# Understanding Softmax Confidence and Uncertainty

**Tim Pearce**
TSAIL
Tsinghua University

**Alexandra Brintrup**
MAG
University of Cambridge

**Jun Zhu**
TSAIL
Tsinghua University

## Abstract

It is often remarked that neural networks fail to increase their uncertainty when predicting on data far from the training distribution. Yet naively using softmax confidence as a proxy for uncertainty achieves modest success in tasks exclusively testing for this, e.g., out-of-distribution (OOD) detection. This paper investigates this contradiction, identifying two implicit biases that do encourage softmax confidence to correlate with epistemic uncertainty: 1) Approximately optimal decision boundary structure, and 2) Filtering effects of deep networks. It describes why low-dimensional intuitions about softmax confidence are misleading. Diagnostic experiments quantify reasons softmax confidence can fail, finding that extrapolations are less to blame than overlap between training and OOD data in final-layer representations. Pre-trained/fine-tuned networks reduce this overlap.

## 1 Introduction

Many papers studying uncertainty and robustness in deep neural networks observe that unmodified networks fail to be uncertain when predicting on data far from the training distribution, or that they are unable to capture epistemic uncertainty (example critiques in section B.2). Intuition for this deficiency can be communicated through visualisations of low-dimensional input spaces where a network's softmax nonsensically extrapolates with increasing confidence, as in Fig. 1a. This is one motivation for research into methods such as ensembling and Bayesian neural networks. When evaluating the quality of these methods empirically, softmax confidence is invariably used as a baseline – what's surprising is that despite the rhetoric, and without any obvious mechanism enabling capture of epistemic uncertainty, softmax confidence nevertheless performs moderately well.

For example, out-of-distribution (OOD) detection requires predicting if a data point is from the training distribution or not (fig. 1 contains example distributions), where success depends on good epistemic uncertainty estimates. A system incapable of capturing epistemic uncertainty should do no better than random guessing – an AUROC of 50%. Yet directly interpreting softmax confidence as uncertainty can score from 75% to 99% across datasets and architectures (section B.1). Meanwhile, adding modifications purposefully designed to capture uncertainty typically improves AUROC by just 1% to 5%. This suggests that softmax confidence may be more useful as an indicator of epistemic uncertainty than widely thought.

**Contributions:** The goal of this paper is to understand the value of softmax confidence as a proxy for epistemic uncertainty in common vision benchmarks for detection of non-adversarial OOD data. This is valuable in improving our understanding of when, why, and how much we can trust deep learning's most accessible uncertainty estimates. We defer mathematical results (section A), extended related work (section B), and full experimental details (section C) to the appendix, using the main paper to intuitively communicate results. Concrete contributions are as follows:

- **Uncertain regions of the softmax (section 3).** We analytically study the softmax layer, defining regions on the final-layer that will be correctly labelled as OOD.

a) Low-dimensional input, shallow network | b) High-dimensional input, deep network

Input space    Final-layer PCA visualisation      Example inputs    Final-layer PCA visualisation
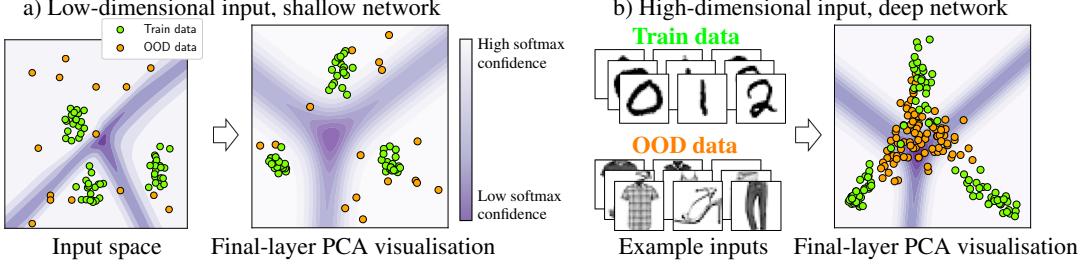
Figure 1: a) In low-dimensions, softmax confidence is unreliable as a measure of uncertainty – OOD data far from the training distribution falls into high confidence extrapolation regions. However, with more complex tasks and deeper networks, implicit biases encourage higher uncertainty for data outside the training distribution. b) For a LeNet trained on three classes of MNIST, OOD data from Fashion MNIST is reliably mapped to the low confidence region of the softmax layer.

- **Decision boundary structure (section 4).** Good structuring of a softmax's decision boundaries is the first important factor in its OOD detection ability, since it enlarges the region that an OOD input is allowed to fall within. We derive a theoretically optimal structuring scheme, finding it closely matches the structure found in trained networks.
- **Deep networks filter for task-specific features (section 5).** Deep networks extract features useful for a specific training task – this serves as an implicit filter of OOD features. We demonstrate why our low-dimensional intuitions mislead us, and propose a more representative mental model.
- **Diagnostic experiments (section 6).** We summarise reasons softmax confidence can fail and measure the importance of each in common benchmark OOD tasks. Further experiments find that pre-trained/fine-tuned networks remedy two of these causes.

## 2   Background

Consider a neural network taking an input $\mathbf{x} \in \mathbb{R}^D$, producing final-layer activations $\mathbf{z} = \psi(\mathbf{x}) \in \mathbb{R}^H$, which are passed through a softmax function to give predictions, $\hat{\mathbf{y}} = \sigma(\mathbf{z}) \in [0,1]^K$,

$$\sigma(\mathbf{z})_i = \frac{\exp \mathbf{w}_i^\mathsf{T} \mathbf{z}}{\sum_{j=1}^K \exp \mathbf{w}_j^\mathsf{T} \mathbf{z}}, \tag{1}$$

where the final-layer weight $\mathbf{W} \in \mathbb{R}^{H \times K}$ can be indexed so $\mathbf{w}_i$ represents column $i$. One can think of appending a one element to $\mathbf{z}$ in place of a bias. Targets for $K$ classes use a one-hot-encoding.

The training distribution is written, $\mathbf{x} \sim \mathcal{D}_{\text{in}}$ or $p_{\text{in}}(\mathbf{x})$. A distribution over the final-layer activations is induced, $\mathbf{z} \sim \mathcal{D}_{\text{in}}$, or, $p_{\text{in}}(\mathbf{z})$. The outlier, or OOD distribution is given by $\mathbf{x} \sim \mathcal{D}_{\text{out}}$ or $p_{\text{out}}(\mathbf{x})$, and similarly for $\mathbf{z}$. A dataset assumes $N$ data points drawn i.i.d..

Two common measures of uncertainty derived from softmax confidence are max predicted probability for any class, and entropy. We define two estimators based on these, ensuring higher is more uncertain. In this paper 'softmax confidence' refers to either estimator.

$$U_{\text{max}}(\mathbf{z}) := -\max_i \sigma(\mathbf{z})_i \qquad U_{\text{entropy}}(\mathbf{z}) := -\sum_{i=1}^K \sigma(\mathbf{z})_i \log \sigma(\mathbf{z})_i \tag{2}$$

We further consider an uncertainty metric operating on the final-hidden layer activations, but not derived from the softmax, similar to Lee et al. [2018]. A probability density is estimated, $\hat{q}(\mathbf{z}) \approx p_{\text{in}}(\mathbf{z})$, and the negative log likelihood is used as an uncertainty score (again higher is more uncertain),

$$U_{\text{density}}(\mathbf{z}) := -\log \hat{q}(\mathbf{z}). \tag{3}$$

In this work we use a Gaussian mixture model (GMM) with $K$ (equal to the number of classes) components to estimate this density (though not restricted to one component per class, and with no restrictions on the covariance structure). Marginal probabilities are given by, $\pi_i := p(\mathbf{y} = i)$. A code snippet in section E.1 summarises the implementation.

$$\hat{q}(\mathbf{z}) = \sum_{i=1}^K \pi_i \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i). \tag{4}$$

**Epistemic vs Aleatoric:** Whilst in everyday usage 'uncertainty' describes the single concept of being unsure, for modelling purposes it is useful to decompose it into two types; aleatoric and epistemic [Der Kiureghian and Ditlevsen, 2008, Gal and Ghahramani, 2015]. **Aleatoric** refers to uncertainty due to classes that overlap in input space, e.g. a network trained on MNIST digits should have aleatoric uncertainty if asked to classify an input appearing between '1' and '7'. **Epistemic** refers to uncertainty about the model or its parameters; if the MNIST network is asked to classify an image of clothing, it should be uncertain due to lack of knowledge about how to handle this type of input. A softmax layer in isolation can learn to output a probability in between 0 and 1 to catch overlapping classes (aleatoric), but fails to decrease its confidence if queried far from the training data (epistemic).

**Related Work:** Neural networks have long combined a softmax output and cross-entropy loss for classification [Bridle, 1989]. Recent work has shown that interpreting softmax confidence as predictive probabilities can have several pitfalls; they are poorly calibrated (generally overconfident) [Guo et al., 2017], and can be easily manipulated by adversarial examples [Nguyen et al., 2015]. It's also been claimed there is no reason to trust them outside of the training distribution – *"[the softmax output] is often erroneously interpreted as model confidence."* [Gal and Ghahramani, 2015].

This paper investigates the last point. Whilst most work on uncertainty and neural networks criticises softmax uncertainties only informally, several works more rigorously demonstrate weaknesses. Hein et al. [2019] prove that any input can be magnified, $\tilde{\mathbf{x}} = \alpha \mathbf{x}$, with large $\alpha > 1$, to produce an arbitrarily high confidence prediction by ReLU networks. Whilst they define $\alpha \mathbf{x}$ as 'far from the training data', our work uses 'far' to mean unfamiliar image datasets with bounded pixel intensity (something trivial to check). Mukhoti et al. [2021] released work concurrently with our own. Whilst they make strong claims about the relationship between softmax confidence and uncertainty, they focus on tasks with high aleatoric uncertainty – we focus on standard datasets which have low aleatoric uncertainty. Section B.3 provides more discussion and a comprehensive literature review.

## 3 Uncertain Regions of the Softmax

This section analytically studies the softmax layer in isolation, defining regions ('valid OOD region') on the softmax layer which an OOD input must fall into to be correctly labelled as OOD.

**Definition 1.** *'Valid OOD region' specifies a region in $\mathbb{R}^H$, such that if an OOD point, $\mathbf{z}' = \psi(\mathbf{x}')$, lies within it, it will be more uncertain than at least $(1-\epsilon)\%$ of the training data (true positive rate, one axis of the ROC). For, $\epsilon \in [0,1]$, some uncertainty estimator $U(\cdot)$, and indicator function $\mathbb{1}(\cdot)$,*

$$\mathcal{R} : \{\mathbf{z}' \in \mathbb{R}^H | \mathbb{E}_{\mathbf{z} \sim \mathcal{D}_{in}} [\mathbb{I}(U(\mathbf{z}') > U(\mathbf{z}))] > 1 - \epsilon\}. \tag{5}$$

Section A.1 provides further insight into softmax confidence, deriving uncertainty vector fields for each of our estimators, formalising common intuition about the danger of softmax extrapolations.

### 3.1 Valid Out-of-Distribution Region Definition

Fig. 2 illustrates valid OOD regions for all estimators. Roughly speaking, for $U_{\text{max}}$ and $U_{\text{entropy}}$, OOD data must fall *closer to a decision boundary* than $(1-\epsilon)\%$ of the training distribution to be in the
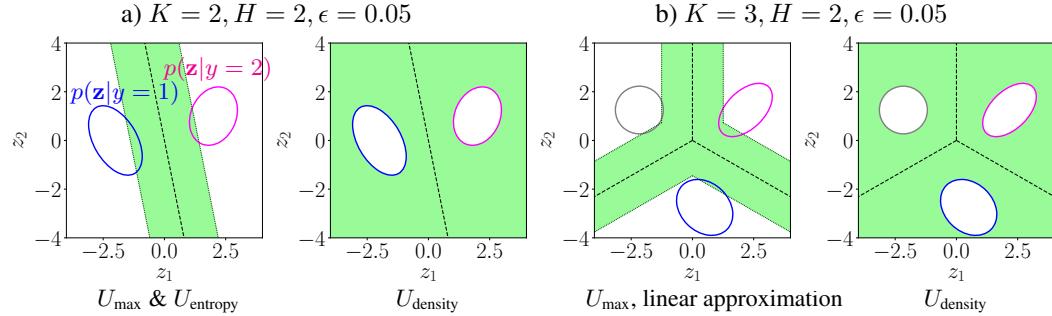


Figure 2: The valid OOD region is shown in green for each uncertainty estimator. Gaussian ellipses capture 95% of clusters. For softmax confidence, OOD data must fall closer to a decision boundary than 95% of training data, whilst the density estimator only demands it fall outside of training clusters.

valid OOD region. For $U_\text{density}$, it is required only that OOD data does not overlap with the training distribution. Observe that if a sample is labelled OOD by $U_\text{max}$ and $U_\text{entropy}$, it will usually also be labelled OOD by $U_\text{density}$, but not vice versa – we exploit this property in section 6's measurements.

These regions are formally defined in theorem 1 for $U_\text{max}$ and $U_\text{entropy}$ for two classes. Making the assumption that final-layer features in the training data follow a mixture of Gaussians allows analytical integration of class clusters in corollary 1.1. Specification of the regions for $U_\text{max}$ and $U_\text{entropy}$ becomes more difficult for higher numbers of classes, since intersections between decision boundary hyperplanes create curved valid OOD regions. However, we define an approximation of the $U_\text{max}$ valid OOD region for general numbers of classes using pairs of linear hyperplanes, offset from the decision boundary by a distance $\alpha(\mathbf{w}_i - \mathbf{w}_j)$ (def. 4). This linear approximation forms a subset of the exact valid OOD region (corollary 2.1), matching the exact valid OOD region for $K = 2$ (proposition 3) and also matching the exact valid OOD region at nearly any point at large magnitudes (theorem 2). Specifying the valid OOD region for $U_\text{density}$ is more straightforward for general classes, $K \geq 2$, as per proposition 2.

## 4 Implicit Bias 1: Approximately Optimal Decision Boundary Structure

The previous section defined the valid OOD region presuming the distribution of training features, $p_\text{in}(\mathbf{z})$, was given. In fact, a neural network learns the transformation, $\mathbf{z} = \psi(\mathbf{x})$, and while $p_\text{in}(\mathbf{x})$ is fixed, the shape and position of $p_\text{in}(\mathbf{z})$ can itself be optimised – this section asks what decision boundary structure emerges when $p_\text{in}(\mathbf{z})$ and $\mathbf{w}_i$'s are jointly optimised. Section 4.1 studies this in an idealised setting, while section 4.2 analyses it empirically in trained networks. Section 4.3 considers how the valid OOD region is affected by structure, and demonstrates experimentally that sub-optimal softmax boundaries can degrade OOD detection.

