.04 ± 0.05</b>	<b>0.97 ± 0.04</b>	94.48 ± 0.44	88.84 ± 0.08	88.45 ± 0.08
	DUQ [65]	JP	Kernel Distance	Kernel Distance	94.05 ± 0.11	1.71 ± 0.07	93.14 ± 0.43	83.87 ± 0.27	84.28 ± 0.26
	SNGP [45]	SN	Predictive Entropy	Predictive Entropy	94.90 ± 0.11	1.01 ± 0.03	93.15 ± 0.85	89.32 ± 0.10	88.96 ± 0.13
	<b>DDU (ours)</b>	SN	Softmax Entropy	GMM Density	94.92 ± 0.06	1 ± 0.04	<b>94.77 ± 0.35</b>	<b>89.98 ± 0.17</b>	<b>89.12 ± 0.13</b>
	5-Ensemble [40]	-	Predictive Entropy	Predictive Entropy Mutual Information	<b>96.06 ± 0.04</b>	1.65 ± 0.07	94.75 ± 0.39	89.87 ± 0.06	88.69 ± 0.05
							94.09 ± 0.20	89.76 ± 0.06	89.04 ± 0.03
					Test Accuracy (↑)	Test ECE (↓)	AUROC SVHN (↑)	AUROC CIFAR-100 (↑)	AUROC Tiny-ImageNet (↑)
CIFAR-100	Softmax	-		Softmax Entropy			81.32 ± 0.65		79.83 ± 0.07
	Energy-based [46]	-	Softmax Entropy	Softmax Density	77.91 ± 0.09	4.32 ± 0.10	82.05 ± 0.69		79.61 ± 0.08
	DUQ [65]	JP	Kernel Distance	Kernel Distance	74.73 ± 0.22	7.68 ± 0.13	82.50 ± 2.09		77.05 ± 0.16
	SNGP [45]	SN	Predictive Entropy	Predictive Entropy					
	<b>DDU (ours)</b>	SN	Softmax Entropy	GMM Density	<b>79.26 ± 0.16</b>	<b>4.07 ± 0.06</b>	<b>87.34 ± 0.64</b>		<b>82.11 ± 0.20</b>
5-Ensemble [40]	-	Predictive Entropy	Predictive Entropy Mutual Information	<b>81.06 ± 0.07</b>	<b>3.54 ± 0.12</b>	83.42 ± 0.89		77.69 ± 0.12	
							84.24 ± 0.90	89.76 ± 0.06	81.59 ± 0.05

Table 6. *OoD detection performance of different baselines using a ResNet-110 architecture with the CIFAR-10 vs SVHN/CIFAR-100/Tiny-ImageNet and CIFAR-100 vs SVHN/Tiny-ImageNet dataset pairs averaged over 25 runs.* Note: SN stands for Spectral Normalisation, JP stands for Jacobian Penalty. We highlight the best deterministic and best method overall in bold for each metric.

Train Dataset	Method	Penalty	Aleatoric Uncertainty	Epistemic Uncertainty	Test Accuracy (↑)	Test ECE (↓)	AUROC SVHN (↑)	AUROC CIFAR-100 (↑)	AUROC Tiny-ImageNet (↑)
CIFAR-10	Softmax	-		Softmax Entropy			93.12 ± 0.44	88.7 ± 0.1	88.07 ± 0.11
	Energy-based [46]	-	Softmax Entropy	Softmax Density	<b>95.08 ± 0.04</b>	1.02 ± 0.04	93.67 ± 0.47	88.60 ± 0.11	88.13 ± 0.11
	DUQ [65]	JP	Kernel Distance	Kernel Distance	94.32 ± 0.17	1.21 ± 0.07	94.02 ± 0.45	86.17 ± 0.35	85.24 ± 0.21
	SNGP [45]	SN	Predictive Entropy	Predictive Entropy	94.85 ± 0.09	1.04 ± 0.02	93.17 ± 0.53	89.23 ± 0.10	88.80 ± 0.12
	<b>DDU (ours)</b>	SN	Softmax Entropy	GMM Density	94.82 ± 0.06	<b>1.01 ± 0.04</b>	<b>95.48 ± 0.30</b>	<b>90.08 ± 0.13</b>	<b>89.18 ± 0.15</b>
	5-Ensemble [40]	-	Predictive Entropy	Predictive Entropy Mutual Information	<b>96.18 ± 0.05</b>	1.57 ± 0.05	95.07 ± 0.45	<b>90.23 ± 0.04</b>	89 ± 0.03
							94.72 ± 0.34	89.69 ± 0.05	88.35 ± 0.05
					Test Accuracy (↑)	Test ECE (↓)	AUROC SVHN (↑)	AUROC CIFAR-100 (↑)	AUROC Tiny-ImageNet (↑)
CIFAR-100	Softmax	-		Softmax Entropy			82.04 ± 0.57		80.13 ± 0.07
	Energy-based [46]	-	Softmax Entropy	Softmax Density	78.65 ± 0.10	3.93 ± 0.13	82.78 ± 0.60		80.01 ± 0.09
	DUQ [65]	JP	Kernel Distance	Kernel Distance	76.16 ± 0.27	6.43 ± 0.75	83.94 ± 0.10		78.54 ± 0.28
	SNGP [45]	SN	Predictive Entropy	Predictive Entropy					
	<b>DDU (ours)</b>	SN	Softmax Entropy	GMM Density	<b>78.89 ± 0.17</b>	<b>3.79 ± 0.07</b>	<b>88.66 ± 0.56</b>		<b>82.58 ± 0.24</b>
5-Ensemble [40]	-	Predictive Entropy	Predictive Entropy Mutual Information	<b>81.80 ± 0.10</b>	<b>3.67 ± 0.11</b>	83.68 ± 0.33		81.12 ± 0.13	
							85.11 ± 0.57	89.76 ± 0.06	81.94 ± 0.06

Table 7. *OoD detection performance of different baselines using a DenseNet-121 architecture with the CIFAR-10 vs SVHN/CIFAR-100/Tiny-ImageNet and CIFAR-100 vs SVHN/Tiny-ImageNet dataset pairs averaged over 25 runs.* Note: SN stands for Spectral Normalisation, JP stands for Jacobian Penalty. We highlight the best deterministic and best method overall in bold for each metric.

Train Dataset	Method	Penalty	Aleatoric Uncertainty	Epistemic Uncertainty	Test Accuracy (↑)	Test ECE (↓)	AUROC SVHN (↑)	AUROC CIFAR-100 (↑)	AUROC Tiny-ImageNet (↑)
CIFAR-10	Softmax	-		Softmax Entropy			94 ± 0.44	87.55 ± 0.11	86.99 ± 0.12
	Energy-based [46]	-	Softmax Entropy	Softmax Density	95.16 ± 0.03	1.10 ± 0.04	94.07 ± 0.54	86.73 ± 0.15	86.43 ± 0.16
	DUQ [65]	JP	Kernel Distance	Kernel Distance	95.02 ± 0.14	1.08 ± 0.08	94.67 ± 0.41	87.38 ± 0.21	86.72 ± 0.14
	SNGP [45]	SN	Predictive Entropy	Predictive Entropy	94.31 ± 0.21	1.08 ± 0.10	94.48 ± 0.34	88.86 ± 0.46	88.40 ± 0.48
	<b>DDU (ours)</b>	SN	Softmax Entropy	GMM Density	<b>95.21 ± 0.03</b>	<b>1.05 ± 0.03</b>	<b>96.21 ± 0.31</b>	<b>90.84 ± 0.06</b>	<b>89.70 ± 0.06</b>
	5-Ensemble [40]	-	Predictive Entropy	Predictive Entropy Mutual Information	<b>96.18 ± 0.05</b>	1.07 ± 0.07	95.78 ± 0.11	90.65 ± 0.03	89.62 ± 0.06
							95.75 ± 0.10	90.71 ± 0.04	89.34 ± 0.06
					Test Accuracy (↑)	Test ECE (↓)	AUROC SVHN (↑)	AUROC CIFAR-100 (↑)	AUROC Tiny-ImageNet (↑)
CIFAR-100	Softmax	-		Softmax Entropy			85.86 ± 0.42		81.10 ± 0.07
	Energy-based [46]	-	Softmax Entropy	Softmax Density	79.02 ± 0.08	4.11 ± 0.08	87.09 ± 0.49		80.84 ± 0.08
	DUQ [65]	JP	Kernel Distance	Kernel Distance	79.15 ± 0.15	6.73 ± 0.10	85.00 ± 0.12		79.76 ± 0.15
	SNGP [45]	SN	Predictive Entropy	Predictive Entropy					
	<b>DDU (ours)</b>	SN	Softmax Entropy	GMM Density	<b>79.15 ± 0.07</b>	<b>4.11 ± 0.06</b>	<b>88.44 ± 0.55</b>		<b>81.85 ± 0.11</b>
5-Ensemble [40]	-	Predictive Entropy	Predictive Entropy Mutual Information	<b>81.01 ± 0.13</b>	4.81 ± 0.05	88.32 ± 0.61		81.45 ± 0.12	
							88.36 ± 0.17	89.76 ± 0.06	81.73 ± 0.06

malisation. Note that VGG-16 does not have residual connections and hence, a VGG-16 does not follow the sensitivity and smoothness (bi-Lipschitz) constraints.

- Effect of sensitivity alone:** Since residual connections make a model sensitive to changes in the input space by lower bounding its Lipschitz constant, we also want to see how a network performs with just the sensitivity con-

straint alone. To observe this, we train a Wide-ResNet-28-10 without spectral normalisation (i.e. no explicit upper bound on the Lipschitz constant of the model).

- Metrics for aleatoric and epistemic uncertainty:** With the above combinations, we try to observe how different metrics for aleatoric and epistemic uncertainty perform. To quantify aleatoric uncertainty, we use the softmax

Table 8. OoD detection performance of different ablations trained on CIFAR-10 using Wide-ResNet-28-10 and VGG-16 architectures with SVHN, CIFAR-100 and Tiny-ImageNet as OoD datasets averaged over 25 runs. Note: SN stands for Spectral Normalisation. We highlight the best deterministic and best method overall in bold for each metric.

Architecture	Ablations		Aleatoric Uncertainty		Epistemic Uncertainty		Test Accuracy (↑)	Test ECE (↓)	AUROC SVHN (↑)	AUROC CIFAR-100 (↑)	AUROC Tiny-ImageNet (↑)
	Ensemble	Residual Connections	SN	GMM							
Wide-ResNet-28-10	x	✓	x	x	Softmax Entropy	Softmax Entropy	<b>95.98 ± 0.02</b>	0.85 ± 0.02	94.44 ± 0.43	89.39 ± 0.06	88.42 ± 0.05
				✓	Softmax Entropy	GMM Density	95.98 ± 0.02	0.85 ± 0.02	96.08 ± 0.25	90.94 ± 0.03	90.62 ± 0.05
				x	Softmax Entropy	Softmax Entropy	95.97 ± 0.03	0.85 ± 0.04	94.05 ± 0.26	90.02 ± 0.07	89.07 ± 0.06
				✓	Softmax Entropy	Softmax Entropy	94.31 ± 0.33		94.31 ± 0.33	89.78 ± 0.08	88.96 ± 0.07
				✓	<b>Softmax Entropy</b>	<b>GMM Density</b>	95.97 ± 0.03	<b>0.85 ± 0.04</b>	<b>97.86 ± 0.19</b>	<b>91.34 ± 0.04</b>	<b>91.07 ± 0.05</b>
✓	✓	x	x	Predictive Entropy	Predictive Entropy	<b>96.59 ± 0.02</b>	<b>0.76 ± 0.03</b>	97.73 ± 0.31	<b>92.13 ± 0.02</b>	90.06 ± 0.03	
VGG-16	x	✓	x	x	Softmax Entropy	Softmax Entropy	93.63 ± 0.04	1.64 ± 0.03	85.76 ± 0.84	82.48 ± 0.14	83.07 ± 0.12
				✓	Softmax Entropy	GMM Density	93.63 ± 0.04	1.64 ± 0.03	84.24 ± 1.04	81.91 ± 0.17	82.82 ± 0.14
				x	Softmax Entropy	Softmax Entropy	93.63 ± 0.04	1.78 ± 0.04	89.25 ± 0.36	86.55 ± 0.10	86.78 ± 0.09
				✓	Softmax Entropy	Softmax Entropy	93.62 ± 0.04	1.78 ± 0.04	87.54 ± 0.41	82.71 ± 0.09	83.33 ± 0.08
				✓	Softmax Entropy	GMM Density	93.62 ± 0.04	1.78 ± 0.04	86.28 ± 0.51	82.15 ± 0.11	83.07 ± 0.10
✓	✓	x	x	Predictive Entropy	Predictive Entropy	92.80 ± 0.18	2.03 ± 0.03	89.01 ± 0.08	87.66 ± 0.08	87.66 ± 0.08	
					Mutual Information	91 ± 0.22		88.43 ± 0.08	88.74 ± 0.05		

Table 9. OoD detection performance of different ablations trained on CIFAR-100 using Wide-ResNet-28-10 and VGG-16 architectures with SVHN and Tiny-ImageNet as the OoD dataset averaged over 25 runs. Note: SN stands for Spectral Normalisation. We highlight the best deterministic and best method overall in bold for each metric.