### 4.1 Optimal Decision Boundary Structure

**Definition 2.** *'Optimal decision boundary structure' (also 'optimal structure') refers to softmax weight vectors that have the following properties. 1) All weight magnitudes are equal, $||\mathbf{w}_i|| = c_1 \; \forall i$. 2) Bias values (absorbed into $\mathbf{w}_i$) are zero. 3) $\mathbf{w}_i$'s are 'evenly distributed'; if $\theta_{i,j}$ represents the angle between $\mathbf{w}_i$ & $\mathbf{w}_j$, then, $\cos\theta_{i,j} = \frac{-1}{K-1} \; \forall i \neq j$. Where relevant, it further means the training distribution features, $p_{in}(\mathbf{z})$, have the following properties. Let $\boldsymbol{\mu}_i := \mathbb{E}_{\mathbf{z}\sim p(\mathbf{z}|\mathbf{y}=i)}[\mathbf{z}]$. 1) Mean vectors for each class are of constant magnitude, $||\boldsymbol{\mu}_i|| = c_2 \; \forall i$, and as large as possible subject to regularisation. 2) Weight and mean vectors of the same class point in the same direction, $\boldsymbol{\mu}_i = c_3 \mathbf{w}_i$. 3) The variance of each cluster is small, $\mathbb{V}ar_{\mathbf{z}\sim p(\mathbf{z}|\mathbf{y}=i)}[\mathbf{z}] \approx \mathbf{0}$.*

Fig. 3 illustrates an optimal structure, compared to a general (non-optimal) structure. Under the assumption that class weights are balanced, $\pi_i = 1/K$, and mean magnitudes are fixed, $||\boldsymbol{\mu}_i|| = c_3$, theorem 3 shows that the optimal structure minimises a regularised cross-entropy loss. This is rather intuitive – training classes should be positioned in small tight clusters as far as possible from each other to minimise this loss. There is inherent symmetry, meaning weights must be evenly distributed. The angle, $\cos\theta_{i,j} = \frac{-1}{K-1}$, arises in spherical coding problems. We will later empirically evidence that optimal structures produce lower cross-entropy losses (fig. 5).



Figure 3: Optimal decision boundary structure for $H = 2, K = 3$. The valid OOD region is overlaid in green.

### 4.2 Empirical Measurement of Decision Boundary Structures in Trained Networks

Analysing properties of final-layer weights in trained networks reveals three key properties of decision boundary structure. Fig. 4 plots histograms of $||\mathbf{w}||$, $\mathbf{b}$, and $\cos\theta_{i,j} \forall i \neq j$ for various architectures and datasets (details in section C.2). 1) Bias values tend to be small. 2) Weight vectors all have similar magnitude. 3) Weight vectors are approximately evenly distributed, with $\cos\theta_{i,j} \approx \frac{-1}{K-1}$.
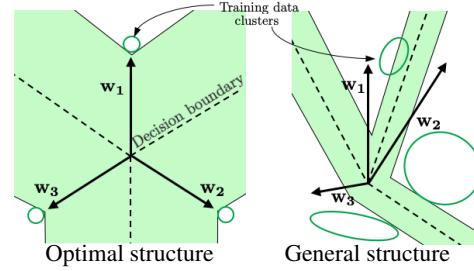
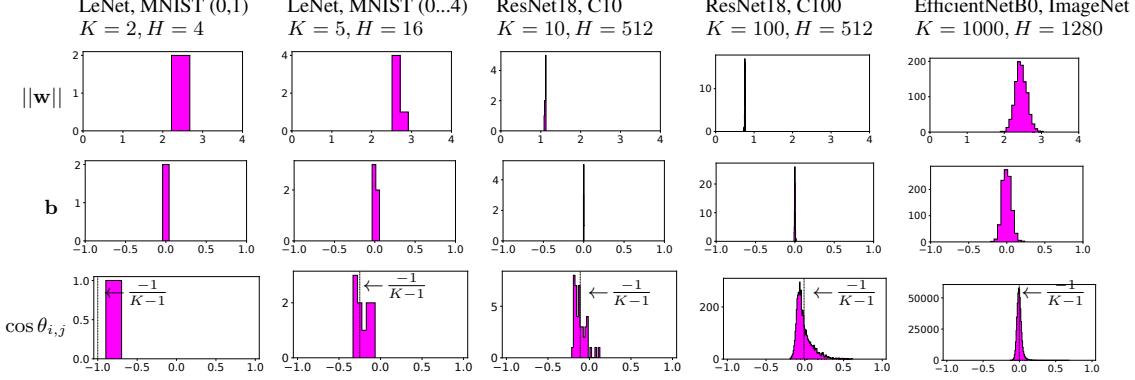| LeNet, MNIST (0,1) $K=2, H=4$ | LeNet, MNIST (0...4) $K=5, H=16$ | ResNet18, C10 $K=10, H=512$ | ResNet18, C100 $K=100, H=512$ | EfficientNetB0, ImageNet $K=1000, H=1280$ |

Figure 4: Histograms of final-layer weight properties for various architectures/datasets show reasonable agreement with the theoretically optimal structure.

These properties match those of the optimal structure reasonably well, though for larger numbers of classes there tends to be feature sharing – there are weak positive correlations between related classes, e.g. different dog breeds, which produces a positively skewed histogram.

Further evidence for agreement can be found in other experiments run in this paper. PCA visualisations of final-layer features in fig. 1 show an evenly-spaced decision boundary, with well-separated data clusters (though these have not collapsed to a point). Fig 15 shows that $\max_i \cos \theta_{\mathbf{z},i} \approx 1$ for most of the training data, verifying that weight vectors and training data clusters are closely aligned, $\boldsymbol{\mu}_i \approx c_3 \mathbf{w}_i$.

### 4.3 The Effect of Decision Boundary Structure on OOD Detection

Having described the decision boundary structure in neural networks, we show why this is a key factor in improving softmax confidence's OOD detection ability – first by quantifying it's effect on the size of the valid OOD region, secondly by running OOD detection experiments with softmax weights frozen in sub-optimal structures.

**Effect on valid OOD region.** Fig. 3 depicts valid OOD regions for optimal and non-optimal structures. The volume of these regions is larger for the optimal structure (corollary 5.1 verifies). This is useful for OOD detection since it creates an increased opportunity for OOD data to fall into these valid regions.

**Counterfactual experiments.** To empirically verify that good structuring is important for good OOD detection, we present a short set of experiments. We take a LeNet architecture ($H = 16$) and freeze the softmax weight vectors in various configurations, then train the rest of the network on a subset of three MNIST classes ($K = 3$) ($\mathcal{D}_{\text{in}}$), measuring OOD detection using Fashion MNIST as OOD ($\mathcal{D}_{\text{out}}$). Trainable refers to not freezing the final layer weights. Further details in section C.3.

Fig. 5 shows results along with illustrations of the structures tested. Whilst structure has limited impact on test accuracy – all structures achieve around 99.5% – there are large differences in OOD detection ability, with only the trainable and optimal structures achieving high AUROC's. Cross-entropy loss is larger for non-optimal structures, supporting our theoretical result that the optimal structure minimises this.

## 5 Implicit Bias 2: Deep Networks Filter for Task-Specific Features

Whilst the decision boundary structure typically learnt by neural networks is helpful for OOD detection, it is not sufficient by itself. This is illustrated by fig. 1; in both tasks shown, the final-layer visualisations suggest good decision boundary structure has been learnt, yet in fig. 1b, the network maps OOD data to the low confidence region of the softmax, whilst in fig. 1a it does not. This section studies a second implicit bias that emerges in deep networks also important in making softmax correlate with epistemic uncertainty.
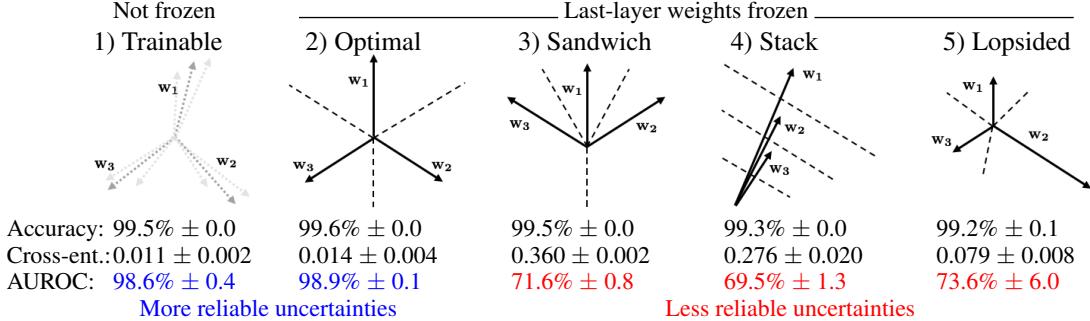
| Not frozen | | Last-layer weights frozen | | |
|---|---|---|---|---|
| 1) Trainable | 2) Optimal | 3) Sandwich | 4) Stack | 5) Lopsided |

| | 1) Trainable | 2) Optimal | 3) Sandwich | 4) Stack | 5) Lopsided |
|---|---|---|---|---|---|
| Accuracy: | $99.5\% \pm 0.0$ | $99.6\% \pm 0.0$ | $99.5\% \pm 0.0$ | $99.3\% \pm 0.0$ | $99.2\% \pm 0.1$ |
| Cross-ent.: | $0.011 \pm 0.002$ | $0.014 \pm 0.004$ | $0.360 \pm 0.002$ | $0.276 \pm 0.020$ | $0.079 \pm 0.008$ |
| AUROC: | $98.6\% \pm 0.4$ | $98.9\% \pm 0.1$ | $71.6\% \pm 0.8$ | $69.5\% \pm 1.3$ | $73.6\% \pm 6.0$ |
| | | More reliable uncertainties | | Less reliable uncertainties | |

Figure 5: Illustration of various softmax structures, along with results when $\mathcal{D}_{\text{in}}$ is $K = 3$ classes of MNIST, and $\mathcal{D}_{\text{out}}$ is Fashion MNIST. Whilst all structures achieve similar test accuracy, AUROC is largest when the softmax weight vectors are either trainable, or fixed in an optimal configuration. Fig. 13 provides PCA visualisations of the trained networks. Mean $\pm 1$ std. err. over three runs.

**What final-layer features are learnt?** Early-layer features of a deep network are known to be generic, while later features evolve in specialised roles to detect and separate the classes present in the training distribution [Zeiler and Fergus, 2014] – the goal of final-layer features is to create clusters that can be linearly separated by a softmax layer. Information about the input that does not directly help with the classification task is thought to be compressed out [Tishby and Zaslavsky, 2015, Saxe et al., 2018]. We therefore view a deep network as a filter, optimised over a training distribution to filter out everything except distinguishing features of the classes.

**How do final-layer features respond to OOD data?** Convolutional filters operate on patches of activations from prior layers. The output of the convolutional filter is maximised when the pattern of patch activations matches the pattern of weights in the filter. Since OOD data is unlikely to contain the distinguishing features that the convolutional filters are trained to extract, OOD activations tend to be of *lower magnitude*. Where some of these distinguishing features are present, it will rarely be in the same combinations seen in training – this results in *unusual patterns* of final-layer activations.

To illustrate this, we train a LeNet with $H = 64$ final-layer neurons. Fig. 6 plots final-layer activations for examples of each training class and OOD inputs, with $\mathcal{D}_{\text{in}}$ and $\mathcal{D}_{\text{out}}$ as denoted. Inputs of the same training class tend to produce consistent patterns of final-layer activations, while OOD inputs produce activation patterns that are of lower magnitude, and/or in unusual combinations. Final-layer activation magnitude is quantified by, $||\mathbf{z}||$, and the activation 'familiarity' by, $\max_i \cos \theta_{i,\mathbf{z}}$.

More comprehensively, fig. 15 plots histograms of $||\mathbf{z}||$ and $\max_i \cos \theta_{i,\mathbf{z}}$ for inputs from the training distribution vs. OOD inputs, for the ResNet18's trained on various datasets as described later in section 6. These plots show $||\mathbf{z}||$ and $\max_i \cos \theta_{i,\mathbf{z}}$ tend to reduce for OOD data.

**What do these OOD final-layer activations do to softmax confidence?** We can write the softmax using, $\mathbf{w}_i^\mathsf{T} \mathbf{z} = ||\mathbf{z}|| \, ||\mathbf{w}_i|| \cos \theta_{i,\mathbf{z}}$, to gain insight into the effect of $||\mathbf{z}||$ & $\cos \theta_{i,\mathbf{z}}$ on softmax confidence,

$$\sigma(\mathbf{z})_i = \frac{\exp ||\mathbf{w}_i|| \, ||\mathbf{z}|| \cos \theta_{i,\mathbf{z}}}{\sum_j \exp ||\mathbf{w}_j|| \, ||\mathbf{z}|| \cos \theta_{j,\mathbf{z}}}. \tag{6}$$

Here, $||\mathbf{z}||$ plays an identical role to the temperature parameter in Platt/temperature scaling, where logits are warmed, $\sigma(\mathbf{z}/T)_i$ for $T > 1$, to produce softmax distributions with lower $U_{\text{max}}$ and $U_{\text{entropy}}$, [Guo et al., 2017] (formally shown in proposition 4 for any decision boundary structure). The role of $\max_i \cos \theta_{i,\mathbf{z}}$ is less straightforward and for general softmax weight structures there is no guarantee decreasing this will result in decreased softmax confidence (proposition 5 provides an example). But crucially, if combined with an optimal decision boundary structure, lower $\max_i \cos \theta_{i,\mathbf{z}}$ does reduce softmax confidence – as proved in theorem 4 for $K = 2$ & $K = 3$, and theorem 5 for large $K$.

## 5.1 Empirical Demonstration of the Filtering Effect

Consider a literal interpretation of a neural network as a filter; when bombarded with *any* allowable input, softmax confidence should be highest on inputs similar to that which the network has been trained on. Typical OOD benchmarks test this over a tiny fraction of the input space covered by
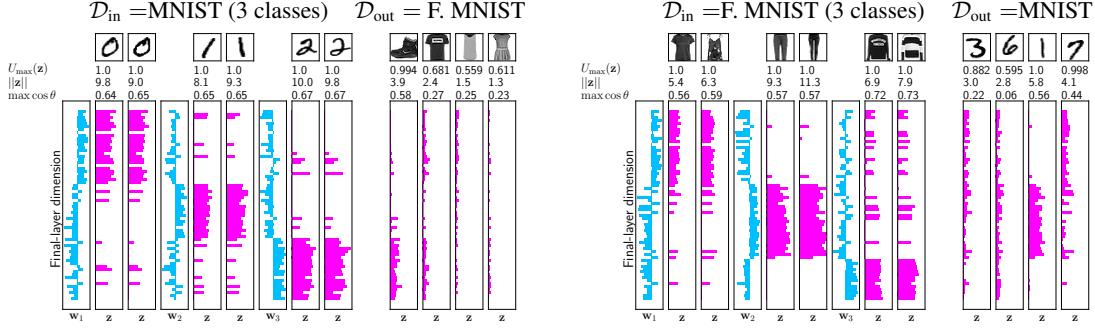
6

Figure 6: Inputs from the same training class produce consistent patterns of final-layer activations, roughly aligning with their weight vectors. OOD inputs tend to produce final-layer activations of lower magnitude and/or in unusual combinations. When combined with good decision boundary structure, this often causes OOD data to fall into the valid OOD region.

some unrelated dataset. We now consider a more exhaustive test; defining a small, bounded training distribution, and in brute-force style generating OOD data sampled uniformly from the input space and examining which inputs are assigned lowest uncertainty by the softmax.

We use MNIST digits, downsampled to $9\times9$ pixels (if smaller the digits become unrecognisable) and binarised, as a training distribution. The entire input space consists of $2^{81}$ possible patterns. We proceed to sample $2^{27} \approx 100$ million inputs uniformly from this space, and collect the most confident examples from each class. Fig. 7 shows that, although noisy, these examples all contain features reminiscent of the training data, evidencing the presence of the filtering effect.

In section E.3 we find empirically that depth is important in strengthening this filtering effect.

## 5.2   Misleading Intuition

Consider again the low-dimensional example in fig. 1a. Such visualisations reinforce the intuition that softmax confidence cannot be trusted 'far from the training data'. Should we be concerned that softmax confidence fails in such a simple setting? We argue that it fails *because* it's so simple.

Firstly, for the filtering effect to emerge, there must be an opportunity to do feature extraction, but in low-dimensions there is no such opportunity – final-layer features are trivial transformations of the input space. Secondly, many of the OOD inputs in fig. 1a can be viewed as magnified training examples, so for any training data point, $\mathbf{x}_i$, an OOD input can be created, $\alpha\mathbf{x}_i$ for $\alpha \gg 1$. Due to ReLU being a homogenous function, magnified inputs produce magnified final-layer features, $\alpha\mathbf{z} \approx \psi(\alpha\mathbf{x})$, leading to increased softmax confidence (as observed by Hein et al. [2019]). In low-dimensions, magnification is *the most obvious way* to create outliers. But this type of OOD input does not translate to high dimensions (should MNIST digits with magnified pixel intensity be labelled
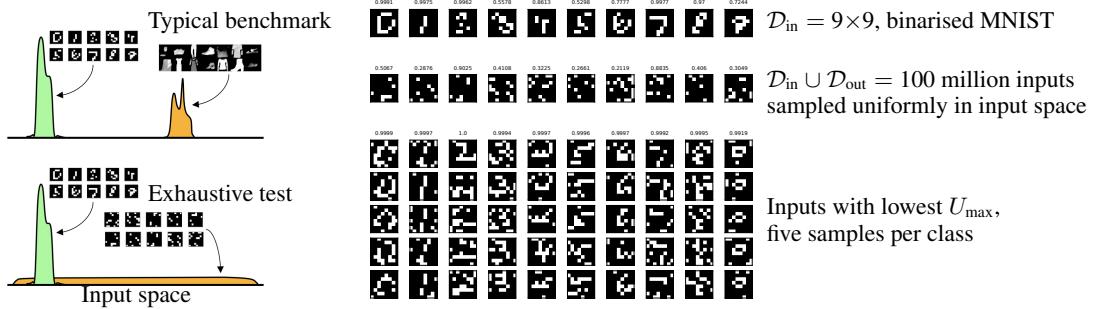


Figure 7: A convolutional network is trained on MNIST, then bombarded with 100 million samples drawn uniformly from the input space. Shown are the top 5 most confident samples for each class. This demonstrates that deep neural networks can filter the entire input space for task-specific features.

7

OOD?) nor does it obviously pose a problem in typical deep learning domains such as images, text, and audio which have bounded input ranges.

**A new mental model.** If low-dimensional examples do not accurately capture the behaviour of softmax confidence in deep networks on complex tasks, how *should* we think about it? We now propose a single interpretable equation summarising this paper's findings.

We take the softmax expressed in terms of magnitudes and angles from eq. 6 and make two strong assumptions, supported by this paper's findings, to give $U_{\text{max mental}}$. 1) The network's decision boundary structure is optimal (section 4.1). 2) $\cos \theta_{i,\mathbf{z}} = \frac{-1}{K-1}$ for all angles except the maximum class. Proposition 6 provides derivation and discussion on assumptions. Note we trade off accuracy for interpretability with these assumptions.

$$U_{\text{max mental}}(\mathbf{z}) := \frac{-1}{1 + (K-1)\exp{-||\mathbf{z}||(\frac{1}{K-1} + \max \cos \theta_{i,\mathbf{z}})}}. \tag{7}$$

Here, $||\mathbf{z}|| \geq 0$ represents the *strength* of distinguishing training features, and $\max \cos \theta_{i,\mathbf{z}} \in [-1, 1]$ represents the *familiarity* of the combination of final-layer features relative to the training data. This model captures that uncertainty increases with a decrease in either of these quantities.

# 6 Failure Causes and their Prevalence

We have so far described two implicit biases in deep networks that encourage OOD data to fall into the softmax's valid OOD region. Softmax confidence is fallible, and we turn our attention to understanding why this mechanism can fail. Figure 16 illustrates a concrete example of failure when trained on specific classes of Fashion MNIST – the network confuses ankle boots for digits. We now define several causes for softmax failure. Some of these issues have been previously noted [Liu et al., 2020, Van Amersfoort et al., 2020, Mukhoti et al., 2021].

**Cause 1) Softmax saturation:** $\max_i \sigma(\mathbf{z})_i$ is typically very high (e.g. fig. 6). Sometimes, the softmax saturates for a proportion of both training and OOD data and $U_{\max}$ is exactly -1.00000; this leads to an inability to create an ordering between samples, and OOD detection deteriorates. Simple fixes include rescaling logits; when problematic we use, $U_{\text{cool}}(\mathbf{z}) := U_{\text{entropy}}(0.1\mathbf{z})$.

**Cause 2a) Softmax extrapolations:** OOD data may be mapped to regions of the softmax of higher confidence than the training distribution – the extrapolation regions.

**Cause 2b) Conflation with aleatoric uncertainty:** Where class distributions within the training data overlap, a network should learn to explicitly map inputs to the appropriate low confidence region of the softmax. This leads to a proportion of the training data legitimately having low softmax confidence. Note that the effect of this is to shrink the valid OOD region, equivalently enlarging the softmax extrapolation region, hence we combine this with cause 2a.
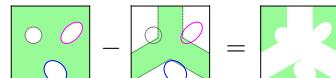
**Cause 3) Feature overlap:** Neural networks are not bijective and there is a risk that a training input and OOD input both map to the same point in the final-layer. Any uncertainty estimator accessing only these final activations will be unable to distinguish these types of input.

**Which of these causes are primarily responsible for overconfident predictions on OOD data?** We first estimate this on typical benchmark tasks – a ResNet18 is trained on one of five datasets (MNIST, Fashion MNIST, SVHN, Cifar10, Cifar100) ($\mathcal{D}_{\text{in}}$) with OOD data taken from the other four datasets ($\mathcal{D}_{\text{out}}$). (Note we convert MNIST and Fashion MNIST to 32×32 RGB format so all datasets are exchangeable.) We secondly test the same protocol, but using an EfficientNetB1 pre-trained on ImageNet, and fine-tuned on five datasets, in each case with the other four used as $\mathcal{D}_{\text{out}}$. Full details are given in section C.5 & C.6.

A perfect uncertainty estimator should achieve 100% AUROC, and we assign responsibility for not reaching this score in the following way. Cause 1 is remedied by using $U_{\text{cool}}$ rather than $U_{\text{entropy}}$. Hence, AUROC using $U_{\text{cool}}$ minus AUROC using $U_{\text{entropy}}$ gives the responsibility of cause 1. Cause 2 is estimated by finding the improvement realised by using $U_{\text{density}}$ rather than $U_{\text{cool}}$, visually:

$$\mathcal{R}_{\text{density}} \bigcap \mathcal{R}'_{\text{entropy}} = \text{Softmax extrapolation region}$$

Any inaccuracy remaining using $U_{\text{density}}$ is attributed to cause 3 – this assumes the density estimator is a good approximation of the true density, $\hat{q}(\mathbf{z}) \approx p_{\text{in}}(\mathbf{z})$. We spent considerable effort tuning the GMM as far as possible, finding it superior to KDE (details in section C.5).

Table 1 displays results. Our numbers in standard training are in line with prior work (table 2). $U_{\text{entropy}}$ is consistently preferable to $U_{\text{max}}$. Saturation is avoided for ResNet18, and only minor for EfficientNet. Importantly, we attribute 4.2% of AUROC drop to softmax extrapolations, and a 7.4% drop to feature overlap – we believe this goes against common intuition in the field. AUROC is high for all estimators in the pre-trained case, and interestingly, $U_{\text{density}}$ provides a close to perfect score across all datasets, largely eliminating feature overlap.

Table 1: AUROC scores for OOD detection, mean $\pm$ 1 standard error. For each $\mathcal{D}_{\text{in}}$ dataset, $\mathcal{D}_{\text{out}}$ comprises the other four datasets. Further breakdown in section C.5 & C.6. Note that cause 1 is the dominant cause of failure in standard training, while pre-trained networks largely avoid all causes.

| | | | | | | Estimated responsibility of each cause | | |
|---|---|---|---|---|---|---|---|---|
| | | $A$ | $B$ | $C$ | $D$ | $C - B$ | $D - C$ | $1.0 - D$ |
| $\mathcal{D}_{\text{in}}$ | Accuracy | $U_{\text{max}}$ | $U_{\text{entropy}}$ | $U_{\text{cool}}$ | $U_{\text{density}}$ | Cause 1 | Cause 2a & 2b | Cause 3 |
| **Standard training**, ResNet18, five runs | | | | | | | | |
| mnist | 99.6% ± 0.0 | 96.3% ± 1.1 | 96.3% ± 1.1 | 96.3% ± 1.1 | 99.5% ± 0.1 | 0.0% | 3.2% | 0.5% |
| f.mnist | 94.5% ± 0.1 | 87.4% ± 1.1 | 88.7% ± 1.1 | 88.7% ± 1.1 | 98.5% ± 0.3 | 0.0% | 9.7% | 1.5% |
| svhn | 96.4% ± 0.1 | 90.7% ± 0.6 | 90.9% ± 0.6 | 90.9% ± 0.6 | 93.2% ± 0.2 | 0.0% | 2.3% | 6.8% |
| c10 | 93.5% ± 0.1 | 85.4% ± 1.1 | 85.7% ± 1.1 | 85.7% ± 1.1 | 89.3% ± 0.4 | 0.0% | 3.6% | 0.7% |
| c100 | 72.0% ± 0.4 | 78.6% ± 0.5 | 80.6% ± 0.5 | 80.6% ± 0.5 | 82.5% ± 0.6 | 0.0% | 2.0% | 17.5% |
| **Mean** | - | 87.7% | 88.4% | 88.4% | 92.6% | 0.0% | **4.2%** | **7.4%** |
| **Pre-trained/fine-tuned**, EfficientNetB1, three runs | | | | | | | | |
| Satellite | 97.7% ±0.1 | 99.7% ±0.0 | 99.9% ±0.0 | 99.9% ±0.0 | 99.9% ±0.0 | 0.0% | 0.0% | 0.1% |
| Cancer | 95.2% ±0.1 | 98.4% ±0.0 | 99.3% ±0.0 | 99.8% ±0.0 | 100.0% ±0.0 | 0.5% | 0.2% | 0.0% |
| Pets | 92.5% ±0.1 | 100.0% ±0.0 | 100.0% ±0.0 | 100.0% ±0.0 | 100.0% ±0.0 | 0.0% | 0.0% | 0.0% |
| Flowers | 94.7% ±0.0 | 99.9% ±0.0 | 99.9% ±0.0 | 99.9% ±0.0 | 100.0% ±0.0 | 0.0% | 0.1% | 0.0% |
| Beans | 94.1% ±0.2 | 96.8% ±0.1 | 98.2% ±0.1 | 99.0% ±0.0 | 100.0% ±0.0 | 0.8% | 1.0% | 0.0% |
| **Mean** | - | 99.0% | 99.5% | 99.7% | 100.0% | 0.3% | **0.2%** | **0.0%** |

# 7   Discussion & Conclusion

**Summary.** This paper sought to understand the value of softmax confidence as a proxy for epistemic uncertainty in the detection of OOD data. It uncovered two implicit biases baked into unmodified networks which together produce a rudimentary mechanism for uncertainty estimation. We argued that common intuition around softmax confidence is misleading. Carefully designed diagnostic experiments on standard benchmark tasks suggested final-layer feature overlap is more responsible for failures than softmax extrapolations.

**Limitations & impact.** This paper has contributed to an understanding of deep learning's most accessible uncertainty estimates; we expect this to be helpful for practitioners and researchers. Whilst our paper has explained why using softmax confidence for uncertainty estimation might have a sounder basis than widely believed, we remind readers of several things. Our results have said nothing about calibration of the softmax probabilities, nor it's effectiveness against adversarial examples, nor about uncertainty in regression tasks. Softmax confidence remains an imperfect measure of uncertainty, and caution should be applied when used in real-world applications. Our study has focused on one specific type of OOD data commonly used in benchmarking OOD detection methods – $\mathcal{D}_{\text{in}}$ is the input distribution of the task of interest (with low aleatoric uncertainty), and $\mathcal{D}_{\text{out}}$ is an arbitrarily selected, unrelated, dataset. Softmax confidence may be less effective on other types of OOD data.

**Conclusions.** We close with three thoughts. 1) The ability to 'capture epistemic uncertainty' is often presented as a binary property, whilst our work described a mechanism providing partial capture – perhaps this property should be interpreted in a softer way. 2) It may be revealing to study ensembling, Bayesian neural networks, and MC Dropout, in the context of strengthening the implicit biases this paper has described, rather than as adding fundamental properties by themselves. 3) Addressing final-layer feature overlap should be a priority for the field. There are at least two approaches; modifying networks to encourage bijective behaviour, e.g. [Liu et al., 2020, Van Amersfoort et al., 2020], or learning more diverse representations, e.g. [Hendrycks et al., 2019a,b, 2020]. Given the trend across machine learning for large-scale pre-training, we see the second route as particularly promising.

# References

Michael A Alcorn, Qi Li, Zhitao Gong, Chengfei Wang, Long Mai, Wei-Shinn Ku, and Anh Nguyen. Strike (with) a Pose: Neural Networks Are Easily Fooled by Strange Poses of Familiar Objects. 2018. URL `https://arxiv.org/abs/1811.11553`.

Bhaskar Bagchi. How to stay away from each other in a spherical universe. *Resonance*, 2(10):38–45, 1997. ISSN 0971-8044. doi: 10.1007/bf02835977.

Eric B Baum and Frank Wilczek. Supervised Learning of Probability Distributions by Neural Networks. *Neural Information Processing Systems*, 1987.

Christopher M Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006. ISBN 9788578110796. doi: 10.1017/CBO9781107415324.004.

T. E. Boult, S. Cruz, A.R. Dhamija, M. Gunther, J. Henrydoss, and W.J. Scheirer. Learning and the Unknown: Surveying Steps toward Open World Recognition. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33:9801–9807, 2019. ISSN 2159-5399. doi: 10.1609/aaai.v33i01. 33019801.

Johann S. Brauchart and Peter J. Grabner. Distributing many points on spheres: Minimal energy and designs. *Journal of Complexity*, 31(3):293–326, 2015. ISSN 10902708. doi: 10.1016/j.jco.2015. 02.003. URL `http://dx.doi.org/10.1016/j.jco.2015.02.003`.

John S. Bridle. Training Stochastic Model Recognition Algorithms as Networks can lead to Maximum Mutual Information Estimation of Parameters. In *Advances in Neural Information Processing Systems 2*, 1989.

John S. Bridle. Probabilistic Interpretation of Feedforward Classification Network Outputs, with Relationships to Statistical Pattern Recognition. *Neurocomputing*, (C):227–236, 1990. doi: 10.1007/978-3-642-76153-9_28.

H. S. M. Coxeter, J. H. Conway, and N. J. A. Sloane. *Sphere Packings, Lattices and Groups.*, volume 96. 1989. ISBN 9781441931344. doi: 10.2307/2323992.