Architecture	Ablations		Aleatoric Uncertainty		Epistemic Uncertainty		Test Accuracy (↑)	Test ECE (↓)	AUROC SVHN (↑)	AUROC Tiny-ImageNet (↑)
	Ensemble	Residual Connections	SN	GMM						
Wide-ResNet-28-10	x	✓	x	x	Softmax Entropy	Softmax Entropy	80.26 ± 0.06	4.62 ± 0.06	77.42 ± 0.57	81.53 ± 0.05
				✓	Softmax Entropy	GMM Density	80.26 ± 0.06	4.62 ± 0.06	78.00 ± 0.63	81.33 ± 0.06
				x	Softmax Entropy	Softmax Entropy	80.98 ± 0.06	4.10 ± 0.08	85.37 ± 0.36	82.57 ± 0.03
				✓	Softmax Entropy	Softmax Entropy	86.41 ± 0.38		86.41 ± 0.38	82.49 ± 0.04
				✓	<b>Softmax Entropy</b>	<b>GMM Density</b>	<b>80.98 ± 0.06</b>	<b>4.10 ± 0.08</b>	<b>87.53 ± 0.62</b>	<b>83.13 ± 0.06</b>
✓	✓	x	x	Predictive Entropy	Predictive Entropy	<b>82.79 ± 0.10</b>	<b>3.32 ± 0.09</b>	79.54 ± 0.91	82.95 ± 0.09	
VGG-16	x	✓	x	x	Softmax Entropy	Softmax Entropy	73.48 ± 0.05	4.46 ± 0.05	76.73 ± 0.72	76.43 ± 0.05
				✓	Softmax Entropy	GMM Density	73.48 ± 0.05	4.46 ± 0.05	77.70 ± 0.86	74.68 ± 0.07
				x	Softmax Entropy	Softmax Entropy	73.58 ± 0.06	4.32 ± 0.06	75.65 ± 0.95	74.32 ± 1.73
				✓	Softmax Entropy	Softmax Entropy	77.21 ± 0.77		77.21 ± 0.77	76.59 ± 0.06
				✓	Softmax Entropy	GMM Density	73.58 ± 0.06	4.32 ± 0.06	77.76 ± 0.90	74.86 ± 0.08
✓	✓	x	x	Predictive Entropy	Predictive Entropy	77.84 ± 0.11	5.32 ± 0.10	75.99 ± 1.23	74.06 ± 1.67	
					Mutual Information	79.62 ± 0.73		72.07 ± 0.48	78.66 ± 0.06	76.27 ± 0.05

Table 10. LDA vs GDA ablation for OoD detection performance using Wide-ResNet-50-2, ResNet-50, Wide-ResNet-50-2 architectures (depending on dataset) on CIFAR-10 vs SVHN/CIFAR-100/TinyImageNet, CIFAR-100 vs SVHN/TinyImageNet, and ImageNet vs ImageNet-O [26]. Best AUROC (↑) scores are marked in bold. GDA performs better, except with SVHN as OoD dataset.

Model ID	WRN-28-10 CIFAR-10			WRN-28-10 CIFAR-100		WRN-50-2 ImageNet	RN-50 ImageNet-O
	SVHN	CIFAR-100	TinyImageNet	SVHN	TinyImageNet		
LDA (Maha [43])	<b>98.41 ± 0.09</b>	82.90 ± 0.23	82.48 ± 0.25	<b>92.53 ± 0.62</b>	68.86 ± 0.13	64.19 ± 0.23	61.68 ± 0.14
GDA (DDU, ours)	97.86 ± 0.19	<b>91.34 ± 0.04</b>	<b>91.07 ± 0.05</b>	87.53 ± 0.62	<b>83.13 ± 0.06</b>	<b>73.12 ± 0.19</b>	<b>71.29 ± 0.08</b>

entropy of the model. On the other hand, to quantify the epistemic uncertainty, we use **i)** the softmax entropy, **ii)** the softmax density [46] or **iii)** the GMM feature density (as described in §3).

For the purposes of comparison, we also present scores obtained by a 5-Ensemble of the respective architectures (i.e. Wide-ResNet-28-10 and VGG-16) in Tab. 8 for CIFAR-10 vs SVHN/CIFAR-100 and in Tab. 9 for CIFAR-100 vs SVHN. Based on these results, we can make the following observations (in addition to the ones we make in §4.2):

**Inductive biases are important for feature density.**

From the AUROC scores in Tab. 8, we can see that using the feature density of a GMM in VGG-16 without the proposed inductive biases yields significantly lower AUROC scores as compared to Wide-ResNet-28-10 with inductive biases. In fact, in none of the datasets is the feature density of a VGG able to outperform its corresponding ensemble. This provides yet more evidence (in addition to Fig. 1) to show that the GMM feature density alone cannot estimate epistemic uncertainty in a model that suffers from feature collapse. We need sensitivity and smoothness conditions (see §5) on the feature space of the model to obtain feature densities that

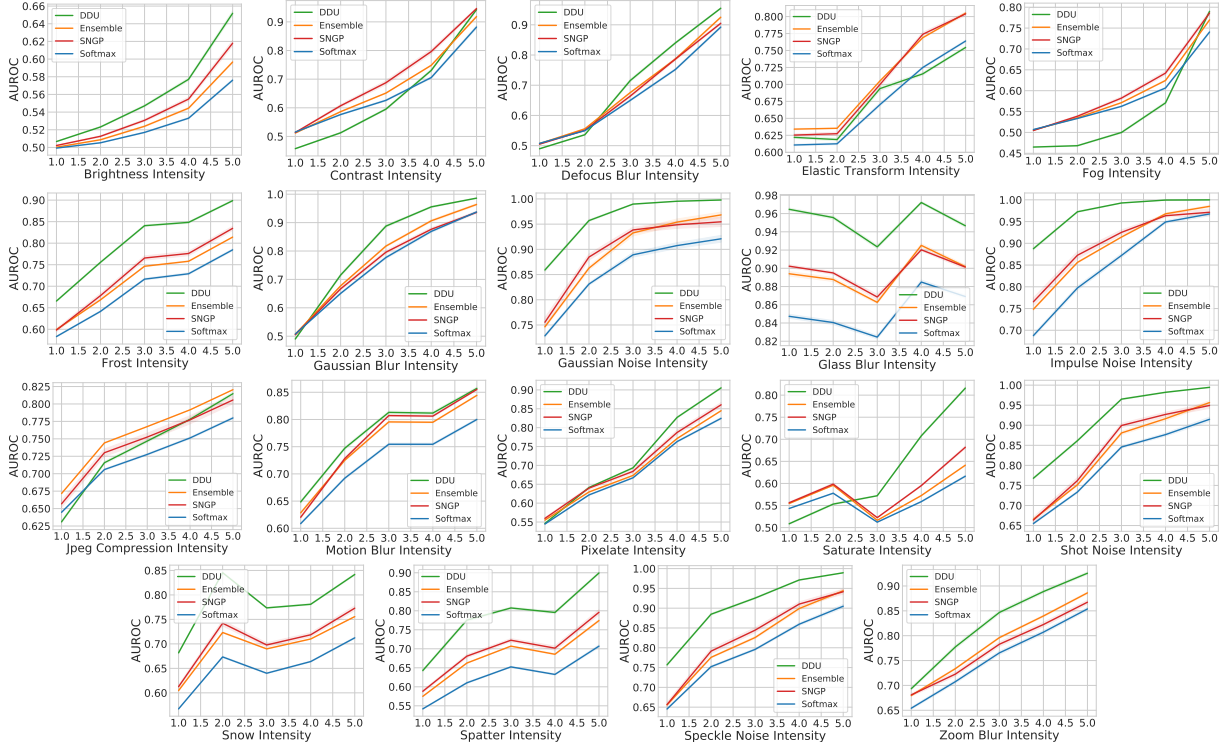


Figure 10. AUROC vs corruption intensity for all corruption types in CIFAR-10-C with Wide-ResNet-28-10 as the architecture and baselines: Softmax Entropy, Ensemble (using Predictive Entropy as uncertainty), SNGP and DDU feature density.

capture epistemic uncertainty.

**Sensitivity creates a bigger difference than smoothness.** We note that the difference between AUROC obtained from feature density between Wide-ResNet-28-10 models with and without spectral normalisation is minimal. Although Wide-ResNet-28-10 with spectral normalisation (i.e. smoothness constraints) still outperforms its counterpart without spectral normalisation, the small difference between the AUROC scores indicates that it might be the residual connections (i.e. sensitivity constraints) that make the model detect OoD samples better. This observation is also intuitive as a sensitive feature extractor should map OoD samples farther from iD ones.

**DDU as a simple baseline.** In DDU, we use the softmax output of a model to get aleatoric uncertainty. We use the GMM’s feature-density to estimate the epistemic uncertainty. Hence, DDU does not suffer from miscalibration as the softmax outputs can be calibrated using post-hoc methods like temperature scaling. At the same time, the feature-densities of the model are not affected by temperature scaling and capture epistemic uncertainty well.

## F. Additional Ablations & Toy Experiments

Here, we provide details for toy experiments from the main paper which are visualized in Fig. 1, Fig. 15 and Fig. 2.

Table 11. *ECE* for *Dirty-MNIST* test set and *AUROC* for *Dirty-MNIST* vs *Fashion-MNIST* as proxies for aleatoric and epistemic uncertainty quality respectively.

Model	ECE	AUROC for Softmax Entropy ( $\uparrow$ )	AUROC for Feature Density ( $\uparrow$ )
LeNet	2.22	84.23	71.41
VGG-16	<b>2.11</b>	84.04	89.01
ResNet-18+SN (DDU)	2.34	83.01	<b>99.91</b>

## F.1. Motivational Example in Figure 1

In Fig. 1 we train a LeNet [41], a VGG-16 [60] and a ResNet-18 with spectral normalisation [22; 52] (ResNet-18+SN) on *Dirty-MNIST*, a modified version of MNIST [41] with additional ambiguous digits (Ambiguous-MNIST). *Ambiguous-MNIST* contains samples with multiple plausible labels and thus higher aleatoric uncertainty (see Fig. 1a). We refer to §B for details on how this dataset was generated. With ambiguous data having various levels of aleatoric uncertainty, *Dirty-MNIST* is more representative of real-world datasets compared to well-cleaned curated datasets, like MNIST and CIFAR-10, commonly used for benchmarking in ML [13; 39]. Moreover, *Dirty-MNIST* also poses a challenge for recent uncertainty estimation methods, which often confound aleatoric and epistemic uncertainty [65]. Figure 1b shows that the softmax entropy of a deterministic

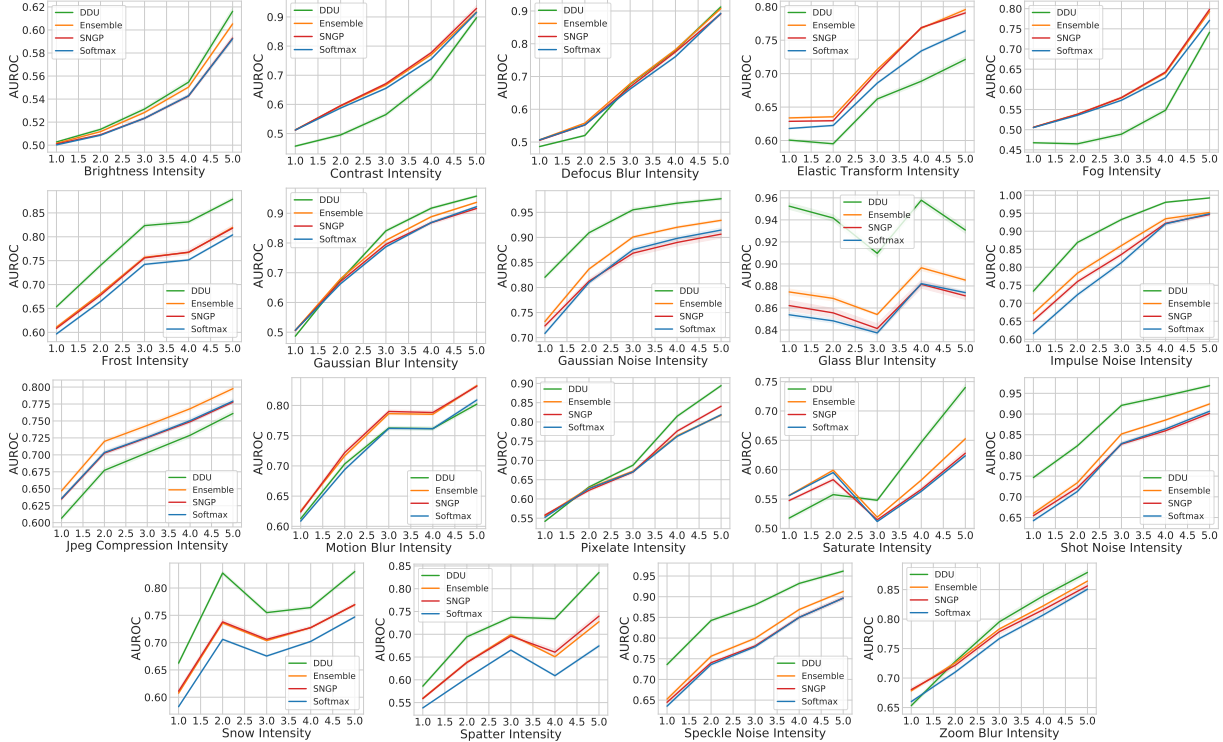


Figure 11. AUROC vs corruption intensity for all corruption types in CIFAR-10-C with ResNet-50 as the architecture and baselines: Softmax Entropy, Ensemble (using Predictive Entropy as uncertainty), SNGP and DDU feature density.

model is unable to distinguish between iD (Dirty-MNIST) and OoD (Fashion-MNIST [70]) samples as the entropy for the latter heavily overlaps with the entropy for Ambiguous-MNIST samples. However, the feature-space density of the model with our inductive biases in Fig. 1c captures epistemic uncertainty reliably and is able to distinguish iD from OoD samples. The same cannot be said for LeNet or VGG in Fig. 1c, whose densities are unable to separate OoD from iD samples. This demonstrates the importance of the inductive bias to ensure the sensitivity and smoothness of the feature space as we further argue below. Finally, Fig. 1b and Fig. 1c demonstrate that our method is able to separate aleatoric from epistemic uncertainty: samples with low feature density have high epistemic uncertainty, whereas those with both high feature density and high softmax entropy have high aleatoric uncertainty—note the high softmax entropy for the most ambiguous Ambiguous-MNIST samples in Fig. 1b.