A. Der Kiureghian and O. Ditlevsen. Aleatoric or epistemic ? Does it matter ? *Special Workshop on Risk Acceptance and Risk Communication*, pages 1–13, 2008.

Yarin Gal and Zoubin Ghahramani. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. In *Proceedings of the 33rd International Conference on Machine Learning*, 2015. ISBN 1506.02142. doi: 10.1109/TKDE.2015.2507132. URL `http://arxiv.org/abs/1506.02142`.

Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep Bayesian Active Learning with Image Data. In *ICML*, 2017. ISBN 9781510855144. doi: 10.17863/CAM.11070. URL `http://arxiv.org/abs/1703.02910`.

Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. *34th International Conference on Machine Learning, ICML 2017*, 2017.

Matthias Hein, Maksym Andriushchenko, and Julian Bitterwolf. Why ReLU networks yield high-confidence predictions far away from the training data and how to mitigate the problem. In *CVPR*, 2019.

Dan Hendrycks and Kevin Gimpel. A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks. In *ICLR*, 2017a. URL `http://arxiv.org/abs/1610.02136`.

Dan Hendrycks and Kevin Gimpel. a Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks. *ICLR*, 2017b. URL `https://arxiv.org/pdf/1610.02136.pdf`.

Dan Hendrycks, Kimin Lee, and Mantas Mazeika. Using pre-training can improve model robustness and uncertainty. *36th International Conference on Machine Learning, ICML 2019*, 2019a.

Dan Hendrycks, Mantas Mazeika, Saurav Kadavath, and Dawn Song. Using Self-Supervised Learning Can Improve Model Robustness and Uncertainty. (NeurIPS), 2019b. URL `http://arxiv.org/abs/1906.12340`.

Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam Dziedzic, Rishabh Krishnan, and Dawn Song. Pretrained Transformers Improve Out-of-Distribution Robustness. In *ACL*, 2020.

Tom Heskes. Practical confidence and prediction intervals. In *Advances in Neural Information Processing Systems 9*, 1996.

Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles. In *NIPS*, 2017.

Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Advances in Neural Information Processing Systems*, 2018. ISSN 10495258.

Jeremiah Zhe Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax-Weiss, and Balaji Lakshminarayanan. Simple and principled uncertainty estimation with deterministic deep learning via distance awareness. *NeurIPS*, 2020. ISSN 23318422.

Zhiyun Lu, Eugene Ie, and Fei Sha. Uncertainty Estimation with Infinitesimal Jackknife, Its Distribution and Mean-Field Approximation. *Arxiv 2006.07584*, 2021. URL `http://arxiv.org/abs/2006.07584`.

David J. C. MacKay. A Practical Bayesian Framework for Backpropagation Networks. *Neural Computation*, 4(3):448–472, 1992. ISSN 0899-7667. doi: 10.1162/neco.1992.4.3.448. URL `http://www.mitpressjournals.org/doi/10.1162/neco.1992.4.3.448`.

Andrey Malinin and Mark Gales. Predictive Uncertainty Estimation via Prior Networks. In *Neural Information Processing Systems*, 2018. URL `http://arxiv.org/abs/1802.10501`.

Jishnu Mukhoti, Andreas Kirsch, Joost van Amersfoort, Philip H. S. Torr, and Yarin Gal. Deterministic Neural Networks with Appropriate Inductive Biases Capture Epistemic and Aleatoric Uncertainty. *ArXiv*, 2021. URL `http://arxiv.org/abs/2102.11582`.

Oleg R. Musin and Alexey S. Tarasov. The Tammes problem for N = 14. *Experimental Mathematics*, 24(4):460–468, 2015. ISSN 1944950X. doi: 10.1080/10586458.2015.1022842.

Gregory Naitzat, Andrey Zhitnikov, and Lek Heng Lim. Topology of deep neural networks. *Journal of Machine Learning Research*, 21:1–40, 2020. ISSN 15337928.

Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 07-12-June:427–436, 2015. ISSN 10636919. doi: 10.1109/CVPR.2015.7298640.

Emin Orhan. Robustness properties of Facebook's ResNeXt WSL models. *ArXiv*, 2019.

Utku Ozbulak, Wesley de Neve, and Arnout van Messem. How the softmax output is misleading for evaluating the strength of adversarial examples. *NeurIPS*, pages 1–9, 2018. ISSN 23318422.

Tim Pearce. *Uncertainty in Neural Networks; Bayesian Ensembles, Priors & Prediction Intervals*. PhD thesis, University of Cambridge, 2021.

Michael D. Richard and Richard P. Lippmann. Neural Network Classifiers Estimate Bayesian a posteriori Probabilities. *Neural Computation*, 3(4):461–483, 1991. ISSN 0899-7667. doi: 10.1162/neco.1991.3.4.461.

Andrew M. Saxe, Yamini Bansal, Joel Dapello, Madhu Advani, Artemy Kolchinsky, Brendan D. Tracey, and David D. Cox. On the information bottleneck theory of deep learning. *ICLR*, 2018. ISSN 00029378. doi: 10.1016/0002-9378(95)90441-7.

Alireza Shafaei, Mark Schmidt, and James J. Little. A less biased evaluation of out-of-distribution sample detectors. *British Machine Vision Conference*, 2019. ISSN 23318422.

Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *2nd International Conference on Learning Representations, ICLR 2014 - Conference Track Proceedings*, pages 1–10, 2014.

Robert Tibshirani. A Comparison of Some Error Estimates for Neural Network Models. *Neural Computation*, 8:152–163, 1996.

Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. *2015 IEEE Information Theory Workshop, ITW 2015*, 2015. doi: 10.1109/ITW.2015.7133169.

Joost Van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Simple and scalable epistemic uncertainty estimation using a single deep deterministic neural network. *ICML*, 2020. ISSN 23318422.

Florian Wenzel, Kevin Roth, Bastiaan S. Veeling, Jakub Świątkowski, Linh Tran, Stephan Mandt, Jasper Snoek, Tim Salimans, Rodolphe Jenatton, and Sebastian Nowozin. How Good is the Bayes Posterior in Deep Neural Networks Really? In *ICML*, 2020. URL `http://arxiv.org/abs/2002.02405`.

Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. *European Conference on Computer Vision, ECCV*, 2014. ISSN 16113349. doi: 10.1007/978-3-319-10590-1_53.

# A Proofs

This section contains all mathematical results and proofs, organised as follows.

- **Section A.1.** Derives uncertainty vector fields for all uncertainty estimators.
- **Section A.2.** Derives valid OOD regions for each estimator under various assumptions:
  - Theorem 1: $U_{\text{max}}$ and $U_{\text{entropy}}$ for $K = 2$ (assume weight vectors, $\mathbf{w}_1 = -\mathbf{w}_2$).
  - Corollary 1.1: $U_{\text{max}}$ and $U_{\text{entropy}}$ for $K = 2$ (Gaussian and linearly separable assumption).
  - Proposition 2: $U_{\text{density}}$ for $K \geq 2$ (Gaussian assumption).
  - Definition 4: Describes a linear approximation to the exact valid OOD region for general $K \geq 2$.
  - Proposition 3: Shows the linear approximation is equal to the exact valid OOD region for $K = 2$ (assume weight vectors, $\mathbf{w}_1 = -\mathbf{w}_2$).
  - Theorem 2: For $K \geq 3$, shows the exact valid OOD region converges to the linear approximation at large magnitudes at nearly all points (except close to intersections between decision hyperplanes).
  - Corollary 2.1: For $K \geq 3$, shows that the linear approximate valid OOD region is a subset of the exact valid OOD region.
- **Section A.3.** Considers the optimal decision boundary structure and related results.
  - Theorem 3: Shows that the optimal structure minimises a regularised cross-entropy loss. Assumes class clusters follow a Gaussian distribution, $H \geq K - 1$, and $||\boldsymbol{\mu}_i|| = c_2 \ \forall i$.
  - Proposition 4: For any structure, $U_{\text{max}}$ increases with a decrease in $||\mathbf{z}||$.
  - Theorem 4: Shows a decrease in $\max_i \cos \theta_{\mathbf{z},i}$ always results in an increase in $U_{\text{max}}$ for $K = 2$ and $K = 3$.
  - Theorem 5: Shows a decrease in $\max_i \cos \theta_{\mathbf{z},i}$ always results in an increase in $U_{\text{max}}$ for large $K$ (assumes that $\frac{-1}{K-1} \approx 0$).
  - Corollary 5.1: Shows that the optimal structure maximises the linear approximate valid OOD region for $U_{\text{max}}$. Assumes that $U_{\text{max}}$ increases for a decrease in $\max_i \cos \theta_{\mathbf{z},i}$ for all $K \geq 2$, and also that $||\boldsymbol{\mu}_i|| = c_2 \ \forall i$.
  - Proposition 5: Shows that a decrease in $\max_i \cos \theta_{\mathbf{z},i}$ does not always produce an increase in $U_{\text{max}}$ for structures that are not optimal.
- **Section A.4.** Derives the proposed mental model $U_{\text{max mental}}$ (proposition 6) under two strong assumptions, which are discussed.

### A.1 Vector Fields

**Proposition 1.** *The following equations calculate the direction of maximum uncertainty given any point, $\mathbf{z}$, in the final-hidden layer space for each uncertainty estimator.*

$$\frac{\partial U_{max}(\mathbf{z})}{\partial \mathbf{z}} = \sigma(\mathbf{z})_i \sum_{j=1}^{K} \sigma(\mathbf{z})_j (\mathbf{w}_j - \mathbf{w}_i) \ \ where, \ i = argmax_i \sigma(\mathbf{z})_i \tag{8}$$

$$\frac{\partial U_{entropy}(\mathbf{z})}{\partial \mathbf{z}} = \sum_{i=1}^{K} (\log \sigma(\mathbf{z})_i + 1) \sigma(\mathbf{z})_i \sum_{j=1}^{K} \sigma(\mathbf{z})_j (\mathbf{w}_j - \mathbf{w}_i) \tag{9}$$

$$\frac{\partial U_{density}(\mathbf{z})}{\partial \mathbf{z}} \propto \frac{1}{\hat{q}(\mathbf{z})} \sum_{i=1}^{K} \pi_i \exp \left[ -\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu}_i)^\mathsf{T} \boldsymbol{\Sigma}_i^{-1}(\mathbf{z} - \boldsymbol{\mu}_i) \right] \boldsymbol{\Sigma}_i^{-1}(\mathbf{z} - \boldsymbol{\mu}_i) \tag{10}$$

*Proof.* These are derived through application of the chain rule, applied to the uncertainty estimators defined in the main text. □

**Remark.** *Fig. 8 visualises the equations from proposition 1. $\partial U_{max}(\mathbf{z})/\partial \mathbf{z}$ is a linear combination of all weight vectors, which is the root cause of softmax confidence being inexpressive, with $-\mathbf{w}_i$ its primary direction. $\partial U_{entropy}(\mathbf{z})/\partial \mathbf{z}$ has similar limitations, producing a smoothed version of $\partial U_{max}(\mathbf{z})/\partial \mathbf{z}$. Meanwhile, $\partial U_{density}(\mathbf{z})/\partial \mathbf{z}$, always points away from the data clusters, $\mathbf{z} - \boldsymbol{\mu}_i$, weighted by the distance from each component.*
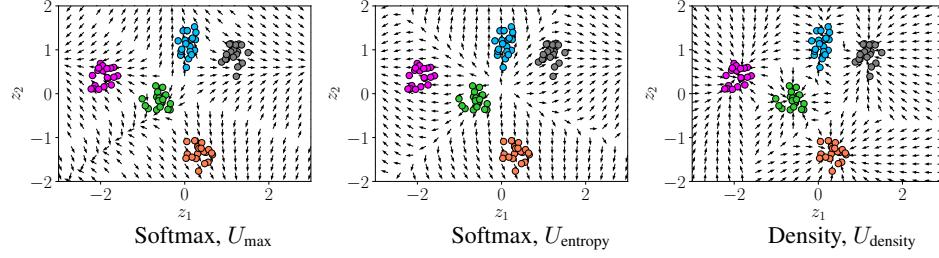


Figure 8: Vector fields showing direction of maximum certainty (magnitude ignored).

## A.2 Valid OOD Regions

**Theorem 1.** *Consider $K = 2$ classes, and softmax weight vectors, $\mathbf{w}_1 = -\mathbf{w}_2$. The valid OOD region (def. 1) for $U_{max}$ and $U_{entropy}$ is the region between two hyperplanes, each parallel to the decision boundary,*

$$\mathcal{R}_{K=2} : \{\mathbf{z} \in \mathbb{R}^{\mathbb{H}} | \mathbf{w}_1^{\mathsf{T}} (\mathbf{z}_0 - \alpha \mathbf{w}_1) < \mathbf{w}_1^{\mathsf{T}} \mathbf{z} < \mathbf{w}_1^{\mathsf{T}} (\mathbf{z}_0 + \alpha \mathbf{w}_1)\}, \tag{11}$$

*where $\mathbf{z}_0$ is some point on the decision boundary. The offset, $\alpha > 0$, is determined by solving,*

$$\pi_1 p_{in}(\mathbf{z} \notin \mathcal{R}_{K=2} | \mathbf{y} = 1) + \pi_2 p_{in}(\mathbf{z} \notin \mathcal{R}_{K=2} | \mathbf{y} = 2) = 1 - \epsilon. \tag{12}$$

**Remark.** *Note that by using a softmax with two weight vectors for two classes, we have an over-parameterised model, and it seems reasonable to assume that the optimised weight vectors will be, $\mathbf{w}_1 \approx -\mathbf{w}_2$ (a softmax with this constraint enforced is equivalent to a scaled sigmoid). Note also that the weight vectors $\mathbf{w}_1$ and $\mathbf{w}_2$ are perpendicular to the decision boundary hyperplane [Bishop, 2006].*

*Proof.* We first consider $U_{\max}$. Begin observing that for the case of two classes, $\sigma(\mathbf{z})_1 = 1 - \sigma(\mathbf{z})_2$. If, $\mathbf{w}_1 = -\mathbf{w}_2$, and if $1 = \text{argmax}_i \sigma(\mathbf{z})_i$, eq. 8 reduces to,

$$\frac{\partial U_{\max}(\mathbf{z})}{\partial \mathbf{z}} = \sigma(\mathbf{z})_1 \left[ (1 - \sigma(\mathbf{z})_1)(\mathbf{w}_2 - \mathbf{w}_1) \right] \tag{13}$$

$$= 2\sigma(\mathbf{z})_1 (1 - \sigma(\mathbf{z})_1)(-\mathbf{w}_1). \tag{14}$$

This shows that for the region where $\sigma(\mathbf{z})_1 > \sigma(\mathbf{z})_2$, uncertainty always increases in the direction of $-\mathbf{w}_1$ (toward the decision boundary). The opposite direction, $-\mathbf{w}_2$, is true for $\sigma(\mathbf{z})_2 > \sigma(\mathbf{z})_1$.

Observe also that $U_{\max}$ is symmetric about the decision boundary (lemma 1.1).

Consider two hyperplanes, parallel to the decision surface, and offset by some perpendicular vector $\alpha \mathbf{w}_1$ and $-\alpha \mathbf{w}_1$, defined by the equations,

$$\mathbf{w}_1^{\mathsf{T}} \mathbf{z} = \mathbf{w}_1^{\mathsf{T}} (\mathbf{z}_0 + \alpha \mathbf{w}_1) \tag{15}$$

$$\mathbf{w}_1^{\mathsf{T}} \mathbf{z} = \mathbf{w}_1^{\mathsf{T}} (\mathbf{z}_0 - \alpha \mathbf{w}_1) \tag{16}$$

where $\mathbf{z}_0$ describes some point on the decision boundary, and $\alpha \in \mathbb{R}_+$.

Define the region enclosed by these two planes as $\mathcal{R}_{K=2}$,

$$\mathcal{R}_{K=2} : \{\mathbf{z} \in \mathbb{R}^{\mathbb{H}} | \mathbf{w}_1^{\mathsf{T}} (\mathbf{z}_0 - \alpha \mathbf{w}_1) < \mathbf{w}_1^{\mathsf{T}} \mathbf{z} < \mathbf{w}_1^{\mathsf{T}} (\mathbf{z}_0 + \alpha \mathbf{w}_1)\}. \tag{17}$$

Due to the way we have designed this volume,

$$U_{\max}(\mathbf{z}_i) > U_{\max}(\mathbf{z}_j) \quad \forall \mathbf{z}_i \in \mathcal{R}_{K=2}, \ \forall \mathbf{z}_j \notin \mathcal{R}_{K=2}. \tag{18}$$

In words; any point inside $\mathcal{R}_{K=2}$ will be more uncertain than any point outside of it. We can use this observation to define the valid OOD region for $U_{\max}$, by choosing $\alpha$ so that the proportion of the training distribution outside the region, $p_{\text{in}}(\mathbf{z} \notin \mathcal{R}_{K=2})$, and hence with lower uncertainty than any point inside the region, is $1 - \epsilon$.

Since $p_{\text{in}}(\mathbf{z} \notin \mathcal{R}_{K=2})$ will always decrease with increasing $\alpha$, this creates a requirement finding $\alpha$ so that,

$$p_{\text{in}}(\mathbf{z} \notin \mathcal{R}_{K=2}) = \pi_1 p_{\text{in}}(\mathbf{z} \notin \mathcal{R}_{K=2} | \mathbf{y} = 1) + \pi_2 p_{\text{in}}(\mathbf{z} \notin \mathcal{R}_{K=2} | \mathbf{y} = 2) = 1 - \epsilon \tag{19}$$

This proof can be extended to $U_{\text{entropy}}$ via lemma 1.2.

$\square$

**Lemma 1.1.** *Assume two classes, and, $\mathbf{w}_1 = -\mathbf{w}_2$. The softmax is symmetric about the decision boundary, so for some point on the decision boundary, $\mathbf{z}_0$, such that $\sigma(\mathbf{z}_0)_1 = \sigma(\mathbf{z}_0)_2 = 0.5$, it holds that $\sigma(\mathbf{z}_0 + \alpha \mathbf{w}_1)_1 = \sigma(\mathbf{z}_0 - \alpha \mathbf{w}_1)_2$, for any $\alpha \in \mathbb{R}$.*

*Proof.* Simply substitute $\mathbf{z}_0 + \alpha \mathbf{w}_1$ and $\mathbf{z}_0 - \alpha \mathbf{w}_1$ into eq. 1 to verify. $\square$

**Lemma 1.2.** *Theorem 1 also holds for $U_{\text{entropy}}$.*

*Proof.* The property required of $U_{\text{max}}$ to prove theorem 1, was that when $\sigma(\mathbf{z})_1 > \sigma(\mathbf{z})_2$, uncertainty always decreases in the direction of $-\mathbf{w}_1$. From eq. 8 and 9, it's straightforward to see that provided $\mathbf{w}_1 = -\mathbf{w}_2$ (as it is by assumption), this also holds for $U_{\text{entropy}}$.

$\square$

**Definition 3.** *We relax the definition of 'linearly separable' slightly. Generally two sets of points, $\mathbf{Z}_1$ and $\mathbf{Z}_2$, are linearly separably if at least one hyperplane exists dividing the sets, i.e. $\mathbf{w}$ exists such that, $\mathbf{w}^\intercal \mathbf{z}_i > c \;\; \forall \mathbf{z}_i \in \mathbf{Z}_1$ and $\mathbf{w}^\intercal \mathbf{z}_j < c \;\; \forall \mathbf{z}_j \in \mathbf{Z}_2$, for some constant $c \in \mathbb{R}$.*

*Our results assume densities rather than finite sets of data points, so this does not obviously apply, but we relax the concept slightly. For two densities, $p(\mathbf{z}|\mathbf{y} = 1)$ and $p(\mathbf{z}|\mathbf{y} = 2)$ to be linearly separable, there must exist a hyperplane such that the density of each class falling on the opposite side of the hyperplane is negligible, i.e. $\mathbf{w}$ and $c$ exist such that, $p(\mathbf{w}^\intercal \mathbf{z} < c|\mathbf{y} = 2) \approx 1$ and $p(\mathbf{w}^\intercal \mathbf{z} > c|\mathbf{y} = 2) \approx 0$, while $p(\mathbf{w}^\intercal \mathbf{z} > c|\mathbf{y} = 1) \approx 1$ and $p(\mathbf{w}^\intercal \mathbf{z} < c|\mathbf{y} = 1) \approx 0$.*

**Corollary 1.1.** *This corollary extends theorem 1. Under conditions as before, and further assuming that each class is represented by a single multivariate Gaussian, $p(\mathbf{z}|\mathbf{y} = i) = \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$, and linearly separable (def. 3), the value of $\alpha$ required to define the valid OOD region for $U_{\text{max}}$ and $U_{\text{entropy}}$ is given by solving the below,*

$$\pi_1 \text{erf}\left( \frac{\mathbf{w}_1^\intercal \boldsymbol{\mu}_i - \mathbf{w}_1^\intercal(\mathbf{z}_0 + \alpha \mathbf{w}_1)}{\sqrt{2\mathbf{w}_1^\intercal \boldsymbol{\Sigma}_i \mathbf{w}_1}} \right) + \pi_2 \text{erf}\left( \frac{\mathbf{w}_1^\intercal \boldsymbol{\mu}_i - \mathbf{w}_1^\intercal(\mathbf{z}_0 + \alpha \mathbf{w}_1)}{\sqrt{2\mathbf{w}_1^\intercal \boldsymbol{\Sigma}_i \mathbf{w}_1}} \right) = 1 - 2\epsilon. \quad (20)$$

**Remark.** *Justification for the Gaussian assumption comes from topological studies finding class components are well modelled by single connected components [Naitzat et al., 2020]. Secondly, similar assumptions have been popularised in prior work, e.g. [Lee et al., 2018].*

*Proof.* This result requires deriving an analytical expression for eq. 19. By assumption we have that $p_{\text{in}}(\mathbf{z}|\mathbf{y} = i) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$, which allows us to analytically evaluate the proportion of a component falling either side of a hyperplane – technically, integrating over a half space. Our assumption on the densities being linearly separable (def. 3) simplifies this further, since we need only consider the integral over a single Gaussian component on each side of the decision boundary, with negligible impact on the result.

The required computation to integrate over the training data, $\mathbf{z} \sim \mathcal{D}_{\text{in}}$, for one component is[1] ,

$$p_{\text{in}}(\mathbf{z} \notin \mathcal{R}_{K=2}|\mathbf{y} = 1) = \frac{1}{(2\pi)^{H/2}(\det \boldsymbol{\Sigma}_1)^{1/2}} \int_{\mathbb{R}_\alpha} \exp\left[ -\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu}_1)^\intercal \boldsymbol{\Sigma}^{-1}(\mathbf{z} - \boldsymbol{\mu}_1) \right] d\mathbf{z} \quad (21)$$

$$= \frac{1}{2}\left[ 1 - \text{erf}\left( -\frac{\mathbf{w}_1^\intercal \boldsymbol{\mu}_i - \mathbf{w}_1^\intercal(\mathbf{z}_0 + \alpha \mathbf{w}_1)}{\sqrt{2\mathbf{w}_1^\intercal \boldsymbol{\Sigma}_i \mathbf{w}_1}} \right) \right]. \quad (22)$$

Combining with the second component, $p_{\text{in}}(\mathbf{z} \notin \mathcal{R}_{K=2}|\mathbf{y} = 2)$, and rearranging, produces the main result.

$\square$

**Proposition 2.** *Consider $K \geq 2$ classes, each represented by a single multivariate Gaussian. The valid OOD region (def. 1) for $U_{\text{density}}$ is given by,*

$$\mathcal{R}_{density} : \{ \mathbf{z} \in \mathbb{R}^H | \bigcap_i (\mathbf{z} - \boldsymbol{\mu}_i)^\intercal \boldsymbol{\Sigma}_i (\mathbf{z} - \boldsymbol{\mu}_i) > c_i \} \quad (23)$$

*where each $c_i$ depends on $\epsilon, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i$.*

*Proof.* Similar to eq. 19, we require the contour defining,

$$p_{\text{in}}(\mathbf{z} \notin \mathcal{R}_{\text{density}}) = \sum_i \pi_i p_{\text{in}}(\mathbf{z} \notin \mathcal{R}_{\text{density}}|\mathbf{y} = i) = 1 - \epsilon. \quad (24)$$

The Gaussian assumption allows this to be done analytically. For a single Gaussian component, the $1 - \epsilon\%$ contour is a hyperellipsoid described by $(\mathbf{z} - \boldsymbol{\mu}_i)^\intercal \boldsymbol{\Sigma}_i (\mathbf{z} - \boldsymbol{\mu}_i) > c_i$, where $c_i$ depends on $\epsilon, \boldsymbol{\mu}_i, \Sigma_i$. Our valid OOD region is then just the intersection over all components. $\square$

---

[1] https://math.stackexchange.com/questions/556977/gaussian-integrals-over-a-half-space

**Definition 4.** *We define the 'linear approximation of the valid OOD region' (also 'linear approximation') for $U_{max}$ as an approximation to the exact valid OOD region (def. 1), using only linear hyperplanes with structure as described following. This applies to multiple classes, $K \geq 2$, and general weight structures.*

*We begin with a description in words: take each decision hyperplane and create a pair of hyperplanes, offset in opposite directions by some perpendicular distance. The union of the regions enclosed by each pair will form the linear approximation of the valid OOD region. This is visualised in fig. 9.*
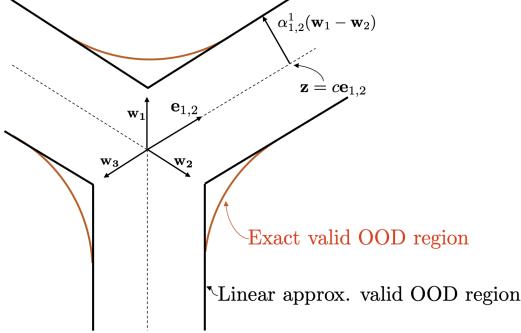


Figure 9: Geometry of the linear approximate region for $H = 2$, $K = 3$.

*The hyperplane separating any two classes requires, $\sigma(\mathbf{z})_i = \sigma(\mathbf{z})_j$. For this to be true, we require $\mathbf{w}_i^\mathsf{T}\mathbf{z} = \mathbf{w}_j^\mathsf{T}\mathbf{z} \implies (\mathbf{w}_i - \mathbf{w}_j)^\mathsf{T}\mathbf{z} = 0$. This produces the hyperplane,*

$$\mathcal{R}_{H,ij} : \{\mathbf{z} \in \mathbb{R}^H | (\mathbf{w}_i - \mathbf{w}_j)^\mathsf{T}\mathbf{z} = 0\}. \tag{25}$$

*This hyperplane will only form part of the decision boundary where either of the classes $i, j$ are the maximum of all classes. We define another region where class $i$ is maximum,*

$$\mathcal{R}_{\hat{\mathbf{y}}=i} : \{\mathbf{z} \in \mathbb{R}^H | i = argmax\, \sigma(\mathbf{z})_i\}, \tag{26}$$

*and we note for $\mathcal{R}_{H,ij}$ we will be interested in the union of the region for both class $i$ and $j$, $\mathcal{R}_{\hat{\mathbf{y}}=i} \cup \mathcal{R}_{\hat{\mathbf{y}}=j}$.*

*As per the description, we are interested in a pair of hyperplanes offset in opposite directions, parallel to $\mathcal{R}_{H,ij}$. Similar to the case of $K = 2$, this results in a region in between the hyperplane pair given by,*

$$\mathcal{R}_{Hvol,ij} : \{\mathbf{z} \in \mathbb{R}^H | (\mathbf{w}_i - \mathbf{w}_j)^\mathsf{T}(\mathbf{z}_0 - \alpha_{i,j}^i(\mathbf{w}_i - \mathbf{w}_j)) < (\mathbf{w}_i - \mathbf{w}_j)^\mathsf{T}\mathbf{z} < (\mathbf{w}_i - \mathbf{w}_j)^\mathsf{T}(\mathbf{z}_0 + \alpha_{i,j}^j(\mathbf{w}_i - \mathbf{w}_j))\}. \tag{27}$$

*The overall linear approximation of the valid OOD region can be built using a combination of these volumes as follows,*

$$\mathcal{R}_{linear\,approx.} := \bigcup_{\forall i \neq j} (\mathcal{R}_{\hat{\mathbf{y}}=i} \cup \mathcal{R}_{\hat{\mathbf{y}}=j}) \cap \mathcal{R}_{Hvol,ij} \tag{28}$$

*It remains to select $\alpha_{i,j}^i$ and $\alpha_{i,j}^j$, to set the distance separating each hyperplane pair. We choose these $\alpha$'s to match the width of the exact valid OOD region at a specific point as follows. Let $\mathbf{e}_{i,j}$ be the vector following the decision boundary between classes $i$ and $j$, in the plane containing the two weight vectors of interest ($\mathbf{e}_{i,j}$ is perpendicular to $(\mathbf{w}_1 - \mathbf{w}_2)$). We choose the point on $\mathbf{e}_{i,j}$ to be far from the origin, so $\mathbf{z} = c\mathbf{e}_{i,j}$, for, $c \to \infty$.*

*Note that the exact valid OOD region will follow contours of constant softmax confidence. For general (non-optimal) weight structures it will **not** be true that, $\sigma(\mathbf{e}_{i,j} + \alpha_{i,j}^i(\mathbf{w}_1 - \mathbf{w}_2))_i = \sigma(\mathbf{e}_{i,j} - \alpha_{i,j}^j(\mathbf{w}_1 - \mathbf{w}_2))_j$, if, $\alpha_{i,j}^i = \alpha_{i,j}^j$, hence specifying the linear approximate region requires determining up to $K(K-1)$ $\alpha$ parameters (each class forms a decision boundary with every other class). However, in the special case of optimal weight structure (def. 2), due to symmetries, only a single global $\alpha$ parameter is required. We constrain, $\alpha > 0$.*

**Proposition 3.** *The linear approximate valid OOD region (def. 4) is equal to the exact valid OOD region (as per theorem 1) for $K = 2$ if $\mathbf{w}_1 = -\mathbf{w}_2$.*

*Proof.* Observe that in eq. 28, if we have $K = 2$, there is only one combination of $i, j$, and $\mathcal{R}_{\hat{\mathbf{y}}=1} \cup \mathcal{R}_{\hat{\mathbf{y}}=2}$ covers the whole space $\mathbb{R}^H$. Hence,

$$\mathcal{R}_{\text{linear approx.}} = \mathcal{R}_{H\text{vol},1,2}. \tag{29}$$

It's straightforward to see that eq. 27 is a scaled version of eq. 17 since by assumption $\mathbf{w}_1 = -\mathbf{w}_2$. (Recall scaling does not alter a hyperplane.) Also due to symmetry, $\alpha_{i,j}^i = \alpha_{i,j}^j$.