### F.1.1 Disentangling Epistemic and Aleatoric Uncertainty

We used a simple example in §1 to demonstrate that a single softmax model with a proper inductive bias can reliably capture epistemic uncertainty via its feature-space density and aleatoric uncertainty via its softmax entropy. To recreate

the natural characteristics of uncurated real-world datasets, which contain ambiguous samples, we use MNIST [41] as an iD dataset of unambiguous samples, Ambiguous-MNIST as an iD dataset of ambiguous samples and Fashion-MNIST [70] as an OoD dataset (see Fig. 1a), with more details in §B. We train a LeNet [41], a VGG-16 [60] and a ResNet-18 [22] with spectral normalisation (SN) on Dirty-MNIST (a mix of Ambiguous- and standard MNIST) with the training setup detailed in §D.1.

Table 11 gives a quantitative evaluation of the qualitative results in §1. The AUROC metric reflects the quality of the epistemic uncertainty as it measures the probability that iD and OoD samples can be distinguished, and OoD samples are never seen during training while iD samples are semantically similar to training samples. The ECE metric measures the quality of the aleatoric uncertainty. The softmax outputs capture aleatoric uncertainty well, as expected, and all 3 models obtain similar ECE scores on the Dirty-MNIST test set. However, with an AUROC of around 84% for all the 3 models, on Dirty-MNIST vs Fashion-MNIST, we conclude that softmax entropy is unable to capture epistemic uncertainty well. This is reinforced in Fig. 1b, which shows a strong overlap between the softmax entropy of OoD and ambiguous iD samples. At the same time, the feature-space densities of LeNet and VGG-16, with AUROC scores around

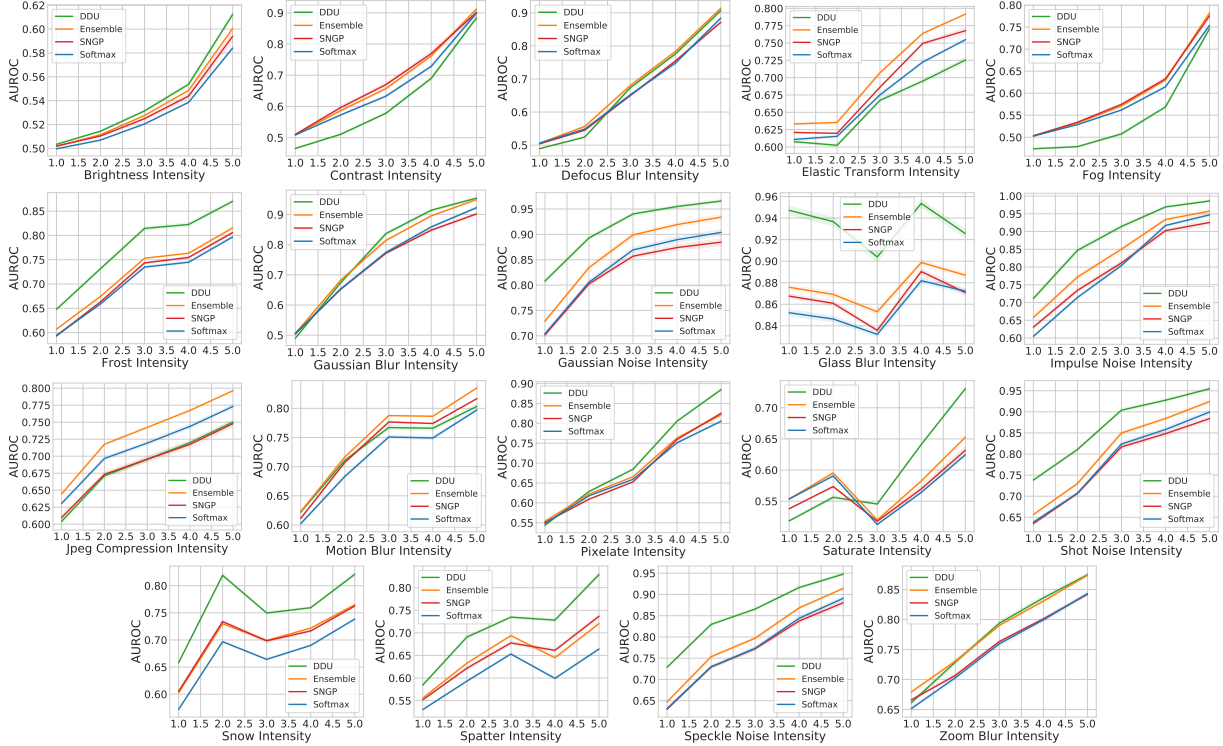


Figure 12. AUROC vs corruption intensity for all corruption types in CIFAR-10-C with ResNet-110 as the architecture and baselines: Softmax Entropy, Ensemble (using Predictive Entropy as uncertainty), SNGP and DDU feature density.

71% and 89% respectively, are unable to distinguish OoD from iD samples, indicating that simply using feature-space density without appropriate inductive biases (as seen in [43]) is not sufficient.

Only by fitting a GMM on top of a feature extractor with appropriate inductive biases (DDU) and using its feature density are we able to obtain performance far better (with AUROC of 99.9%) than the alternatives in the ablation study (see Tab. 11, also noticeable in Fig. 1c). The entropy of a softmax model can capture aleatoric uncertainty, even without additional inductive biases, but it *cannot* be used to estimate epistemic uncertainty (see §5). On the other hand, feature-space density can *only* be used to estimate epistemic uncertainty *when the feature extractor is sensitive and smooth*, as achieved by using a ResNet and spectral normalisation in DDU.

## F.2. Effects of a well-regularized feature space on feature collapse

The effects of a well-regularized feature space on feature collapse are visible in Fig. 1. In the case of feature collapse, we must have *some* OoD inputs for which the features are mapped on top of the features of iD inputs. The distances of these OoD features to each class centroid must be equal to the distances of the corresponding iD inputs to class centroids,

and hence the density for these OoD inputs must be equal to the density of the iD inputs. If the density histograms do not overlap, no feature collapse can be present<sup>3</sup>. We see no overlapping densities in Fig. 1(c, right), therefore we indeed have no feature collapse. First, the effects of having a well-regularized feature space on feature collapse can be analysed from Fig. 1. In case of feature collapse we must have *some* OoD features mapped on top of iD features, therefore the distances of at least some OoD features to the class centroids must be equal to iD’s distances, hence OoD density must overlap with iD density. We see this in Fig. 1(c) (left) in the case without a regularized feature space. On the other hand, when we regularize the feature space, we see the densities do not overlap, i.e. the distances of the features of OoD examples to the centroids are larger than the distances of iD examples (Fig. 1(c) right), hence feature collapse is not present.

## F.3. QUBIQ Challenge

In this section, we evaluate DDU’s performance on the real-world QUBIQ challenge related to biomedical imaging. QUBIQ has a total of 7 binary segmentation tasks in 4

<sup>3</sup>Note though that the opposite (‘if the density histograms overlap then there must be feature collapse’) needs not hold: the histograms can also overlap due to other reasons.



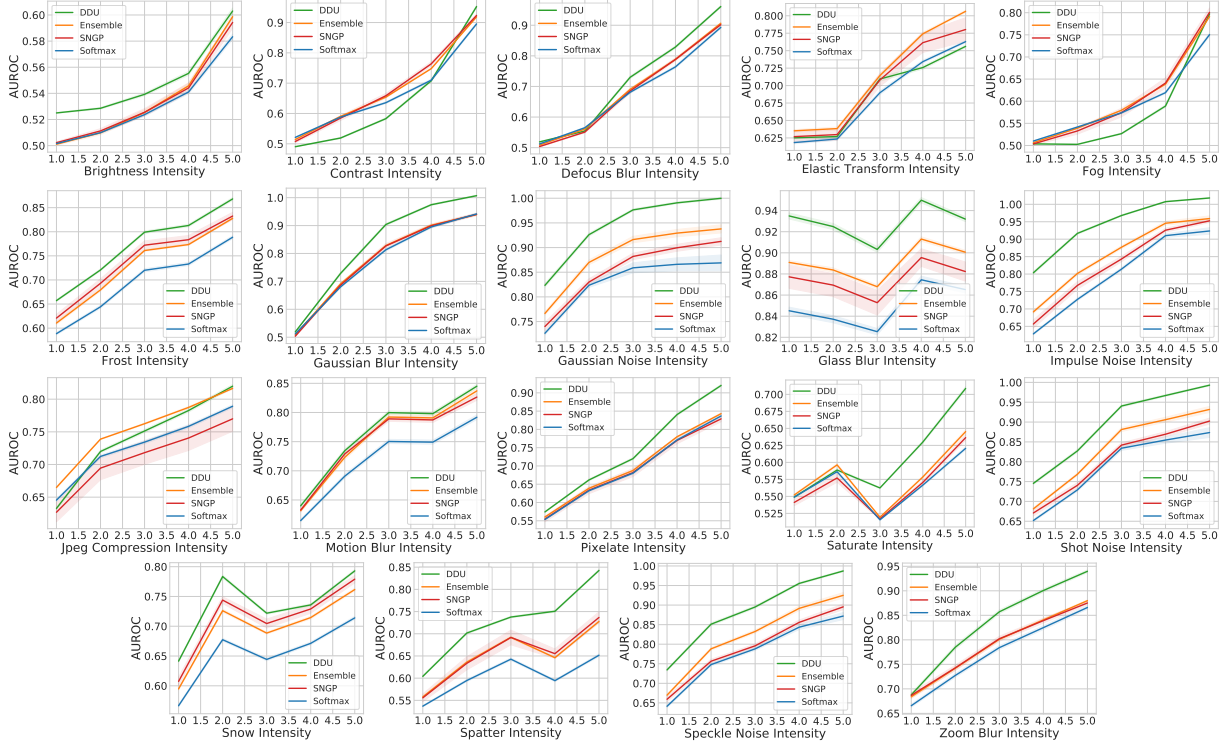


Figure 13. AUROC vs corruption intensity for all corruption types in CIFAR-10-C with DenseNet-121 as the architecture and baselines: Softmax Entropy, Ensemble (using Predictive Entropy as uncertainty), SNGP and DDU feature density.

Softmax	Energy	3-Ensemble PE	DDU
78.4 ± 1.31	77.31 ± 1.5	82.25 ± 0.83	<b>82.63 ± 1.08</b>

Table 12. Dice scores for the QUBIQ 2021 challenge

biomedical imaging datasets with multiple annotations per image. The task is to predict the distribution of source labels with a mask of values between 0 and 1. For evaluation, the annotations are averaged to provide a continuous ground-truth. The prediction mask and continuous ground-truth are binarized by thresholding between  $[0, 1]$  and a Dice score is computed between the resulting binary masks. The average dice score across thresholds, images and tasks is reported. Note that in the continuous ground-truth, 0.5 indicates maximum uncertainty and values above or below indicate lower uncertainty. Thus, for our comparison, we scale all uncertainty values to the range  $u \in [0, 0.5]$  and use  $p + u$  if  $p = 0$  and  $p - u$  if  $p = 1$ , where  $p$  is the binary prediction. We use a UNet model with a ResNet encoder and report the Dice scores averaged over 5 runs in Tab. 12. Even on this real-world dataset, DDU performs as well as ensembles and outperforms Softmax and Energy baselines.

#### F.4. Additional Baselines

In this section, we provide an ablation with additional baselines on OoD detection for comparison with DDU. In particular, we provide comparisons with Feature Space Sin-

gularity (FSSD) [31], Batch Ensemble (BE) [67] and SWAG [49] using Wide-ResNet-28-10 as additional recent baselines. We also use the Wide-ResNet-28-10 feature extractor trained using SNGP loss and fit DDU (i.e., GDA) on its feature space. Since SNGP also uses a sensitive smooth feature space with residual connections and spectral normalization, its feature space makes for a good candidate to apply DDU. In Tab. 13, we provide the AUROC scores for models trained on CIFAR-10 (C10) and CIFAR-100 (C100). Broadly, DDU outperforms all competitive baselines. Additionally, we observe a broad improvement in AUROC when DDU is applied on the SNGP feature extractor as compared to vanilla SNGP. However, DDU on a feature extractor trained using softmax loss is still superior to DDU on the SNGP feature extractor.