$\square$

**Theorem 2.** *The linear approximate valid OOD region is an increasingly close approximation of the $U_{max}$ exact valid OOD region at large magnitudes, for any $K \geq 3$ and general weight structures, for all points except close to the intersection between decision boundary hyperplanes.*

*Proof.* We begin by referring back to fig. 9, which illustrates intuitively what we will show. While the exact valid OOD region is curved at the intersection between decision boundary hyperplanes (both close to the origin and where two decision boundary hyperplanes intersect), far away from these intersections, the exact valid OOD region does become linear.

Consider some $\mathbf{z}$ vector that moves along the decision hyperplane between classes 1 and 2 in any direction, $\beta \mathbf{d}_{1,2}$, and perpendicular to the hyperplane, $\alpha(\mathbf{w}_1 - \mathbf{w}_2)$,

$$\mathbf{z} = \beta \mathbf{d}_{1,2} + \alpha(\mathbf{w}_1 - \mathbf{w}_2), \tag{30}$$

illustrated in fig. 10. Note that, $(\mathbf{w}_1 - \mathbf{w}_2)^\mathsf{T} \mathbf{d}_{1,2} = 0$, and that $\mathbf{e}_{1,2}$ (from def. 4) is a special case of $\mathbf{d}_{1,2}$, where it lives on the plane of $\mathbf{w}_1$ and $\mathbf{w}_2$.
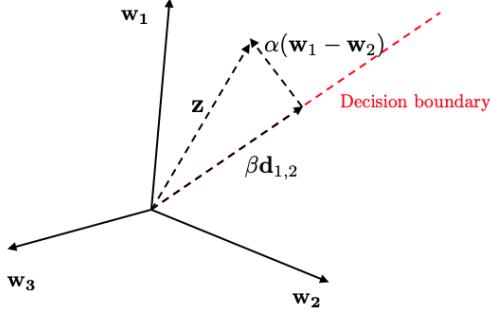


Figure 10: Vector geometry for $H = 2$, $K = 3$, and general (non-optimal) weight structuring.

Note that by specifying $\mathbf{d}_{1,2}$, $\beta$ and $\alpha$, this $\mathbf{z}$ vector can span the entire $\mathbb{R}^H$ space (though we are only interested in its use over the range $\mathcal{R}_{\hat{\mathbf{y}}=1} \cup \mathcal{R}_{\hat{\mathbf{y}}=2}$). We now push this vector through the softmax and simplify, presuming we are interested in $\beta > 0$, so the max class is 1.

$$\sigma(\mathbf{z})_1 = \frac{\exp^{\mathbf{w}_1^\mathsf{T}(\beta \mathbf{d}_{1,2} + \alpha(\mathbf{w}_1 - \mathbf{w}_2))}}{\exp^{\mathbf{w}_1^\mathsf{T}(\beta \mathbf{d}_{1,2} + \alpha(\mathbf{w}_1 - \mathbf{w}_2))} + \exp^{\mathbf{w}_2^\mathsf{T}(\beta \mathbf{d}_{1,2} + \alpha(\mathbf{w}_1 - \mathbf{w}_2))} + \sum_{j=3}^{K} \exp^{\mathbf{w}_j^\mathsf{T}(\beta \mathbf{d}_{1,2} + \alpha(\mathbf{w}_1 - \mathbf{w}_2))}} \tag{31}$$

$$= \frac{\exp^{\alpha \mathbf{w}_1^\mathsf{T}(\mathbf{w}_1 - \mathbf{w}_2)}}{\exp^{\alpha \mathbf{w}_1^\mathsf{T}(\mathbf{w}_1 - \mathbf{w}_2)} + \exp^{\mathbf{w}_2^\mathsf{T}(\beta \mathbf{d}_{1,2} + \alpha(\mathbf{w}_1 - \mathbf{w}_2)) - \mathbf{w}_1^\mathsf{T} \beta \mathbf{d}_{1,2}} + \sum_{j=3}^{K} \exp^{\mathbf{w}_j^\mathsf{T}(\beta \mathbf{d}_{1,2} + \alpha(\mathbf{w}_1 - \mathbf{w}_2)) - \mathbf{w}_1^\mathsf{T} \beta \mathbf{d}_{1,2}}} \tag{32}$$

$$= \frac{\exp^{\alpha \mathbf{w}_1^\mathsf{T}(\mathbf{w}_1 - \mathbf{w}_2)}}{\exp^{\alpha \mathbf{w}_1^\mathsf{T}(\mathbf{w}_1 - \mathbf{w}_2)} + \exp^{\alpha \mathbf{w}_2^\mathsf{T}(\mathbf{w}_1 - \mathbf{w}_2)} + \sum_{j=3}^{K} \exp^{\beta(\mathbf{w}_j^\mathsf{T} \mathbf{d}_{1,2} - \mathbf{w}_1^\mathsf{T} \mathbf{d}_{1,2}) + \alpha \mathbf{w}_j^\mathsf{T}(\mathbf{w}_1 - \mathbf{w}_2)}}. \tag{33}$$

Since we know that, $\mathbf{w}_j^\mathsf{T}\mathbf{d}_{1,2} < \mathbf{w}_1^\mathsf{T}\mathbf{d}_{1,2}$ (the logit for class 1 must be largest), we can say that provided $\mathbf{w}_1^\mathsf{T}\mathbf{d}_{1,2} = \mathbf{w}_2^\mathsf{T}\mathbf{d}_{1,2} \neq \mathbf{w}_j^\mathsf{T}\mathbf{d}_{1,2} \; \forall j \geq 3$,

$$\lim_{\beta \to \infty} \sum_{j=3}^{K} \exp^{\beta(\mathbf{w}_j^\mathsf{T}\mathbf{d}_{1,2} - \mathbf{w}_1^\mathsf{T}\mathbf{d}_{1,2}) + \alpha\mathbf{w}_j^\mathsf{T}(\mathbf{w}_1 - \mathbf{w}_2)} = 0. \tag{34}$$

Hence, when $\beta$ grows large (at large magnitudes / at points far from the origin), eq. 33 depends on $\alpha$ in a constant way for nearly all vectors $\mathbf{d}_{1,2}$.

Our linear approximate region def. 4 is defined to match the width of the exact valid OOD region, far from the origin. Therefore, at increasing magnitudes, the linear approximate valid OOD region is an increasingly close approximation of the exact valid OOD region for most points.

However, note that the rate of this convergence is slower when $\mathbf{d}_{1,2}$ approaches meeting the criteria; $\mathbf{w}_1^\mathsf{T}\mathbf{d}_{1,2} = \mathbf{w}_2^\mathsf{T}\mathbf{d}_{1,2} = \mathbf{w}_j^\mathsf{T}\mathbf{d}_{1,2}$ for some $j \geq 3$. This convergence happens at all orientations except where $\mathbf{d}_{1,2}$ intersects with a second decision boundary hyperplane. For example, if $\mathbf{w}_1^\mathsf{T}\mathbf{d}_{1,2} = \mathbf{w}_2^\mathsf{T}\mathbf{d}_{1,2} = \mathbf{w}_3^\mathsf{T}\mathbf{d}_{1,2}$, then eq. 34 does not reduce to zero, and an additional term remains in the denominator of $\sigma(\mathbf{z})_i$, that departs from the linear plane over the rest of the space.

$\square$

**Corollary 2.1.** *This corollary is an extension of theorem 2. Consider the linear approximate valid OOD region (def. 4), $\mathcal{R}_{linear\ approx.}$, for any number of classes $K \geq 3$ and general weight structure. The linear approximate valid OOD region forms a subset of the exact valid OOD region (def. 1), $\mathcal{R}_{linear\ approx.} \subseteq \mathcal{R}_{exact}$.*

*Proof.* The linear approximate valid OOD region chooses all $\alpha$'s by assuming that $\beta = \infty$ in eq. 33. In general $\beta < \infty$ so we can say,

$$\sum_{j=3}^{K} \exp^{\beta(\mathbf{w}_j^\mathsf{T}\mathbf{d}_{1,2} - \mathbf{w}_1^\mathsf{T}\mathbf{d}_{1,2}) + \alpha\mathbf{w}_j^\mathsf{T}(\mathbf{w}_1 - \mathbf{w}_2)} > 0, \quad \forall 0 \leq \beta < \infty, \quad \forall \alpha > 0. \tag{35}$$

Hence, for some $\mathbf{z}$ with $0 \leq \beta < \infty$, $\sigma(\mathbf{z})_1$ can only be lower than the linear approximate valid OOD region, and $U_{\max}(\mathbf{z})$ can only be higher. Therefore the linear region must be a subset of the exact valid OOD region $\square$

### A.3 Optimal Structure

**Theorem 3.** *Consider $K$ equally weighted components in $\mathbb{R}^H$, where $H \geq K - 1$, each following a Gaussian distribution, $p(\mathbf{z}|\mathbf{y} = i) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$. Assume mean vectors have fixed magnitude, $||\boldsymbol{\mu}_i|| = c_2 \; \forall i$. These clusters are to be separated with a softmax with weights $\mathbf{w}_i$. In order to minimise a regularised cross-entropy loss,*

$$\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i, \mathbf{w}_i = argmin_{\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i, \mathbf{w}_i} - \mathbb{E}_{p(\mathbf{z})} \left[ \sum_i^K p(\mathbf{y} = i|\mathbf{z}) \log \sigma(\mathbf{z})_i \right] + \sum_i^K \lambda_1 ||\mathbf{w}_i||^2, \quad (36)$$

*the parameters $\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i, \mathbf{w}_i$, must have an optimal decision boundary structure (def. 2).*

**Remark.** *The assumption that 'mean vectors have fixed magnitude' may not be strictly necessary for this proof, but it reduces the space of possible structures we must consider, making it more straightforward.*

*The assumption can be justified in two ways. 1) Empirically, $||\mathbf{z}||$ does not vary greatly within the training distribution with no obvious signs of multi-modal behaviour in fig. 15. 2) Earlier in the network, weight initialisations, batchnormalisation and regularisation might be expected to encourage activations for different classes to be of similar magnitudes.*

*Proof.* In order to align with the optimal structure, the following criteria must be fulfilled.

- Criteria 1. Weights are of constant magnitude, $||\mathbf{w}_i|| = c_1 \; \forall i$, and as large as possible subject to regularisation.

- Criteria 2. Bias value are zero.

- Criteria 3a. Weight vectors are evenly distributed.

- Criteria 3b. $\cos \theta_{i,j} = \frac{-1}{K-1} \; \forall i \neq j$.

- Criteria 4. Weights point toward mean vectors, $\boldsymbol{\mu}_i = c_3 \mathbf{w}_i$.

- Criteria 5. Variance of clusters is small, $\boldsymbol{\Sigma}_i \approx \mathbf{0}$.

First observe that, since we have full control over placement of class clusters, we should choose them to be linearly separable (def. 3), which simplifies our objective,

$$\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i, \mathbf{w}_i = argmin_{\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i, \mathbf{w}_i} - \sum_i^K \mathbb{E}_{p(\mathbf{z}|y=i)} \left[ \log \sigma(\mathbf{z})_i \right] + \lambda_1 ||\mathbf{w}_i||^2. \quad (37)$$

We now consider the effect of the covariance matrix on our objective. We use an analytical approximation for passing a normally distributed random vector through a softmax [Lu et al., 2021],

$$\int_{\mathbf{z}} \sigma(\mathbf{z})_i \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) d\mathbf{z} \approx \frac{1}{1 + \sum_{j \neq i} \exp - \frac{\mathbf{w}_j^\mathsf{T} \mathbf{z} - \mathbf{w}_i^\mathsf{T} \mathbf{z}}{\sqrt{1 + \frac{\pi}{8}(\mathbf{w}_i^\mathsf{T} \boldsymbol{\Sigma}_i \mathbf{w}_i + \mathbf{w}_j^\mathsf{T} \boldsymbol{\Sigma}_i \mathbf{w}_j)}}}. \quad (38)$$

From this we can see that softmax confidence always decreases with a increasing variance, and hence is maximised when, $\boldsymbol{\Sigma}_i \approx \mathbf{0}$ – this proves **criteria 5**. Class clusters are therefore delta functions, and we can rewrite our objective,

$$\max_{\boldsymbol{\mu}_i, ||\mathbf{w}_i||, \cos \theta_{\boldsymbol{\mu}_i, i} \forall i} \sum_i \log \left( \frac{\exp ||\mathbf{w}_i|| c_2 \cos \theta_{\boldsymbol{\mu}_i, i}}{\sum_j \exp ||\mathbf{w}_j|| c_2 \cos \theta_{\boldsymbol{\mu}_i, j}} \right) - \lambda_1 ||\mathbf{w}_i||^2. \quad (39)$$

Since $||\boldsymbol{\mu}_i||$ is fixed, there is no option to 'stack' boundaries (fig. 5), and so $\boldsymbol{\mu}_i$'s must be arranged on a hypersphere. There creates inherent symmetry in our problem, and weights should be of constant magnitude. $\sigma(\mathbf{z})_i$ is increasing in $||\mathbf{w}_i||$ (proposition 4), which should be as large as possible subject to the regularisation penalty (**criteria 1**).

$$\max_{\boldsymbol{\mu}_i, \cos \theta_{\boldsymbol{\mu}_i, i} \forall i} \sum_i \log \left( \frac{\exp c_1 c_2 \cos \theta_{\boldsymbol{\mu}_i, i}}{\sum_j \exp c_1 c_2 \cos \theta_{\boldsymbol{\mu}_i, j}} \right) - \lambda_1 ||\mathbf{w}_i||^2 \quad (40)$$

20

It remains to determine how these weight vectors should be arranged. One would like to maximise $\cos\theta_{\boldsymbol{\mu}_i,i}$, and minimise $\cos\theta_{\boldsymbol{\mu}_i,j}$ $\forall i \neq j$. Maximising $\cos\theta_{\boldsymbol{\mu}_i,i}$ is straightforward; by setting $\boldsymbol{\mu}_i = c_3\mathbf{w}_i$ for $c_3 > 0$ (**criteria 4**), we find $\cos\theta_{\boldsymbol{\mu}_i,i} = 1$.

To minimise $\cos\theta_{\boldsymbol{\mu}_i,j}$ $\forall i \neq j$, we turn to an area know as spherical codes [Coxeter et al., 1989, Brauchart and Grabner, 2015, Musin and Tarasov, 2015]. A general problem considered in this area requires arranging $K$ unit vectors in $H$ dimensions, so that the minimum angle between any pair of vectors, is maximised (see 'the kissing number problem' and 'Tammes problem'). This matches our own objective, where we'd like to minimise the cosine angle.

For general $H$ and $K$, this problem is unsolved, however when $H \geq K - 1$, there are more degrees of freedom allowed by the space than there are vectors, and the problem becomes more straightforward. Note that in neural networks, typically $H \gg K$.

The best thing one can do if $H \geq K - 1$, is to evenly distribute weight vectors [Bagchi, 1997] (**criteria 3a**), this implies, $\sum_i \mathbf{w}_i = \mathbf{0}$, i.e. directly summing all vectors returns to the origin, and also, $\cos\theta_{i,j} = \text{const.}$ $\forall i \neq j$, i.e. the angle between all pairs of vectors is equal.

For any vector, $\mathbf{a} \in \mathbb{R}^H$, we can therefore say,

$$\mathbf{a}^\mathsf{T} \sum_i \mathbf{w}_i = \sum_i \mathbf{w}_i^\mathsf{T} \mathbf{a} = 0. \tag{41}$$

This allows us to deduce the result in **criteria 3b**,

$$0 = \sum_i \sum_j \mathbf{w}_i^\mathsf{T} \mathbf{w}_j \tag{42}$$

$$= \sum_j \mathbf{w}_j^\mathsf{T} \mathbf{w}_j + \sum_{i \neq j} \mathbf{w}_i^\mathsf{T} \mathbf{w}_j \tag{43}$$

$$= K c_1^2 + \sum_{i \neq j} c_1^2 \cos\theta_{i,j} \tag{44}$$

$$= K + 2 \sum_{i > j} \cos\theta_{i,j} \tag{45}$$

$$= K + 2 \binom{K}{2} \cos\theta_{i,j} \tag{46}$$

$$= K + 2 \frac{K!}{2!(K-2)!} \cos\theta_{i,j} \tag{47}$$

$$= K + K(K-1) \cos\theta_{i,j} \tag{48}$$

$$\cos\theta_{i,j} = \frac{-1}{K-1}. \tag{49}$$

Finally we think of the bias, absorbed into $\mathbf{w}_i$, as having a corresponding final-layer feature of a constant 1 across classes. Since this can not be modified as part of $\boldsymbol{\mu}_i$, it only introduces a penalty via the regularisation term, $\lambda_1 ||\mathbf{w}_i||^2$. Hence it should be set to zero (**criteria 2**).

$\square$

**Proposition 4.** *Consider a softmax with weight vectors, $\mathbf{w}_i$, and $K$ classes. For any weight structure, $U_{max}(\mathbf{z})$ increases with decreasing $||\mathbf{z}||$.*

*Proof.* We have,

$$\sigma(\mathbf{z})_i = \frac{\exp \mathbf{w}_i^\mathsf{T} \mathbf{z}}{\sum_j \exp \mathbf{w}_j^\mathsf{T} \mathbf{z}} \tag{50}$$

$$= \frac{1}{1 + \sum_{j \neq i} \exp \mathbf{w}_j^\mathsf{T} \mathbf{z} - \mathbf{w}_i^\mathsf{T} \mathbf{z}} \tag{51}$$

$$= \frac{1}{1 + \sum_{j \neq i} \exp ||\mathbf{w}_j|| \, ||\mathbf{z}|| \cos \theta_{\mathbf{z},j} - ||\mathbf{w}_i|| \, ||\mathbf{z}|| \cos \theta_{\mathbf{z},i}} \tag{52}$$

$$= \frac{1}{1 + \sum_{j \neq i} \exp ||\mathbf{z}|| \left( ||\mathbf{w}_j|| \cos \theta_{\mathbf{z},j} - ||\mathbf{w}_i|| \cos \theta_{\mathbf{z},i} \right)} \tag{53}$$

$$\tag{54}$$

which clearly decreases with decreasing $||\mathbf{z}||$. Hence, $U_{\max}(\mathbf{z}) = \max_i -\sigma(\mathbf{z})_i$ increases with decreasing $||\mathbf{z}||$.

$\square$

**Remark.** *Ahead of presenting theorem 4 & 5, we first provide a simple result that intuits why we'd expect their result to be true for any number of classes, $K \geq 2$. We then explain why the proof is more involved than might initially be expected, and outline our approach.*

*We can say that the sum of the changes in cosine angles must be zero,*

$$\sum_{i=1}^{K} \triangle \cos \theta_{\mathbf{z},i} = 0. \tag{55}$$

*Decreasing $\max_i \cos \theta_{\mathbf{z},i}$ requires setting $\triangle \cos \theta_{\mathbf{z},1} < 0$. This means the cumulative change in all other cosine angles must be positive, $- \triangle \cos \theta_{\mathbf{z},1} = \sum_{j=2}^{K} \triangle \cos \theta_{\mathbf{z},j} > 0$. Writing this out in a form similar to the softmax but without exponentials clearly decreases with decreasing $\triangle \cos \theta_{\mathbf{z},1}$,*

$$\frac{\triangle \cos \theta_{\mathbf{z},1}}{\triangle \cos \theta_{\mathbf{z},1} + \sum_{j=2}^{K} \triangle \cos \theta_{\mathbf{z},j}}. \tag{56}$$

*Unfortunately, for general $a_i$'s and $b_i$'s,*

$$\sum a_i > \sum b_i \nRightarrow \sum \exp^{a_i} > \exp^{b_i}. \tag{57}$$

*Hence, introducing the exponentials to create a softmax makes the proof challenging, and we must prove there are properties of the $\cos \theta_{\mathbf{z},i}$'s that make this true. We have taken two approaches:*

*1) **Theorem 4.** For low numbers of classes, it is possible to keep track of how angles between $\mathbf{z}$ and every $\mathbf{w}_i$ change with a change in $\mathbf{z}$, and to enumerate all possible rotations of a vector. However, this approach becomes cumbersome beyond $K = 3$.*

*2) **Theorem 5.** On the other hand, we have, $\cos \theta_{i,j} = \frac{-1}{K-1}$, and for large numbers of classes, $\cos \theta_{i,j} \approx 0$, i.e. vectors are approximately mutually orthogonal to each other. This simplifies our proof, since we no longer need to keep track of angles between $\mathbf{z}$ and every $\mathbf{w}_i$, only those which it is moving away/towards from (the rest can be assumed to remain at $\cos \theta_{i,j} \approx 0$).*

**Theorem 4.** *Consider a softmax with weight vectors, $\mathbf{w}_i$, and $K = 2$ or $K = 3$ classes. For the optimal structure (def. 2), $U_{max}(\mathbf{z})$ increases with decreasing $\max_i \cos \theta_{\mathbf{z},i}$ (but fixed $||\mathbf{z}||$).*

*Proof.* We first observe two useful equations.

Similar to eq. 41, for general classes,

$$0 = \sum_i \mathbf{w}_i^\intercal \mathbf{z} \tag{58}$$

$$= \sum_i ||\mathbf{w}_i|| \, ||\mathbf{z}|| \cos \theta_{\mathbf{z},i} \tag{59}$$

$$= \sum_i \cos \theta_{\mathbf{z},i}. \tag{60}$$

It will simplify notation to write the softmax with weight vectors relabelled, so $\mathbf{w}_1$ now represents the weight vector closest (largest cosine similarity) to $\mathbf{z}$, $\cos \theta_{\mathbf{z},1} := \max_i \cos \theta_{\mathbf{z},i}$, $\mathbf{w}_2$ the second closest and so on. We also have,

$$\max_i \sigma(\mathbf{z})_i = \frac{\exp ||\mathbf{w}_1|| \, ||\mathbf{z}|| \cos \theta_{\mathbf{z},1}}{\sum_j \exp ||\mathbf{w}_j|| \, ||\mathbf{z}|| \cos \theta_{\mathbf{z},j}} \tag{61}$$

$$= \frac{\exp c \cos \theta_{\mathbf{z},1}}{\exp c \cos \theta_{\mathbf{z},1} + \sum_{j=2}^K \exp c \cos \theta_{\mathbf{z},j}}. \tag{62}$$

for some constant, $||\mathbf{w}_i|| \, ||\mathbf{z}|| = c > 0$. The specific value is unimportant here, it's effect being evaluated independently in lemma 4.

**Two classes**

For two classes we have, $\cos \theta_{\mathbf{z},1} + \cos \theta_{\mathbf{z},2} = 0$.

Any decrease in $\cos \theta_{\mathbf{z},1}$ therefore leads to an increase in $\cos \theta_{\mathbf{z},2}$. Inspecting the softmax,

$$-U_{\max}(\mathbf{z}) = \max_i \sigma(\mathbf{z})_i = \frac{\exp c \cos \theta_{\mathbf{z},1}}{\exp c \cos \theta_{\mathbf{z},1} + \exp c \cos \theta_{\mathbf{z},2}} \tag{63}$$

$$= \frac{1}{1 + \exp c(\cos \theta_{\mathbf{z},2} - \cos \theta_{\mathbf{z},1})} \tag{64}$$

$$= \frac{1}{1 + \exp -2c \cos \theta_{\mathbf{z},1}} \tag{65}$$

reveals that $\max_i \sigma(\mathbf{z})_i$ decreases with decreasing $\cos \theta_{\mathbf{z},1}$, so $U_{\max}(\mathbf{z})$ increases.

**Three classes**

For three classes, $\cos \theta_{\mathbf{z},1} + \cos \theta_{\mathbf{z},2} + \cos \theta_{\mathbf{z},3} = 0$. The subspace covered by the three weight vectors is actually a 2D plane (in fact the space covered by evenly distributed vectors is always a subspace, $\mathbb{R}^{K-1}$ [Bagchi, 1997] hence our requirement for $H \geq K - 1$). On this 2D plane, $\cos \theta_{i,j} = \frac{-1}{K-1} = \frac{-1}{2} \; \forall i \neq j \implies \theta_{i,j} = \frac{2}{3}\pi$.

First consider the case when $\mathbf{z}$ is also on this plane. We can capture all possible movement by considering the change of a single parameter, $\theta_{\mathbf{z},1} \in [0, \pi/3]$ (noting beyond $\pi/3$ the maximum class changes, and all other possibilities are symmetries of this).

In this 2D planar case, we can write, $\theta_{\mathbf{z},2} = \frac{2}{3}\pi - \theta_{\mathbf{z},1}$ and $\theta_{\mathbf{z},3} = \frac{2}{3}\pi + \theta_{\mathbf{z},1}$. We now express the softmax in these terms,

$$-U_{\max}(\mathbf{z}) = \max_i \sigma(\mathbf{z})_i = \frac{\exp c \cos \theta_{\mathbf{z},1}}{\exp c \cos \theta_{\mathbf{z},1} + \exp c \cos \theta_{\mathbf{z},2} + \exp c \cos \theta_{\mathbf{z},3}} \tag{66}$$

$$= \frac{\exp c \cos \theta_{\mathbf{z},1}}{\exp c \cos \theta_{\mathbf{z},1} + \exp c \cos(\frac{2}{3}\pi - \theta_{\mathbf{z},1}) + \exp c \cos(\frac{2}{3}\pi + \theta_{\mathbf{z},1})} \tag{67}$$

This can be plugged into a derivative calculator, as visualised in fig. 11, to verify that, $\frac{\partial \max_i \sigma(\mathbf{z})_i}{\partial \cos \theta_{\mathbf{z},1}} < 0$ across the evaluated range, hence $\max_i \sigma(\mathbf{z})_i$ is decreasing in $\cos \theta_{\mathbf{z},1}$, or equivalently, $U_{\max}(\mathbf{z})$ is increasing in $\cos \theta_{\mathbf{z},1}$.

It remains to describe what happens when $\cos \theta_{\mathbf{z},1}$ is decreased by lifting the $\mathbf{z}$ vector away from the 2D subspace, rather than moved within the plane as above. Consider some modified vector rotated off the 2D plane by some angle $\theta_r$, which can be created by, $\mathbf{z}_r = \mathbf{z} \cos \theta_r + \mathbf{n} \sin \theta_r$ (here $\mathbf{z}$ lives on
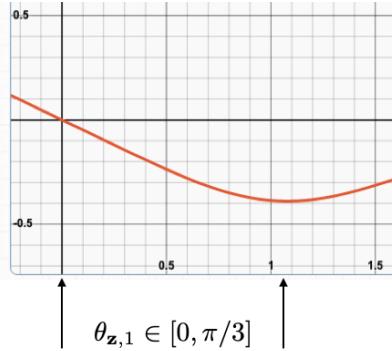
$$\theta_{\mathbf{z},1} \in [0, \pi/3]$$

Figure 11: Plot of first derivative of eq. 67, $\frac{\partial \max_i \sigma(\mathbf{z})_i}{\partial \cos \theta_{\mathbf{z},1}}$.

the 2D plane), for some vector normal to the 2D plane, $\mathbf{n}$. Now consider passing this vector through the softmax,

$$-U_{\max}(\mathbf{z}) = \max_i \sigma(\mathbf{z}_r)_i = \frac{\exp \mathbf{w}_1^\mathsf{T} \mathbf{z}_r}{\exp \mathbf{w}_1^\mathsf{T} \mathbf{z}_r + \exp \mathbf{w}_2^\mathsf{T} \mathbf{z}_r + \exp \mathbf{w}_3^\mathsf{T} \mathbf{z}_r} \tag{68}$$

$$= \frac{\exp \mathbf{w}_1^\mathsf{T} (\mathbf{z} \cos \theta_r + \mathbf{n} \sin \theta_r)}{\exp \mathbf{w}_1^\mathsf{T} (\mathbf{z} \cos \theta_r + \mathbf{n} \sin \theta_r) + \exp \mathbf{w}_2^\mathsf{T} (\mathbf{z} \cos \theta_r + \mathbf{n} \sin \theta_r) + \cdots} \tag{69}$$

$$= \frac{\exp \mathbf{w}_1^\mathsf{T} (\mathbf{z} \cos \theta_r)}{\exp \mathbf{w}_1^\mathsf{T} (\mathbf{z} \cos \theta_r) + \exp \mathbf{w}_2^\mathsf{T} (\mathbf{z} \cos \theta_r) + \exp \mathbf{w}_3^\mathsf{T} (\mathbf{z} \cos \theta_r)} \tag{70}$$

where the last line follows since $\mathbf{n}$ is normal to *all* weight vectors, so $\mathbf{w}_i^\mathsf{T} \mathbf{n} = 0$. Here $\cos \theta_r$ has a similar cooling effect to lowering the magnitude as described in lemma 4. Hence, decreasing $\cos \theta_{\mathbf{z},1}$ by rotating it off the plane also reduces softmax confidence, and hence increases $U_{\max}(\mathbf{z})$.

$\square$

**Theorem 5.** *Consider a softmax with weight vectors, $\mathbf{w}_i$, and some large number of $K$ classes, so it holds that $\cos \theta_{i,j} = \frac{-1}{K-1} \approx 0$. For the optimal structure (def. 2), $U_{max}(\mathbf{z})$ increases with decreasing $\max_i \cos \theta_{\mathbf{z},i}$.*

*Proof.* We borrow notation and several results from theorem 4; $\cos \theta_{\mathbf{z},1}$ will refer to the largest cosine angle between any weight vector $\mathbf{w}_i$ and $\mathbf{z}$. We also use the conservation of $\cos \theta_{\mathbf{z},i}$'s,

$$\sum_{i=1}^K \cos \theta_{\mathbf{z},i} = 0. \tag{71}$$

We will consider the change in $U_{\max}$ for two vectors which are modified in specific ways to decrease $\max \cos \theta_{\mathbf{z},1}$. For $\mathbf{z}_{\text{orig.}}$ and $\mathbf{z}_{\text{mod.}}$, we study the change in uncertainty estimator, $U_{\max}(\mathbf{z}_{\text{orig.}}) - U_{\max}(\mathbf{z}_{\text{mod.}})$. Where obvious from context we drop the orig. and mod. subscript. Note that the changes in $\cos \theta_{\mathbf{z},i}$'s are also conserved,

$$\sum_{i=1}^K \cos \theta_{\mathbf{z}_{\text{orig.}},i} - \cos \theta_{\mathbf{z}_{\text{mod.}},i} = \sum_{i=1}^K \triangle \cos \theta_{\mathbf{z},i} = 0. \tag{72}$$

To lighten notation, we assume (without loss of generality) that, $||\mathbf{w}_i|| = ||\mathbf{z}_{\text{orig.}}|| = ||\mathbf{z}_{\text{mod.}}|| = 1$ (this can be achieved with a scaling of the space).