#### F.5. Additional Calibration Metrics

In Tab. 14, in addition to ECE%, we provide additional calibration error scores: temperature scaled Thresholded Adaptive Calibration Error (TACE) % and Negative Log Likelihood (NLL) for Wide-ResNet-28-10 trained on CIFAR-10/100 (main results for this can be found in Tab. 1). The results for TACE and NLL are consistent with what we see for ECE. Ensembles produce the most calibrated models and among deterministic baselines, DDU is the best calibrated.

Train dataset	Method	AUROC (†)		
		SVHN	CIFAR-100/10	Tiny-ImageNet
C10	FSSD [31]	97.24	89.88	90.23
	BE [67]	95.36	87.63	88.14
	SWAG [49]	96.37	90.33	90.24
	SNGP [45]	94.0 ± 1.3	91.13 ± 0.15	89.97 ± 0.19
	5-Ensemble [40]	97.73 ± 0.31	<b>92.13 ± 0.02</b>	90.06 ± 0.03
	SNGP + DDU	96.47 ± 0.7	89.97 ± 0.13	90.3 ± 0.12
	<b>DDU (Ours)</b>	<b>97.86 ± 0.19</b>	91.34 ± 0.04	<b>91.07 ± 0.05</b>
C100		SVHN		Tiny-ImageNet
	FSSD [31]	<b>87.64</b>		82.2
	BE [67]	86.44		78.33
	SWAG [49]	81.41		81.67
	SNGP [45]	85.71 ± 0.81		78.85 ± 0.43
	5-Ensemble [40]	79.54 ± 0.91		82.95 ± 0.09
	SNGP + DDU	87.34 ± 0.76		79.62 ± 0.36
<b>DDU(Ours)</b>	87.53 ± 0.62		<b>83.13 ± 0.06</b>	

Table 13. OoD detection ablation with WRN-28-10 model with additional baselines, FSSD [31], Batch Ensemble (BE) [67] and SWAG [49] as well as using DDU with a feature extractor trained on SNGP. For comparison, we also provide performance for vanilla SNGP, deep ensemble and DDU.

Dataset	Metric	Softmax & Energy	DUQ	SNGP	DDU	5-Ensemble
C10	ECE%	0.85 ± 0.02	1.55 ± 0.08	1.8 ± 0.1	0.85 ± 0.04	0.76 ± 0.03
	TACE%	0.63 ± 0.01	0.84 ± 0.03	0.9 ± 0.04	0.61 ± 0.01	0.48 ± 0.01
	NLL	0.18 ± 0.06	0.23 ± 0.07	0.27 ± 0.08	0.16 ± 0.06	0.11 ± 0.02
C100	ECE%	4.62 ± 0.06	-	4.33 ± 0.01	4.1 ± 0.08	3.32 ± 0.09
	TACE%	1.31 ± 0.02	-	1.23 ± 0.04	1.06 ± 0.03	0.58 ± 0.03
	NLL	1.17 ± 0.13	-	0.92 ± 0.16	0.86 ± 0.14	0.73 ± 0.09

Table 14. Calibration error scores ECE%, TACE% and NLL for WRN-28-10.

## F.6. Two Moons

In this section, we evaluate DDU’s performance on a well-known toy setup: the Two Moons dataset. We use scikit-learn’s *datasets* package to generate 2000 samples with a noise rate of 0.1. We use a 4-layer fully connected architecture, ResFFN-4-128 with 128 neurons in each layer and a residual connection, following [45]. As an ablation, we also train using a 4-layer fully connected architecture with 128 neurons in each layer, but *without the residual connection*. We name this architecture FC-Net. The input is 2-dimensional and is projected into the 128 dimensional space using a fully connected layer. Using the ResFFN-4-128 architecture we train 3 baselines:

1. **Softmax:** We train a single softmax model and use the softmax entropy as the uncertainty metric.
2. **3-ensemble:** We train an ensemble of 3 softmax models and use the predictive entropy of the ensemble as the measure of uncertainty.
3. **DDU:** We train a single softmax model applying spectral normalization on the fully connected layers and using the feature density as the measure of model confidence.

Each model is trained using the Adam optimiser for 150 epochs. In Fig. 14, we show the uncertainty results for all the above 3 baselines. It is clear that both the softmax entropy as well as the predictive entropy of the ensemble is uncertain only along the decision boundary between the two classes whereas DDU is confident only on the data distribution and is not confident anywhere else. It is worth mentioning that even DUQ and SNGP perform well in this setup and deep ensembles have been known to underperform in the Two-Moons setup primarily due to the simplicity of the dataset causing all the ensemble components to generalise in the same way. Fi-

nally, also note that the feature space density of FC-Net without residual connections is not able to capture uncertainty well (see Fig. 14d), thereby reaffirming our claim that proper inductive biases are indeed a necessary component to ensure that feature space density captures uncertainty reliably. Thus, in addition to its excellent performance in active learning, CIFAR-10 vs SVHN/CIFAR-100/Tiny-ImageNet/CIFAR-10-C, CIFAR-100 vs SVHN/Tiny-ImageNet, and ImageNet vs ImageNet-O, we note that DDU captures uncertainty reliably even in a small 2D setup like Two Moons.

## F.7. 5-Ensemble Visualisation

In Fig. 15, we provide a visualisation of a 5-ensemble (with five deterministic softmax networks) to see how softmax entropy fails to capture epistemic uncertainty precisely because the mutual information (MI) of an ensemble does not (see §5). We train the networks on 1-dimensional data with binary labels 0 and 1. The data is shown in Fig. 15a. From Fig. 15a and Fig. 15b, we find that the softmax entropy is high in regions of ambiguity where the label can be both 0 and 1 (i.e.  $x$  between  $-4.5$  and  $-3.5$ , and between  $3.5$  and  $4.5$ ). This indicates that softmax entropy can capture aleatoric uncertainty. Furthermore, in the  $x$  interval  $(-2, 2)$ , we find that the deterministic softmax networks disagree in their predictions (see Fig. 15a) and have softmax entropies which can be high, low or anywhere in between (see Fig. 15b) following our claim in §5. In fact, this disagreement is the very reason why the MI of the ensemble is high in the interval  $(-2, 2)$ , thereby reliably capturing epistemic uncertainty. Finally, note that the predictive entropy (PE) of the ensemble is high both in the OoD interval  $(-2, 2)$  as well as at points of ambiguity (i.e. at  $-4$  and  $4$ ). This indicates that the PE of a Deep Ensemble captures both epistemic and aleatoric uncertainty well. From these visualisations, we draw the conclusion that the softmax entropy of a deterministic softmax model cannot capture epistemic uncertainty precisely because the MI of a Deep Ensemble can.

## F.8. Feature-Space Density & Epistemic Uncertainty vs Softmax Entropy & Aleatoric Uncertainty

To empirically verify the connection between feature-space density and epistemic uncertainty on the one hand and the connection between softmax entropy and aleatoric uncertainty on the other hand, we train ResNet-18+SN models on increasingly large subsets of DirtyMNIST and CIFAR-10 and evaluate the epistemic and aleatoric uncertainty on the corresponding test sets using the feature-space density and softmax entropy, respectively. Moreover, we also train a 5-ensemble on the same subsets of data and use the ensemble’s mutual information as a baseline measure of epistemic uncertainty.

In Fig. 2, Fig. 16 and Tab. 15, we see that with larger

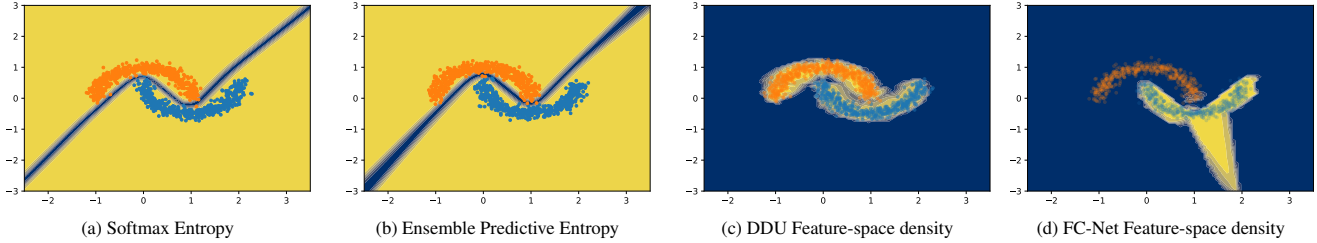


Figure 14. *Uncertainty on Two Moons dataset.* Blue indicates high uncertainty and yellow indicates low uncertainty. Both the softmax entropy of a single model as well as the predictive entropy of a deep ensemble are uncertain only along the decision boundary whereas the feature-space density of DDU is uncertain everywhere except on the data distribution (the ideal behaviour). However, the feature density of a normal fully connected network (FC-Net) without any inductive biases can't capture uncertainty properly.

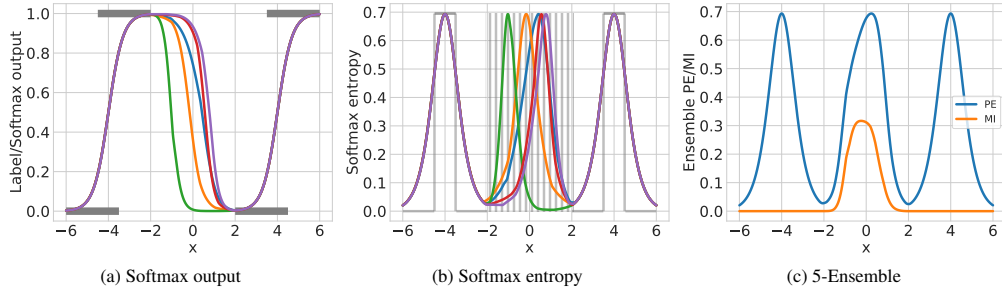


Figure 15. *Softmax outputs & entropies for 5 softmax models along with the predictive entropy (PE) and mutual information (MI) for the resulting 5-Ensemble.* (a) and (b) show that the softmax entropy is only reliably high for ambiguous iD points ( $\pm 3.5$ – $4.5$ ), whereas it can be low or high for OoD points ( $-2$ – $2$ ). The different colors are the different ensemble components. Similarly, (c) shows that the MI of the ensemble is only high for OoD, whereas the PE is high for both OoD and for regions of ambiguity. See §F.7.

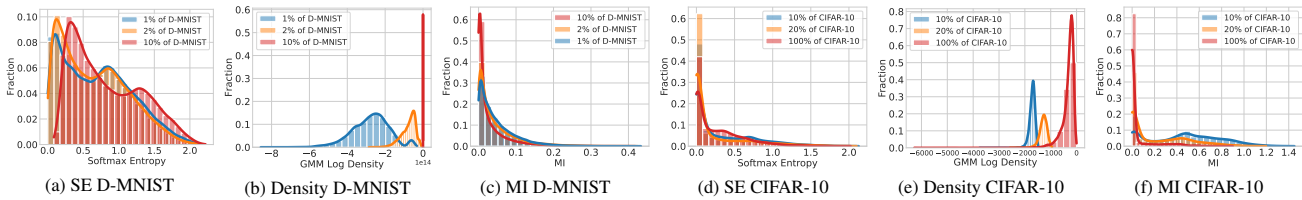


Figure 16. *Comparison of epistemic and aleatoric uncertainty captured by (ResNet-18+SN) on increasingly large subsets of Dirty-MNIST and CIFAR-10.* Clearly, feature density captures epistemic uncertainty which reduces when the model is trained on increasingly large subsets of training data, whereas softmax entropy (SE) does not. For comparison, we also plot a deep-ensemble's epistemic uncertainty, through mutual information (MI) for the same settings. For more details, see Tab. 15.

training sets, the average feature-space density increases which is consistent with the epistemic uncertainty decreasing as more data is available as reducible uncertainty. This is also evident from the consistent strong positive correlation between the negative log density and mutual information of the ensemble.

On the other hand, the softmax entropy stays roughly the same which is consistent with aleatoric uncertainty as irreducible uncertainty, which is independent of the training data. Importantly, all of this is also consistent with the experiments comparing epistemic and aleatoric uncertainty on increasing training set sizes in Table 3 of [33].

## F.9. Objective Mismatch Ablation with Wide-ResNet-28-10 on CIFAR-10

We further validate Proposition 5.3 by running an ablation on Wide-ResNet-28-10 on CIFAR-10. Table 16 shows that the feature-space density estimator indeed performs worse than the softmax layer for aleatoric uncertainty (accuracy and ECE).

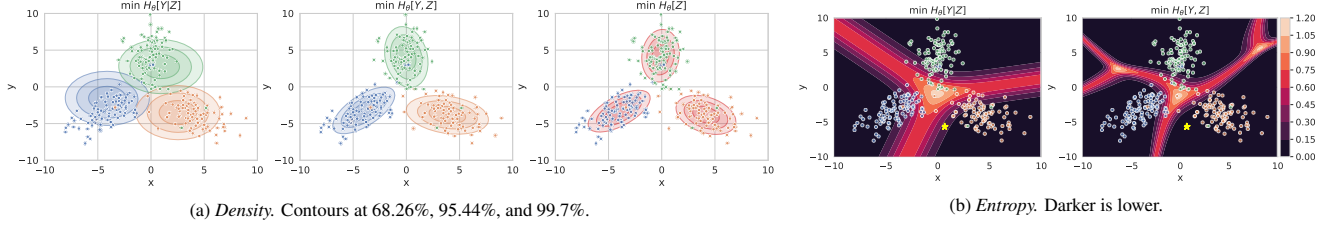


Figure 17. 3-component GMM fitted to a synthetic dataset with 3 different classes (differently colored) with 4% label noise using different objectives. (a): The optimas for conditional log-likelihood  $\mathbb{H}_\theta[Y | Z]$ , joint log-likelihood  $\mathbb{H}_\theta[Y, Z]$ , and marginalised log-likelihood  $\mathbb{H}_\theta[Z]$  all differ. Hence, the best calibrated model ( $\mathbb{H}_\theta[Y | Z]$ ) will not provide the best density estimate ( $\mathbb{H}_\theta[Z]$ ), and vice-versa. (b): A mixture model that optimizes  $\mathbb{H}_\theta[Y, Z]$  (GDA) does not have calibrated decision boundaries for aleatoric uncertainty: the ambiguous sample (due to label noise) marked by the yellow star has no aleatoric uncertainty under the GDA model. See §G.2.3 for details.

Table 15. Average softmax entropy (SE) and feature-space density of the test set for models trained on different amounts of the training set (Dirty-MNIST and CIFAR-10) behave consistently with aleatoric and epistemic uncertainty. Aleatoric uncertainty for individual samples does not change much as more data is added to the training set while epistemic uncertainty decreases as more data is added. This is also consistent with Table 3 in [33]. Finally, we observe a consistent strong positive correlation between the negative log feature space density and the mutual information (MI) of a deep ensemble trained on the same subsets of data for both Dirty-MNIST and CIFAR-10. However, the correlation between softmax entropy and MI is not consistent.

Training Set	Avg Test SE (±)	Avg Test Log GMM Density (†)	Avg Test MI	Correlation(SE    MI)	Correlation(-Log GMM Density    MI)
1% of D-MNIST	0.7407	-2.7268e + 14	0.0476		
2% of D-MNIST	0.6580	-7.8633e + 13	0.0447	-0.79897	0.8132
10% of D-MNIST	0.8295	-1279.1753	0.0286		
10% of CIFAR-10	0.3189	-1715.3516	0.4573		
20% of CIFAR-10	0.2305	-1290.1726	0.2247	0.5663	0.9556
100% of CIFAR-10	0.2747	-324.8040	0.0479		

Table 16. Objective Mismatch Ablation with WideResNet-28-10 models with and without spectral normalisation on CIFAR-10. While GMMs perform much better than Softmax Entropy for feature-space density/epistemic uncertainty estimation, they underperform for aleatoric uncertainty estimation: both accuracy and in particular ECE are much worse than a regular softmax layer. Averaged over 25 runs.

Model	Prediction Source	Accuracy in % (†)	ECE (‡)
WideResNet-28-10	Softmax	95.98 ± 0.02	2.29 ± 0.02
	GMM	95.86 ± 0.02	4.13 ± 0.02
WideResNet-28-10+SN	Softmax	95.97 ± 0.03	2.23 ± 0.03
	GMM	95.88 ± 0.02	4.12 ± 0.02

## G. Theoretical Results

### G.1. Softmax entropy “cannot” capture epistemic uncertainty because Deep Ensembles “can”

#### G.1.1 Qualitative Statement

We start with a proof of Proposition 5.2, which quantitatively examines the qualitative statements that given the same predictive entropy, higher epistemic uncertainty for one point

than another will cause some ensemble members to have lower softmax entropy.

**Proposition 5.2.** *Let  $x_1$  and  $x_2$  be points such that  $x_1$  has higher epistemic uncertainty than  $x_2$  under the ensemble:  $\mathbb{I}[Y_1; \omega | x_1, \mathcal{D}] > \mathbb{I}[Y_2; \omega | x_2, \mathcal{D}] + \delta, \delta \geq 0$ . Further assume both have similar predictive entropy  $|\mathbb{H}[Y_1 | x_1, \mathcal{D}] - \mathbb{H}[Y_2 | x_2, \mathcal{D}]| \leq \epsilon, \epsilon \geq 0$ . Then, there exist sets of ensemble members  $\Omega$  with  $p(\Omega | \mathcal{D}) > 0$ , such that for all softmax models  $\omega \in \Omega$  the softmax entropy of  $x_1$  is lower than the softmax entropy of  $x_2$ :  $\mathbb{H}[Y_1 | x_1, \omega] < \mathbb{H}[Y_2 | x_2, \omega] - (\delta - \epsilon)$ .*

*Proof.* From Eq. (1), we obtain

$$|\mathbb{I}[Y_1; \omega | x_1, \mathcal{D}] + \mathbb{E}_{p(\omega|\mathcal{D})} [\mathbb{H}[Y_1 | x_1, \omega]] - \mathbb{I}[Y_2; \omega | x_2, \mathcal{D}] - \mathbb{E}_{p(\omega|\mathcal{D})} [\mathbb{H}[Y_2 | x_2, \omega]]| \leq \epsilon. \quad (2)$$

and hence we have

$$\mathbb{E}_{p(\omega|\mathcal{D})} [\mathbb{H}[Y_1 | x_1, \omega]] - \mathbb{E}_{p(\omega|\mathcal{D})} [\mathbb{H}[Y_2 | x_2, \omega]] + \underbrace{(\mathbb{I}[Y_1; \omega | x_1, \mathcal{D}] - \mathbb{I}[Y_2; \omega | x_2, \mathcal{D}])}_{> \delta} \leq \epsilon. \quad (3)$$

We rearrange the terms:

$$\mathbb{E}_{p(\omega|\mathcal{D})} [\mathbb{H}[Y_1 | x_1, \omega]] < \mathbb{E}_{p(\omega|\mathcal{D})} [\mathbb{H}[Y_2 | x_2, \omega]] - (\delta - \epsilon). \quad (4)$$

Now, the statement follows by contraposition: if  $\mathbb{H}[Y_1 | x_1, \omega] \geq \mathbb{E}_{p(\omega|\mathcal{D})} [\mathbb{H}[Y_2 | x_2, \omega]] - (\delta - \epsilon)$  for all  $\omega$ , the monotonicity of the expectation would yield  $\mathbb{E}_{p(\omega|\mathcal{D})} [\mathbb{H}[Y_1 | x_1, \omega]] \geq \mathbb{E}_{p(\omega|\mathcal{D})} [\mathbb{H}[Y_2 | x_2, \omega]] - (\delta - \epsilon)$ . Thus, there is a non-null-set  $\Omega'$  with  $p(\Omega') > 0$ , such that

$$\mathbb{H}[Y_1 | x_1, \omega] < \mathbb{H}[Y_2 | x_2, \omega] - (\delta - \epsilon), \quad (5)$$

for all  $\omega \in \Omega'$ .  $\square$

While this statement provides us with an intuition for why ensemble members and thus deterministic models cannot provide epistemic uncertainty reliably through their softmax entropies, we can examine this further by establishing some upper bounds.

### G.1.2 Infinite Deep Ensemble

There are two interpretations of the ensemble parameter distribution  $p(\omega | \mathcal{D})$ : we can view it as an empirical distribution given a specific ensemble with members  $\omega_{i \in \{1, \dots, K\}}$ , or we can view it as a distribution over all possible trained models, given: random weight initializations, the dataset, stochasticity in the minibatches and the optimization process. In that case, any Deep Ensemble with  $K$  members can be seen as finite Monte-Carlo sample of this posterior distribution. The predictions of an ensemble then are an unbiased estimate of the predictive distribution  $\mathbb{E}_{p(\omega | \mathcal{D})} [p(y|x, \omega)]$ , and similarly the expected information gain computed using the members of the Deep Ensemble is just a (biased) estimator of  $\mathbb{I}[Y; \omega | x, \mathcal{D}]$ .

### G.1.3 Analysis of Softmax Entropy of a Single Deterministic Model on OoD Data using Properties of Deep Ensembles

Based on the interpretation of Deep Ensembles as a distribution over model parameters, we can walk backwards and, given *some value* for the predictive distribution and epistemic uncertainty of a Deep Ensemble, estimate what the softmax entropies from each ensemble component must have been. I.e. if we observe Deep Ensembles to have high epistemic uncertainty on OoD data, we can deduce from that what the softmax entropy of deterministic neural nets (the ensemble components) must look like. More specifically, given a predictive distribution  $p(y | x)$  and epistemic uncertainty, that is expected information gain  $\mathbb{I}[Y; \omega | x]$ , of the infinite Deep Ensemble, we estimate the expected softmax entropy from a single deterministic model, considered as a sample  $\omega \sim p(\omega | \mathcal{D})$  and model the variance. Empirically, we find the real variance to be higher by a large amount for OoD samples, showing that softmax entropies do not capture epistemic uncertainty well for samples with high epistemic uncertainty.

We will need to make several strong assumptions that limit the generality of our estimation, but we can show that our analysis models the resulting softmax entropy distributions appropriately. This will show that deterministic softmax models can have widely different entropies and confidence values.

Given the predictive distribution  $p(y | x)$  and epistemic uncertainty  $\mathbb{I}[Y; \omega | x]$ , we can approximate the distribution over softmax probability vectors  $p(y|x, \omega)$  for different  $\omega$  using its maximum-entropy estimate: a Dirichlet distribution  $(Y_1, \dots, Y_K) \sim \text{Dir}(\alpha)$  with non-negative concentration parameters  $\alpha = (\alpha_1, \dots, \alpha_K)$  and  $\alpha_0 := \sum \alpha_i$ . Note that the Dirichlet distribution is used *only as an analysis tool*, and at no point do we need to actually fit Dirichlet distributions to our data.

### Preliminaries

Before we can establish our main result, we need to look more closely at Dirichlet-Multinomial distributions. Given a Dirichlet distribution  $\text{Dir}(\alpha)$  and a random variable  $\mathbf{p} \sim \text{Dir}(\alpha)$ , we want to quantify the expected entropy  $\mathbb{E}_{\mathbf{p} \sim \text{Dir}(\alpha)} \mathbb{H}_{Y \sim \text{Cat}(\mathbf{p})}[Y]$  and its variance  $\text{Var}_{\mathbf{p} \sim \text{Dir}(\alpha)} \mathbb{H}_{Y \sim \text{Cat}(\mathbf{p})}[Y]$ . For this, we need to develop more theory. In the following,  $\Gamma$  denotes the Gamma function,  $\psi$  denotes the Digamma function,  $\psi'$  denotes the Trigamma function.

**Lemma G.1.** *Given a Dirichlet distribution and random variable  $\mathbf{p} \sim \text{Dir}(\alpha)$ , the following hold:*

1. *The expectation  $\mathbb{E}[\log \mathbf{p}_i]$  is given by:*

$$\mathbb{E}[\log \mathbf{p}_i] = \psi(\alpha_i) - \psi(\alpha_0). \quad (6)$$

2. *The covariance  $\text{Cov}[\log \mathbf{p}_i, \log \mathbf{p}_j]$  is given by*

$$\text{Cov}[\log \mathbf{p}_i, \log \mathbf{p}_j] = \psi'(\alpha_i) \delta_{ij} - \psi'(\alpha_0). \quad (7)$$

3. *The expectation  $\mathbb{E}[\mathbf{p}_i^n \mathbf{p}_j^m \log \mathbf{p}_i]$  is given by:*

$$\begin{aligned} \mathbb{E}[\mathbf{p}_i^n \mathbf{p}_j^m \log \mathbf{p}_i] \\ = \frac{\alpha_i^n \alpha_j^m}{\alpha_0^{n+m}} (\psi(\alpha_i + n) - \psi(\alpha_0 + n + m)), \end{aligned} \quad (8)$$

where  $i \neq j$ , and  $n^{\bar{k}} = n(n+1) \dots (n+k-1)$  denotes the rising factorial.

*Proof.* 1. The Dirichlet distribution is members of the exponential family. Therefore the moments of the sufficient statistics are given by the derivatives of the partition function with respect to the natural parameters. The natural parameters of the Dirichlet distribution are just its concentration parameters  $\alpha_i$ . The partition function is

$$A(\alpha) = \sum_{i=1}^k \log \Gamma(\alpha_i) - \log \Gamma(\alpha_0), \quad (9)$$

the sufficient statistics is  $T(x) = \log x$ , and the expectation  $\mathbb{E}[T]$  is given by

$$\mathbb{E}[T_i] = \frac{\partial A(\alpha)}{\partial \alpha_i} \quad (10)$$

as the Dirichlet distribution is a member of the exponential family. Substituting the definitions and evaluating the partial derivative yields

$$\mathbb{E}[\log \mathbf{p}_i] = \frac{\partial}{\partial \alpha_i} \left[ \sum_{i=1}^k \log \Gamma(\alpha_i) - \log \Gamma\left(\sum_{i=1}^k \alpha_i\right) \right] \quad (11)$$

$$= \psi(\alpha_i) - \psi(\alpha_0) \frac{\partial}{\partial \alpha_i} \alpha_0, \quad (12)$$

where we have used that the Digamma function  $\psi$  is the log derivative of the Gamma function  $\psi(x) = \frac{d}{dx} \ln \Gamma(x)$ . This proves (6) as  $\frac{\partial}{\partial \alpha_i} \alpha_0 = 1$ .

2. Similarly, the covariance is obtained using a second-order partial derivative:

$$\text{Cov}[T_i, T_j] = \frac{\partial^2 A(\alpha)}{\partial \alpha_i \partial \alpha_j}. \quad (13)$$

Again, substituting yields

$$\text{Cov}[\log \mathbf{p}_i, \log \mathbf{p}_j] = \frac{\partial}{\partial \alpha_j} [\psi(\alpha_i) - \psi(\alpha_0)] \quad (14)$$

$$= \psi'(\alpha_i) \delta_{ij} - \psi'(\alpha_0). \quad (15)$$

3. We will make use of a simple reparameterization to prove the statement using Eq. (6). Expanding the expectation and substituting the density  $\text{Dir}(\mathbf{p}; \alpha)$ , we obtain

$$\mathbb{E}[\mathbf{p}_i^n \mathbf{p}_j^m \log \mathbf{p}_i] = \int \text{Dir}(\mathbf{p}; \alpha) \mathbf{p}_i^n \mathbf{p}_j^m \log \mathbf{p}_i d\mathbf{p} \quad (16)$$

$$= \int \frac{\Gamma(\alpha_0)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{k=1}^K \mathbf{p}_k^{\alpha_k-1} \mathbf{p}_i^n \mathbf{p}_j^m \log \mathbf{p}_i d\mathbf{p} \quad (17)$$

$$= \frac{\Gamma(\alpha_i + n) \Gamma(\alpha_j + m) \Gamma(\alpha_0 + n + m)}{\Gamma(\alpha_i) \Gamma(\alpha_j) \Gamma(\alpha_0)} \quad (18)$$

$$\int \text{Dir}(\hat{\mathbf{p}}; \hat{\alpha}) \hat{\mathbf{p}}_i^n \hat{\mathbf{p}}_j^m \log \hat{\mathbf{p}}_i d\hat{\mathbf{p}} \\ = \frac{\alpha_i^n \alpha_j^m}{\alpha_0^{n+m}} \mathbb{E}[\log \hat{\mathbf{p}}_i], \quad (19)$$

where  $\hat{\mathbf{p}} \sim \text{Dir}(\hat{\alpha})$  with  $\hat{\alpha} = (\alpha_0, \dots, \alpha_i + n, \dots, \alpha_j + m, \dots, \alpha_K)$  and we made use of the fact that  $\frac{\Gamma(z+n)}{\Gamma(z)} = z^n$ . Finally, we can apply Eq. (6) on  $\hat{\mathbf{p}} \sim \text{Dir}(\hat{\alpha})$  to show

$$= \frac{\alpha_i^n \alpha_j^m}{\alpha_0^{n+m}} (\psi(\alpha_i + n) - \psi(\alpha_0 + n + m)). \quad (20)$$

□

With this, we can already quantify the expected entropy  $\mathbb{E}_{\mathbf{p} \sim \text{Dir}(\alpha)} \mathbb{H}_{Y \sim \text{Cat}(\mathbf{p})}[Y]$ :

**Lemma G.2.** *Given a Dirichlet distribution and a random variable  $\mathbf{p} \sim \text{Dir}(\alpha)$ , the expected entropy  $\mathbb{E}_{\mathbf{p} \sim \text{Dir}(\alpha)} \mathbb{H}_{Y \sim \text{Cat}(\mathbf{p})}[Y]$  of the categorical distribution  $Y \sim \text{Cat}(\mathbf{p})$  is given by*

$$\mathbb{E}_{\mathbf{p}(\mathbf{p}|\alpha)} \mathbb{H}[Y | \mathbf{p}] = \psi(\alpha_0 + 1) - \sum_{y=1}^K \frac{\alpha_y}{\alpha_0} \psi(\alpha_y + 1). \quad (21)$$

*Proof.* Applying the sum rule of expectations and Eq. (8) from Lemma G.1, we can write

$$\mathbb{E} \mathbb{H}[Y | \mathbf{p}] = \mathbb{E} \left[ - \sum_{i=1}^K \mathbf{p}_i \log \mathbf{p}_i \right] = - \sum_i \mathbb{E}[\mathbf{p}_i \log \mathbf{p}_i] \quad (22)$$

$$= - \sum_i \frac{\alpha_i}{\alpha_0} (\psi(\alpha_i + 1) - \psi(\alpha_0 + 1)). \quad (23)$$

The result follows after rearranging and making use of  $\sum_i \frac{\alpha_i}{\alpha_0} = 1$ . □

With these statements, we can answer a slightly more complex problem:

**Lemma G.3.** *Given a Dirichlet distribution and a random variable  $\mathbf{p} \sim \text{Dir}(\alpha)$ , the covariance  $\text{Cov}[\mathbf{p}_i^n \log \mathbf{p}_i, \mathbf{p}_j^m \log \mathbf{p}_j]$  is given by*

$$\text{Cov}[\mathbf{p}_i^n \log \mathbf{p}_i, \mathbf{p}_j^m \log \mathbf{p}_j] \quad (24)$$

$$= \frac{\alpha_i^n \alpha_j^m}{\alpha_0^{n+m}} ((\psi(\alpha_i + n) - \psi(\alpha_0 + n + m)) \\ (\psi(\alpha_j + m) - \psi(\alpha_0 + n + m)) \\ - \psi'(\alpha_0 + n + m)) \\ + \frac{\alpha_i^n \alpha_j^m}{\alpha_0^n \alpha_0^m} (\psi(\alpha_i + n) - \psi(\alpha_0 + n)) \\ (\psi(\alpha_j + m) - \psi(\alpha_0 + n)), \quad (25)$$

for  $i \neq j$ , where  $\psi$  is the Digamma function and  $\psi'$  is the Trigamma function. Similarly, the covariance  $\text{Cov}[\mathbf{p}_i^n \log \mathbf{p}_i, \mathbf{p}_i^m \log \mathbf{p}_i]$  is given by

$$\text{Cov}[\mathbf{p}_i^n \log \mathbf{p}_i, \mathbf{p}_i^m \log \mathbf{p}_i] \quad (26)$$

$$= \frac{\alpha_i^{n+m}}{\alpha_0^{n+m}} ((\psi(\alpha_i + n + m) - \psi(\alpha_0 + n + m))^2 \\ + \psi'(\alpha_i + n + m) - \psi'(\alpha_0 + n + m)) \\ + \frac{\alpha_i^n \alpha_i^m}{\alpha_0^n \alpha_0^m} (\psi(\alpha_i + n) - \psi(\alpha_0 + n)) \\ (\psi(\alpha_i + m) - \psi(\alpha_0 + n)). \quad (27)$$

Regrettably, the equations are getting large. By abuse of notation, we introduce a convenient shorthand before proving the lemma.

**Definition G.4.** We will denote by

$$\overline{\mathbb{E}[\log \hat{\mathbf{p}}_i^{n,m}]} = \psi(\alpha_i + n) - \psi(\alpha_0 + n + m), \quad (28)$$

and use  $\overline{\mathbb{E}[\log \hat{\mathbf{p}}_i^n]}$  for  $\overline{\mathbb{E}[\log \hat{\mathbf{p}}_i^{n,0}]}$ . Likewise,

$$\overline{\text{Cov}[\log \hat{\mathbf{p}}_i^{n,m}, \log \hat{\mathbf{p}}_j^{n,m}]} = \psi'(\alpha_i + n) \delta_{ij} - \psi'(\alpha_0 + n + m). \quad (29)$$

This notation agrees with the proof of Eq. (6) and (7) in Lemma G.1. With this, we can significantly simplify the previous statements:

**Corollary G.5.** *Given a Dirichlet distribution and random variable  $\mathbf{p} \sim \text{Dir}(\alpha)$ ,*

$$\mathbb{E} [\mathbf{p}_i^n \mathbf{p}_j^m \log \mathbf{p}_i] = \frac{\alpha_i^{\bar{n}} \alpha_j^{\bar{m}}}{\alpha_0^{\bar{n}+\bar{m}}} \mathbb{E} [\log \hat{\mathbf{p}}_i^{n,m}], \quad (30)$$

$$\begin{aligned} \text{Cov}[\mathbf{p}_i^n \log \mathbf{p}_i, \mathbf{p}_j^m \log \mathbf{p}_j] & \quad (31) \\ &= \frac{\alpha_i^{\bar{n}} \alpha_j^{\bar{m}}}{\alpha_0^{\bar{n}+\bar{m}}} \left( \overline{\mathbb{E} [\log \hat{\mathbf{p}}_i^{n,m}] \mathbb{E} [\log \hat{\mathbf{p}}_j^{m,n}]} \right. \\ & \quad \left. \overline{\text{Cov}[\log \hat{\mathbf{p}}_i^{n,m}, \log \hat{\mathbf{p}}_j^{m,n}]} \right) \\ & \quad + \frac{\alpha_i^{\bar{n}} \alpha_j^{\bar{m}}}{\alpha_0^{\bar{n}} \alpha_0^{\bar{m}}} \mathbb{E} [\log \hat{\mathbf{p}}_i^n] \mathbb{E} [\log \hat{\mathbf{p}}_j^m] \quad \text{for } i \neq j, \text{ and} \end{aligned} \quad (32)$$

$$\begin{aligned} \text{Cov}[\mathbf{p}_i^n \log \mathbf{p}_i, \mathbf{p}_i^m \log \mathbf{p}_i] & \quad (33) \\ &= \frac{\alpha_i^{\bar{n}+\bar{m}}}{\alpha_0^{\bar{n}+\bar{m}}} \left( \overline{\mathbb{E} [\log \hat{\mathbf{p}}_i^{n+m}]}^2 \right. \\ & \quad \left. + \overline{\text{Cov}[\log \hat{\mathbf{p}}_i^{n+m}, \log \hat{\mathbf{p}}_i^{n+m}]} \right) \\ & \quad + \frac{\alpha_i^{\bar{n}} \alpha_i^{\bar{m}}}{\alpha_0^{\bar{n}} \alpha_0^{\bar{m}}} \mathbb{E} [\log \hat{\mathbf{p}}_i^n] \mathbb{E} [\log \hat{\mathbf{p}}_i^m]. \end{aligned} \quad (34)$$

*Proof of Lemma G.3.* This proof applies the well-know formula (**cov**)  $\text{Cov}[X, Y] = \mathbb{E}[XY] - \mathbb{E}[X] \mathbb{E}[Y]$  once forward and once backward (**rcov**)  $\mathbb{E}[XY] = \text{Cov}[X, Y] + \mathbb{E}[X] \mathbb{E}[Y]$  while applying Eq. (8) several times:

$$\text{Cov}[\mathbf{p}_i^n \log \mathbf{p}_i, \mathbf{p}_j^m \log \mathbf{p}_j] \quad (35)$$

$$\stackrel{\text{cov}}{=} \mathbb{E} [\mathbf{p}_i^n \log(\mathbf{p}_i) \mathbf{p}_j^m \log(\mathbf{p}_j)] \quad (36)$$

$$- \mathbb{E} [\mathbf{p}_i^n \log \mathbf{p}_i] \mathbb{E} [\mathbf{p}_j^m \log \mathbf{p}_j] \quad (37)$$

$$= \frac{\alpha_i^{\bar{n}} \alpha_j^{\bar{m}}}{\alpha_0^{\bar{n}+\bar{m}}} \mathbb{E} [\log(\hat{\mathbf{p}}_i^{i,j}) \log(\hat{\mathbf{p}}_j^{i,j})] \quad (37)$$

$$- \mathbb{E} [\log \hat{\mathbf{p}}_i^i] \mathbb{E} [\log \mathbf{p}_j^j] \quad (38)$$

$$\stackrel{\text{rcov}}{=} \frac{\alpha_i^{\bar{n}} \alpha_j^{\bar{m}}}{\alpha_0^{\bar{n}+\bar{m}}} \left( \overline{\text{Cov}[\log \hat{\mathbf{p}}_i^{i,j}, \log \hat{\mathbf{p}}_j^{i,j}]} \right. \\ \left. + \mathbb{E} [\log \hat{\mathbf{p}}_i^{i,j}] \mathbb{E} [\log \hat{\mathbf{p}}_j^{i,j}] \right) \quad (38)$$

$$- \frac{\alpha_i^{\bar{n}} \alpha_j^{\bar{m}}}{\alpha_0^{\bar{n}} \alpha_0^{\bar{m}}} \mathbb{E} [\log \hat{\mathbf{p}}_i^i] \mathbb{E} [\log \mathbf{p}_j^j],$$

where  $\mathbf{p}^{i,j} \sim \text{Dir}(\alpha^{i,j})$  with  $\alpha^{i,j} = (\dots, \alpha_i + n, \dots, \alpha_j + m, \dots)$ .  $\mathbf{p}^{i/j}$  and  $\alpha^{i/j}$  are defined analogously. Applying Eq. (7) and Eq. (6) from Lemma G.1 yields the statement. For  $i = j$ , the proof follows the same pattern.  $\square$

Now, we can prove the theorem that quantifies the variance of the entropy of  $Y$ :

**Theorem G.6.** *Given a Dirichlet distribution and a random variable  $\mathbf{p} \sim \text{Dir}(\alpha)$ , the variance of the entropy  $\text{Var}_{\mathbf{p} \sim \text{Dir}(\alpha)} \mathbb{H}_{Y \sim \text{Cat}(\mathbf{p})}[Y]$  of the categorical distribution  $Y \sim \text{Cat}(\mathbf{p})$  is given by*

$$\begin{aligned} \text{Var}[\mathbb{H}[Y | \mathbf{p}]] & \quad (39) \\ &= \sum_i \frac{\alpha_i^{\bar{2}}}{\alpha_0^{\bar{2}}} \left( \overline{\text{Cov}[\log \hat{\mathbf{p}}_i^2, \log \hat{\mathbf{p}}_i^2]} + \overline{\mathbb{E} [\log \hat{\mathbf{p}}_i^2]^2} \right) \\ & \quad + \sum_{i \neq j} \frac{\alpha_i \alpha_j}{\alpha_0^{\bar{2}}} \left( \overline{\text{Cov}[\log \hat{\mathbf{p}}_i^1, \log \hat{\mathbf{p}}_j^1]} \right. \\ & \quad \left. + \overline{\mathbb{E} [\log \hat{\mathbf{p}}_i^1] \mathbb{E} [\log \hat{\mathbf{p}}_j^1]} \right) \\ & \quad - \sum_{i,j} \frac{\alpha_i \alpha_j}{\alpha_0^{\bar{2}}} \overline{\mathbb{E} [\log \hat{\mathbf{p}}_i^1] \mathbb{E} [\log \hat{\mathbf{p}}_j^1]}. \end{aligned} \quad (40)$$

*Proof.* We start by applying the well-known formula  $\text{Var}[\sum_i X_i] = \sum_{i,j} \text{Cov}[X_i, X_j]$  and then apply Lemma G.3 repeatedly.  $\square$

## Main Result

Given that we can view an ensemble member as a single deterministic model and vice versa, this provides an intuitive explanation for why single deterministic models report inconsistent and widely varying predictive entropies and confidence scores for OoD samples for which a Deep Ensemble would report high epistemic uncertainty (expected information gain) and high predictive entropy.

Assuming that  $p(y|x, \omega)$  only depends on  $\mathbf{p}(y|x)$  and  $\mathbb{I}[Y; \omega | x]$ , we model the distribution of  $p(y|x, \omega)$  (as a function of  $\omega$ ) using a Dirichlet distribution  $\text{Dir}(\alpha)$  which satisfies:

$$\mathbf{p}(y|x) = \frac{\alpha_y}{\alpha_0} \quad (41)$$

$$\mathbb{H}[Y | x] - \mathbb{I}[Y; \omega | x] = \psi(\alpha_0 + 1) \quad (42)$$

$$- \sum_{y=1}^K \mathbf{p}(y|x) \psi(\alpha_0 \mathbf{p}(y|x) + 1).. \quad (43)$$

Then, we can model the softmax distribution using a random variable  $\mathbf{p} \sim \text{Dir}(\alpha)$  as:

$$\mathbf{p}(y|x, \omega) \approx \text{Cat}(\mathbf{p}). \quad (44)$$

The variance  $\text{Var}[\mathbb{H}[Y | x, \omega]]$  of the softmax entropy for different samples  $x$  given  $\mathbf{p}(y|x)$  and  $\mathbb{I}[Y; \omega | x]$  is then approximated by  $\text{Var}[\mathbb{H}[Y | \mathbf{p}]]$ :

$$\text{Var}_{\omega}[\mathbb{H}[Y | x, \omega]] \approx \text{Var}_{\mathbf{p}}[\mathbb{H}[Y | \mathbf{p}]] \quad (45)$$

with the latter term given in eq. (40). We empirically find this to provide a lower bound on the true variance  $\text{Var}_\omega[\mathbb{H}[Y | x, \omega]]$ .

## Empirical Results

We empirically verify that softmax entropies vary considerably in Fig. 18. In Fig. 19, we verify that the predicted softmax entropy variance indeed lower-bounds the empirical softmax entropy variance. Moreover, Fig. 19c shows both **i**) the non-linear relationship between epistemic uncertainty and variance in the softmax entropies and **ii**) that Dirichlet distributions cannot capture it and can only provide a lower bound. Nonetheless, this simple approximation seems to be able to capture the empirical entropy distribution quite well as shown in Fig. 20.

## G.2. Objective Mismatch

In §5, we noted that the objectives that lead to optimal estimators for aleatoric and epistemic uncertainty via softmax entropy and feature-space density do not match, and DDU therefore uses the softmax layer as a discriminative classifier (implicit LDA) to estimate the predictive entropy, while it is using a GMM as generative classifier to estimate the feature-space density. Here we prove this.

### G.2.1 Preliminaries

Before we prove Proposition 5.3, we will introduce some additional notation following [36].

- Definition G.7.** 1.  $\hat{p}(y, z)$  is the data distribution of the  $\mathcal{D}$  in feature space with class labels  $y$  and feature representation  $z$ .  
 2.  $p_\theta(\cdot)$  is a probability distribution parameterized by  $\theta$ .  
 3. Entropies and conditional entropies are over the empirical data distribution  $\hat{p}(\cdot)$ :

$$\mathbb{H}[\cdot] = \mathbb{H}(\hat{p}(\cdot)) = \mathbb{E}_{\hat{p}(\cdot)}[-\log \hat{p}(\cdot)]. \quad (46)$$

4.  $\mathbb{H}[Y | z]$  is the entropy of  $\hat{p}(y | z)$  for a given  $z$ , whereas  $\mathbb{H}[Y | Z]$  is the conditional entropy:

$$\mathbb{H}[Y | Z] = \mathbb{E}_{\hat{p}(z)} \mathbb{H}[Y | z]. \quad (47)$$

5.  $\mathbb{H}(p(y, z) || q(y | z))$  is the cross-entropy of  $q(y | z)$  under  $p(y | z)$  in expectation over  $p(z)$ :

$$\begin{aligned} \mathbb{H}(p(y, z) || q(y | z)) &= \mathbb{E}_{p(z)} \mathbb{H}(p(y | z) || q(y | z)) \\ &= \mathbb{E}_{p(y, z)} [-\log q(y | z)]. \end{aligned}$$

6. Similarly,  $D_{\text{KL}}(p(y, z) || q(y | z))$  is the Kullback-Leibler divergence of  $q(y | z)$  under  $p(y | z)$  in expectation over  $p(z)$ :

$$\begin{aligned} D_{\text{KL}}(p(y, z) || q(y | z)) &= \mathbb{E}_{p(z)} D_{\text{KL}}(p(y | z) || q(y | z)) \\ &= \mathbb{H}(p(y, z) || q(y | z)) - \mathbb{H}[Y | Z] \end{aligned}$$

7. For cross-entropies of  $p_\theta(\cdot)$  under  $\hat{p}(z, y)$ , we use the convenient short-hand  $\mathbb{H}_\theta[\cdot] = \mathbb{H}(\hat{p}(z, y) || p_\theta(\cdot))$ .

Then we can observe the following connection between  $\mathbb{H}_\theta[\cdot]$  and  $\mathbb{H}[\cdot]$ :

**Lemma G.8.** *Cross-entropies upper-bound the respective entropy with equality when  $p_\theta(\cdot) = \hat{p}(\cdot)$ , which is important for variational arguments:*

1.  $\mathbb{H}_\theta[Y, Z] \geq \mathbb{H}[Y, Z]$ ,
2.  $\mathbb{H}_\theta[Z] \geq \mathbb{H}[Z]$ , and
3.  $\mathbb{H}_\theta[Y | Z] \geq \mathbb{H}[Y | Z]$ .

*Proof.* 1.  $\mathbb{H}_\theta[Y, Z] - \mathbb{H}[Y, Z] = D_{\text{KL}}(\hat{p}(y, z) || p_\theta(y, z)) \geq 0$ .

2. follows from Item 1.

3. We expand the expectations and note that inequality commutes with expectations:

$$\mathbb{H}_\theta[Y | Z] - \mathbb{H}[Y | Z] = \mathbb{E}_{\hat{p}(z)} [\mathbb{H}_\theta[Y | z] - \mathbb{H}[Y | z]] \geq 0,$$

because  $\mathbb{H}_\theta[Y | z] - \mathbb{H}[Y | z] \geq 0$  for all  $z$ . The equality conditions follows from the properties of the Kullback-Leibler divergence as well.  $\square$

We also have:

**Lemma G.9.**

$$\mathbb{H}_\theta[Y, Z] = \mathbb{H}_\theta[Y | Z] + \mathbb{H}_\theta[Z] \quad (48)$$

$$= \mathbb{H}_\theta[Z | Y] + \mathbb{H}_\theta[Y]. \quad (49)$$

*Proof.* We substitute the definitions and obtain:

$$\mathbb{H}_\theta[Y, Z] = \mathbb{E}_{p(y, z)} [-\log q(y, z)] \quad (50)$$

$$= \mathbb{E}_{p(y, z)} [-\log q(y | z)] + \mathbb{E}_{p(y, z)} [-\log q(z)] \quad (51)$$

$$= \mathbb{H}_\theta[Y | Z] + \mathbb{H}_\theta[Z]. \quad (52)$$

$\square$

The same holds for entropies:  $\mathbb{H}[Y, Z] = \mathbb{H}[Y | Z] + \mathbb{H}[Z] = \mathbb{H}[Y | Z] + \mathbb{H}[Y]$  [5].

### G.2.2 Proof

We can now prove the observation.

**Proposition 5.3.** *For an input  $x$ , let  $z = f_\theta(x)$  denote its feature representation in a feature extractor  $f_\theta$  with parameters  $\theta$ . Then the following hold:*

1. *A discriminative classifier  $p(y | z)$ , e.g. a softmax layer, is well-calibrated in its predictions when it maximises the conditional log-likelihood  $\log p(y | z)$ ;*
2. *A feature-space density estimator  $q(z)$  is optimal when it maximises the marginalised log-likelihood  $\log q(z)$ ;*



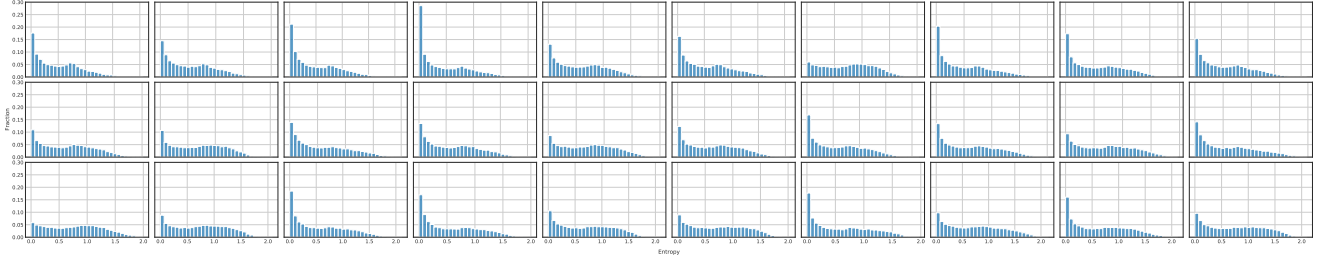


Figure 18. *Softmax entropy histograms of 30 Wide-ResNet-28-10+SN models trained on CIFAR-10, evaluated on SVHN (OoD). The softmax entropy distribution of the different models varies considerably.*

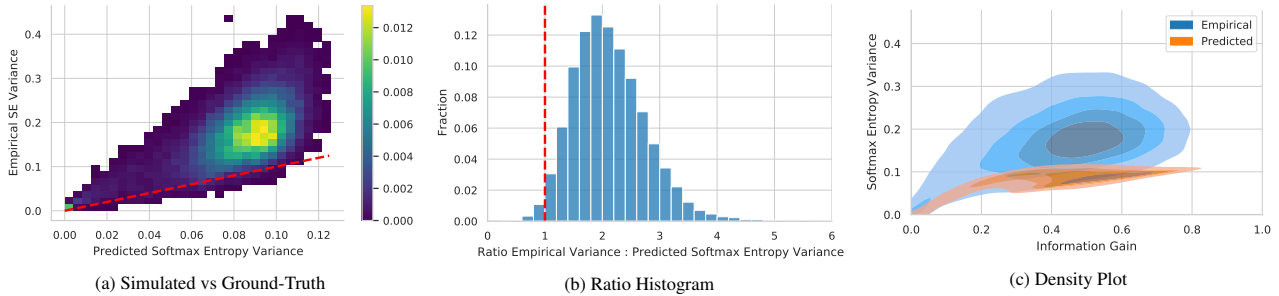


Figure 19. *The variance of softmax entropies can be lower-bounded by fitting Dirichlet distributions on the samples  $p(y | x, \omega)$ . (a) The empirical variance of softmax entropies is lower-bounded by  $\text{Var}[\mathbb{H}[Y | \mathbf{p}]]$ . The red dashed line depicts equality. (b) The ratio histogram shows that there are only few violations due to precision issues ( $< 2\%$ ). (c) The variance of the softmax entropy is not linearly correlated to the epistemic uncertainty. For both high and low epistemic uncertainty, the variance decreases.*

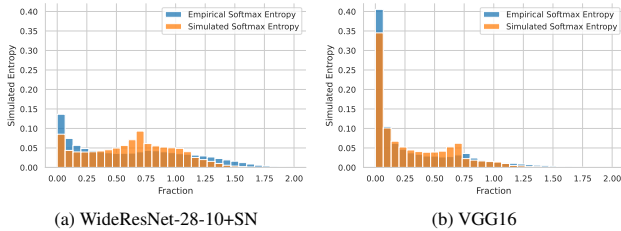


Figure 20. *Simulated vs empirical softmax entropy on WideResNet-28-10+SN and VGG16. Even though the Dirichlet variance approximation lower-bounds the empirical softmax entropy variance, sampling from the fitted Dirichlet distributions does approximate the empirical entropy distribution quite well.*

3. A mixture model  $q(y, z) = \sum_y q(z | y) q(y)$  might not maximise both objectives, conditional log-likelihood and marginalised log-likelihood, at the same time. In the specific instance that a GMM with one component per class does maximise both, the resulting model must be a GDA (but the opposite does not hold).

*Proof.* 1. The conditional log-likelihood is a strictly proper scoring rule [17]. The optimization objective can be rewritten as

$$\max_{\theta} \mathbb{E}_{\log p_{\theta}(y|z)} = \min_{\theta} \mathbb{H}_{\theta}[Y | Z] \geq \mathbb{H}[Y | Z]. \quad (53)$$

An optimal discriminative classifier  $p_{\theta}(y | z)$  would thus capture the true (empirical) distribution everywhere:  $p_{\theta}(y | z) = \hat{p}(y | z)$ . This means the negative conditional log-likelihood will be equal  $\mathbb{H}[Y | Z]$  and  $\mathbb{H}_{\theta}[Y | z] = \mathbb{H}[Y | z]$  for all  $z$ .

2. For density estimation  $q(z)$ , the maximum likelihood  $\mathbb{E}[\log q(z)]$  using the empirical data distribution is maximized. We can rewrite this as

$$\max_{\theta} \mathbb{E}_{\hat{p}(y,z)} \log p_{\theta}(z) = \min_{\theta} \mathbb{H}_{\theta}[Z] \geq \mathbb{H}[Z]. \quad (54)$$

We see that the negative marginalized likelihood of the density estimator upper-bounds the entropy of the feature representations  $\mathbb{H}[Z]$ . We have equality and  $p_{\theta}(z) = \hat{p}(z)$  in the optimum case.

3. Using  $\mathbb{H}_{\theta}[Y, Z] = \mathbb{H}_{\theta}[Y | Z] + \mathbb{H}_{\theta}[Z]$ , we can relate the objectives from Eq. (53) and (54) to each other. First, we characterize a shared optimum, and then we show that both objectives are generally not minimized at the same time. For both objectives to be minimized, we have  $\nabla \mathbb{H}_{\theta}[Y | Z] = 0$  and  $\nabla \mathbb{H}_{\theta}[Z] = 0$ , and we obtain

$$\nabla \mathbb{H}_{\theta}[Y, Z] = \nabla \mathbb{H}_{\theta}[Y | Z] + \nabla \mathbb{H}_{\theta}[Z] = 0. \quad (55)$$

From this we conclude that minimizing both objectives also minimizes  $\mathbb{H}_{\theta}[Y, Z]$ , and that generally the objectives trade-

off with each other at stationary points  $\theta$  of  $\mathbb{H}_\theta[Y, Z]$ :

$$\nabla \mathbb{H}_\theta[Y | Z] = -\nabla \mathbb{H}_\theta[Z] \quad \text{when } \nabla \mathbb{H}_\theta[Y, Z] = 0. \quad (56)$$

This tells us that to construct a case where the optima do not coincide, discriminative classification needs to be opposed better density estimation.

Specially, when we have a GMM with one component per class, minimizing  $\mathbb{H}_\theta[Y, Z]$  on an empirical data distribution is equivalent to Gaussian Discriminant Analysis, as is easy to check, and minimizing  $\mathbb{H}_\theta[Z]$  is equivalent to fitting a density estimator, following Eq. (54). The difference is that using a GMM as a density estimator does not constrain the component assignment, unlike in GDA.

Consequently, we see that *both objectives can be minimized at the same time exactly when the feature representations of different classes are perfectly separated*, such that a GMM fit as density estimator would assign each class's feature representations to a single component.

By the above, we can construct a simple case: if the samples of different classes are not separated in feature-space, optima for the objectives will not coincide, so for example if samples were drawn from the same Gaussian and labeled randomly. On the other hand, if we have classes whose features lie in well-separated clusters, GDA will minimize all objectives.  $\square$

Given that perfect separation is impossible with ambiguous data for a GMM, a shared optimum will be rare with noisy real-world data, but only then would GDA be optimal. In all other cases, GDA does not optimize both objectives, and neither can any other GMM with one component per class. Moreover, Eq. (56) shows that a GMM fit using EM is a better density estimator than GDA, and a softmax layer is a better classifier, as optimizing the softmax objective  $\mathbb{H}_\theta[Y | Z]$  or density objective  $\mathbb{H}_\theta[Z]$  using gradient descent will move away from the GDA optimum.

As can easily be verified, a trivial optimal minimizer  $q^*(y, z)$  for  $\mathbb{H}_\theta[Y, Z]$  given an empirical data distribution  $\hat{p}(y, z)$  is an adapted Parzen estimator:

$$q^*(y, z) = \sum_y \hat{p}(y) \mathbb{E}_{z \sim \hat{p}(z|y)} \mathcal{N}(z; \hat{z}, \sigma^2 \mathbf{I}), \quad (57)$$

for small enough  $\sigma$ . This shows that above proposition is not general.

### G.2.3 Intuitions & Validation with a Toy Example

Figure 17 visualises this on a synthetic 2D dataset with three classes and 4% label noise, which causes the optima to diverge as described in the proof. Label noise is a common issue in real-world datasets. Non-separability even more so. To explain Proposition 5.3 in an intuitive way, we focus on

on a simple 2D toy case and fit a GMM using the different objectives. We sample "latents"  $z$  from 3 Gaussians (each representing a different class  $y$ ) with 4% label noise. Following the construction in the proof, this will lead the objectives to have different optima.

We now discuss the different objectives in Fig. 17 and the resulting scores in more detail:

**min  $\mathbb{H}_\theta[Y | Z]$ .** A softmax linear layer is equivalent to an LDA (Linear Discriminant Analysis) with conditional likelihood as detailed in [54], for example. We optimize an LDA with the usual objective " $\min -1/N \sum \log p(y | z)$ ", i.e. the cross-entropy of  $p(y | z)$  or (average) negative log-likelihood (NLL). Following Definition G.7, we use the short-hand " $\min \mathbb{H}_\theta[Y | Z]$ " for this cross-entropy.

Because we optimize only  $p(y|z)$ ,  $p(z)$  does not affect the objective and is thus not optimized. Indeed, the components do not actually cover the latents well, as can be seen in the first density plot of Fig. 17a. However, it does provide the lowest NLL.

**min  $\mathbb{H}_\theta[Y, Z]$ .** We optimize a GDA for the combined objective " $\min -1/N \sum \log q(y, z)$ ", i.e. the cross-entropy of  $q(y, z)$ . We use the short-hand " $\min \mathbb{H}_\theta[Y | Z]$ " for this.

**min  $\mathbb{H}_\theta[Z]$ .** We optimize a GMM for the objective " $\min -1/N \sum \log q(z)$ ", i.e. the cross-entropy of  $q(z)$ . We use the short-hand " $\min \mathbb{H}_\theta[Z]$ " for this.

We do not provide scores for  $\mathbb{H}_\theta[Y | Z]$  and  $\mathbb{H}_\theta[Y, Z]$  for the third objective  $\min \mathbb{H}_\theta[Z]$  in Tab. 17 as it does not depend on  $Y$ , and hence the different components do not actually model the different classes necessarily. Hence, we also use a single color to visualize the components for this objective in Fig. 17a.

In Tab. 17 and Fig. 17a, we see that each solution minimizes its own objective best. The GMM provides the best density model (best fit according to the entropy), while the LDA (like a softmax linear layer) provides the best NLL for the labels. The GDA provides a density model that is almost as good.

**Entropy.** Looking at the entropy plots in Fig. 17b, we first notice that the LDA solution optimized for  $\min \mathbb{H}_\theta[Y | Z]$  has a wide decision boundary. This is due to the overlap

Table 17. *Realized objective scores (columns) for different optimization objectives (rows) for the synthetic 2D toy example depicted in Fig. 17. Smaller is better. We see that each objectives minimizes its own score while being suboptimal in regards to the other two objectives (when it is possible to compute the scores). This empirically further validates Proposition 5.3.*

Objective	$\mathbb{H}_\theta[Y   Z]$ ( $\downarrow$ )	$\mathbb{H}_\theta[Y, Z]$ ( $\downarrow$ )	$\mathbb{H}_\theta[Z]$ ( $\downarrow$ )
min $\mathbb{H}_\theta[Y   Z]$	0.1794	5.4924	5.2995
min $\mathbb{H}_\theta[Y, Z]$	0.2165	4.9744	4.7580
min $\mathbb{H}_\theta[Z]$	n/a	n/a	4.7073

of the Gaussian components, which is necessary to provide the right aleatoric uncertainty.

Optimizing the negative log-likelihood  $-\log p(y | z)$  is a proper scoring rule, and hence is optimized for calibrated predictions.

Compared to this, the GDA solution (optimized for  $\min \mathbb{H}_\theta[Y, Z]$ ) has a much narrower decision boundary and cannot capture aleatoric uncertainty as well. This is reflected in the higher NLL. Moreover, unlike for LDA, GDA decision boundaries behave differently than one would naively expect due to the untied covariance matrices. They can be curved and the decisions change far away from the data [54].

To show the difference between the two objectives we have marked an ambiguous point near  $(0, -5)$  with a yellow star . Under the first objective  $\min \mathbb{H}_\theta[Y, Z]$ , the point has high aleatoric uncertainty (high entropy), as seen in the left entropy plot while under the second objective ( $\min \mathbb{H}_\theta[Y, Z]$ ) the point is only assigned very low entropy. The GDA optimized for the second objective thus is overconfident.

As above explained above, we do not show an entropy plot of  $Y | Z$  for the third objective  $\min \mathbb{H}_\theta[Z]$  in Fig. 17b because the objective does not depend on  $Y$ , and there are thus no class predictions.

Intuitively, for aleatoric uncertainty, the Gaussian components need to overlap to express high aleatoric uncertainty (uncertain labelling). At the same time, this necessarily provides looser density estimates. On the other hand, the GDA density is much tighter, but this comes at the cost of NLL for classification because it cannot express aleatoric uncertainty that well. Figure 17 visualizes how the objectives trade-off between each other, and why we use the softmax layer trained for  $p(y | z)$  for classification and aleatoric uncertainty, and GDA as density model for  $q(z)$ .

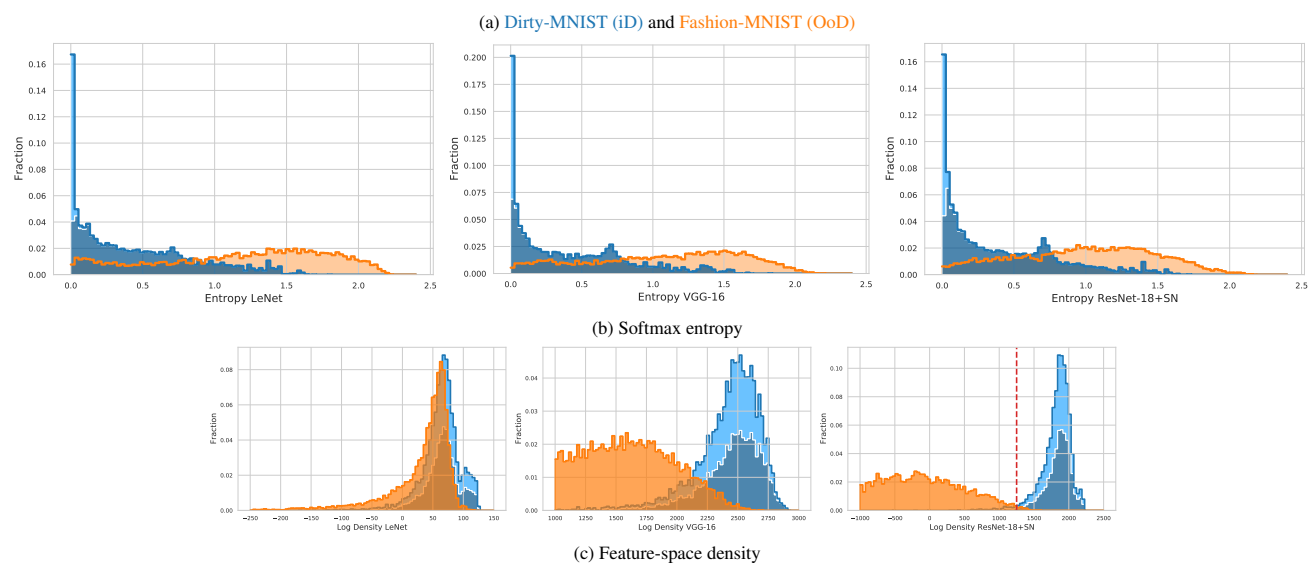


Figure 21. *Disentangling aleatoric and epistemic uncertainty on Dirty-MNIST (iD) and Fashion-MNIST (OoD) (a) requires using softmax entropy (b) and feature-space density (GMM) (c) with appropriate inductive biases (ResNet-18+SN vs LeNet & VGG-16 without them). Enlarged version. (b):* Softmax entropy captures aleatoric uncertainty for iD data (Dirty-MNIST), thereby separating **unambiguous MNIST samples** and **Ambiguous-MNIST samples**. However, iD and OoD are confounded: softmax entropy has arbitrary values for OoD, indistinguishable from iD. (c): With appropriate inductive biases (DDU with ResNet-18+SN), iD and OoD densities do not overlap, capturing epistemic uncertainty. However, without appropriate inductive biases (LeNet & VGG-16), feature density suffers from *feature collapse*: iD and OoD densities overlap. Generally, feature-space density confounds unambiguous and ambiguous iD samples as their densities overlap. **Note:** **Unambiguous MNIST samples** and **Ambiguous-MNIST samples** are shown as stacked histograms with the total fractions adding up to 1 for Dirty-MNIST.