24

Let us begin by writing the softmax for the original and modified vector, with $\delta_i := \triangle \cos \theta_{\mathbf{z},i}$,

$$-U_{\max}(\mathbf{z}_{\text{orig.}}) = \max_i \sigma(\mathbf{z}_{\text{orig.}})_i = \frac{\exp^{\cos \theta_{\mathbf{z},1}}}{\exp^{\cos \theta_{\mathbf{z},1}} + \sum_{j=2}^K \exp^{\cos \theta_{\mathbf{z},j}}} \tag{73}$$

$$-U_{\max}(\mathbf{z}_{\text{mod.}}) = \max_i \sigma(\mathbf{z}_{\text{mod.}})_i = \frac{\exp^{\cos \theta_{\mathbf{z},1}+\delta_1}}{\exp^{\cos \theta_{\mathbf{z},1}+\delta_1} + \sum_{j=2}^K \exp^{\cos \theta_{\mathbf{z},j}+\delta_j}} \tag{74}$$

$$= \frac{\exp^{\cos \theta_{\mathbf{z},1}}}{\exp^{\cos \theta_{\mathbf{z},1}} + \sum_{j=2}^K \exp^{\cos \theta_{\mathbf{z},j}+\delta_j-\delta_1}}. \tag{75}$$

For an increase in $U_{\max}$, or a decrease in $\max_i \sigma$, it must hold that,

$$\sum_{j=2}^K \exp^{\cos \theta_{\mathbf{z},j}} < \sum_{j=2}^K \exp^{\cos \theta_{\mathbf{z},j}+\delta_j-\delta_1}. \tag{76}$$

**Moving between two weight vectors**

We first consider a rotation such that $\cos \theta_{\mathbf{z},1}$ decreases by moving directly toward $\cos \theta_{\mathbf{z},2}$. Hence, $\delta_1 < 0$ and $\delta_2 > 0$. Because rotation is only in this 2D plane, $\delta_j$ will all be equal for $j \geq 3$, say, $\delta_c$, and must conserve the change in $\cos \theta_{\mathbf{z},i}$.

$$\delta_1 + \delta_2 + \sum_{j=3}^K \delta_j = 0 \tag{77}$$

$$\delta_1 + \delta_2 + (K-2)\delta_c = 0 \tag{78}$$

$$\frac{-(\delta_1 + \delta_2)}{K-2} = \delta_c \tag{79}$$

We turn to showing that eq. 76 is true for our case. Note that since we have assumed a large number of classes, $\cos \theta_{\mathbf{z},j} \approx 0 \ \ \forall j \geq 3$.

$$\exp^{\cos \theta_{\mathbf{z},2}} + \sum_{j=3}^K \exp^{\cos \theta_{\mathbf{z},j}} < \exp^{\cos \theta_{\mathbf{z},2}+\delta_2-\delta_1} + \sum_{j=3}^K \exp^{\cos \theta_{\mathbf{z},j}+\delta_j-\delta_1} \tag{80}$$

$$\exp^{\cos \theta_{\mathbf{z},2}} + (K-2) < \exp^{\cos \theta_{\mathbf{z},2}+\delta_2-\delta_1} + (K-2)\exp^{\delta_c-\delta_1} \tag{81}$$

$$(K-2)(1 - \exp^{\delta_c-\delta_1}) < \exp^{\cos \theta_{\mathbf{z},2}+\delta_2-\delta_1} - \exp^{\cos \theta_{\mathbf{z},2}}. \tag{82}$$

It's a little difficult to directly see why this should be true. Formally it follows from the exponential function being strictly increasing and convex. The quantity, $1 - \exp^{\delta_c-\delta_1}$, will be small but it's multiplied by $K-2$, while, $\exp^{\cos \theta_{\mathbf{z},2}+\delta_2-\delta_1} - \exp^{\cos \theta_{\mathbf{z},2}}$, will be large. Fig. 12 articulates the quantities geometrically. Roughly speaking, each side of eq. 82 has a similar total budget in terms of change in $\cos \theta$, but the budget used on the RHS is applied higher up the exponential curve, so is amplified (since it is strictly increasing and convex), resulting in a greater change than on the LHS. Hence, the inequality must be true.

**Moving between multiple weight vectors**

Having shown $U_{\max}$ always increases with a decrease in $\cos \theta_{\mathbf{z},1}$ when moving on the 2D space between two vectors, we now repeat a similar exercise, but allow movement in a 3D space between three vectors, so while $\delta_1 < 0$, both $\delta_2 > 0$ and $\delta_3 > 0$.

$$\exp^{\cos \theta_{\mathbf{z},2}} + \exp^{\cos \theta_{\mathbf{z},3}} + \sum_{j=4}^K \exp^{\cos \theta_{\mathbf{z},j}} < \exp^{\cos \theta_{\mathbf{z},2}+\delta_2-\delta_1} + \exp^{\cos \theta_{\mathbf{z},3}+\delta_3-\delta_1} + \sum_{j=4}^K \exp^{\cos \theta_{\mathbf{z},j}+\epsilon_j-\delta_1} \tag{83}$$

$$\exp^{\cos \theta_{\mathbf{z},2}} + \exp^{\cos \theta_{\mathbf{z},3}} + (K-3) < \exp^{\cos \theta_{\mathbf{z},2}+\delta_2-\delta_1} + \exp^{\cos \theta_{\mathbf{z},3}+\delta_3-\delta_1} + (K-3)\exp^{\delta_c-\delta_1} \tag{84}$$

$$(K-3)(1 - \exp^{\delta_c-\delta_1}) < (\exp^{\cos \theta_{\mathbf{z},2}+\delta_2-\delta_1} - \exp^{\cos \theta_{\mathbf{z},2}}) + (\exp^{\cos \theta_{\mathbf{z},3}+\delta_3-\delta_1} - \exp^{\cos \theta_{\mathbf{z},3}}) \tag{85}$$
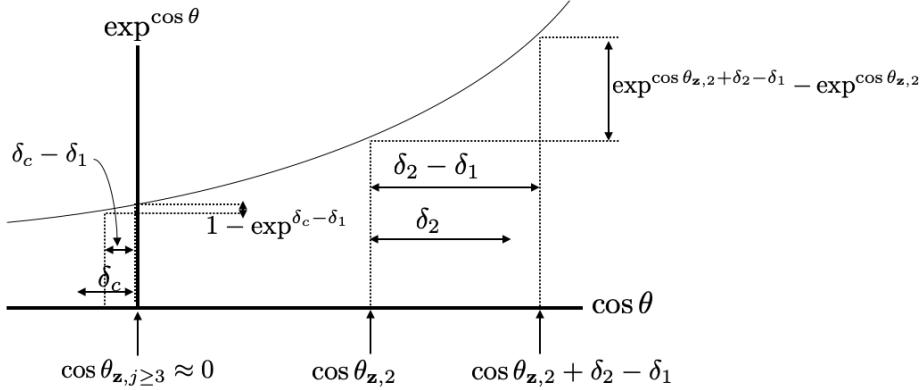
Figure 12: Geometry of eq. 82.

Precisely the same arguments apply in this case as before; since $\cos\theta_{\mathbf{z},2} > \cos\theta_{\mathbf{z},j} \approx 0$ and $\cos\theta_{\mathbf{z},3} > \cos\theta_{\mathbf{z},j} \approx 0$, the exponential amplifies the budget of change in cosine angle on the RHS of the equation, making it true.

Clearly this reasoning can be extended to further weight vectors.

**Moving outside of the subspace covered by all vectors**

The other way to decrease $\max_i \cos\theta_{\mathbf{z},i}$ is to rotate a vector outside of the $(K-1$ dimensional) subspace covered by all weight vectors. It should be straightforward to recognise that this decreases the cosine angle of all weight vectors, increasing $U_{\max}$ – we discuss this in more detail in theorem 4 for $K = 3$.

$\square$

**Corollary 5.1.** *Generalising from theorem 4 & 5, assume that for the optimal structure, for any $K \geq 2$, it holds that $U_{max}$ always increases for a decrease in $\max_i \cos\theta_{\mathbf{z},i}$. Further assume fixed cluster magnitudes, $||\boldsymbol{\mu}_i|| = c_2 \ \forall i$. The optimal structure (def. 2) maximises the linear approximate valid OOD region for $U_{max}$.*

*Proof.* This result can be seen by observing fig. 3. First note that an increase in cluster variance will decrease the valid OOD region. Since the optimal structure assumes, $\boldsymbol{\Sigma} \approx \mathbf{0}$, the valid OOD region is maximised from this perspective. Secondly, if $U_{\max}$ always increases for a decrease in $\max_i \cos\theta_{\mathbf{z},i}$, the valid OOD region will lie between all data clusters.

$\square$

**Proposition 5.** *$U_{max}(\mathbf{z})$ is not guaranteed to increase with decreasing $\max_i \cos\theta_{\mathbf{z},i}$ for general (non-optimal) decision boundary structures.*

*Proof.* We provide one numerical counterexample to verify. Consider a 'sandwich' weight structuring as in fig. 5, with $K = 3$ and $H = 2$. We choose,

$$\mathbf{w}_1 = [0, 1]$$
$$\mathbf{w}_2 = [-1, 0]$$
$$\mathbf{w}_3 = [1, 0]$$
$$\mathbf{z}_1 = [1, 0]$$
$$\mathbf{z}_2 = [0.9, -0.44],$$

so that $||\mathbf{z}_1|| = ||\mathbf{z}_2||$, and, $\max_i \cos\theta_{\mathbf{z}_1,i} > \max_i \cos\theta_{\mathbf{z}_2,i}$ – specifically $\max_i \cos\theta_{\mathbf{z}_1,i} = 1$, $\max_i \cos\theta_{\mathbf{z}_2,i} = 0.8$.

Passing both $\mathbf{z}$ vectors through a softmax gives, $\max\sigma(\mathbf{z}_1)_i = 0.665$ and $\max\sigma(\mathbf{z}_2)_i = 0.700$, hence $U_{\max}(\mathbf{z}_1) > U_{\max}(\mathbf{z}_2)$ despite $\mathbf{z}_1$ having a larger $\max\cos\theta$.

$\square$

## A.4 Mental Model

**Proposition 6.** *Assume a softmax layer separating $K$ classes has an optimal decision boundary structure (def. 2), and also that $\cos\theta_{i,\mathbf{z}} = \frac{-1}{K-1}$ for all angles except the maximum class. An uncertainty estimator that maintains the ordering of the resulting estimator (but not absolute values) is,*

$$U_{max\ mental}(\mathbf{z}) := \frac{-1}{1 + (K-1)\exp -||\mathbf{z}||(\frac{1}{K-1} + \max\cos\theta_{i,\mathbf{z}})}. \tag{86}$$

*Proof.* Begin by writing $U_{\max}$ in terms of magnitudes and angles. Then observe that for the optimal structure, weight magnitudes are constant, $||\mathbf{w}_i|| = c \quad \forall i$. Next substitute $\cos\theta_{i,\mathbf{z}} = \frac{-1}{K-1}$ for non-maximum class angles,

$$-U_{\max}(\mathbf{z}) = \frac{\exp ||\mathbf{w}_i||\, ||\mathbf{z}||\cos\theta_{i,\mathbf{z}}}{\sum_j \exp ||\mathbf{w}_j||\, ||\mathbf{z}||\cos\theta_{j,\mathbf{z}}} \tag{87}$$

$$= \frac{\exp c\, ||\mathbf{z}||\cos\theta_{i,\mathbf{z}}}{\sum_j \exp c\, ||\mathbf{z}||\cos\theta_{j,\mathbf{z}}} \tag{88}$$

$$= \frac{\exp c\, ||\mathbf{z}||\max\cos\theta_{i,\mathbf{z}}}{\exp c\, ||\mathbf{z}||\max\cos\theta_{i,\mathbf{z}} + \sum_{j\neq i}\exp c\, ||\mathbf{z}||\cos\theta_{j,\mathbf{z}}} \tag{89}$$

$$= \frac{\exp c\, ||\mathbf{z}||\max\cos\theta_{i,\mathbf{z}}}{\exp c\, ||\mathbf{z}||\max\cos\theta_{i,\mathbf{z}} + (K-1)\exp c\, ||\mathbf{z}||\frac{-1}{K-1}} \tag{90}$$

$$= \frac{1}{1 + (K-1)\exp\left(c\, ||\mathbf{z}||\frac{-1}{K-1} - c\, ||\mathbf{z}||\max\cos\theta_{i,\mathbf{z}}\right)} \tag{91}$$

$$= \frac{1}{1 + (K-1)\exp -c\, ||\mathbf{z}||\left(\frac{1}{K-1} + \max\cos\theta_{i,\mathbf{z}}\right)} \tag{92}$$

If we are only concerned by the relative ordering produced by our uncertainty estimator, as is the case in OOD detection, we can ignore the constant $c$, to produce our mental model,

$$U_{\max\ mental}(\mathbf{z}) := \frac{-1}{1 + (K-1)\exp -||\mathbf{z}||(\frac{1}{K-1} + \max\cos\theta_{i,\mathbf{z}})}. \tag{93}$$

$\square$

**Remark.** *We are careful to point out that the model we propose in proposition 6 is intended to formalise findings from this paper in a single interpretable equation, where we have traded off accuracy for readability. Nevertheless, we discuss the validity of the assumptions made, arguing that, whilst imperfect, they are not entirely unreasonable.*

*Assumption 1: A trained softmax layer has optimal decision boundary structure. Section 4.2 evaluated how well the optimal structure matched that found in neural networks of various architectures trained on various datasets. Whilst there were differences, we argued it was a reasonable match.*

*Assumption 2: $\cos\theta_{i,\mathbf{z}} = \frac{-1}{K-1}$ for all angles except the maximum class. Under the optimal structure, data from the training distribution should point in the same direction as their corresponding weight vector. In this case, we'd expect. $\max_i \cos\theta_{i,\mathbf{z}} \approx 1$, and for all other $i$'s, $\cos\theta_{i,\mathbf{z}} = \frac{-1}{K-1} \; \forall i \neq argmax_i \cos\theta_{i,\mathbf{z}}$. It's less clear how well this holds for OOD data. Fig. 15 (middle column) plots this empirically for trained ResNet18's, which shows similar patterns for both the OOD and training histograms, suggesting this assumption holds no worse than the optimal structure assumption.*

# B Extended Related Work

## B.1 Comparison Studies in the Literature

In table 2 we list a selection of previous studies reporting AUROC for OOD detection, both for softmax confidence, and other reference methods.

Note softmax confidence AUROC never drops below 78% (Shafaei et al. [2019] report accuracy), and achieves up to 98.8%, hence we report 75% to 99% in the main text, also in line with our own results in fig. 14. Modifications typically improve over softmax by 1-5% with two exceptions, one of which used an unconventional 'Dirty-MNIST' training dataset (testing regular MNIST vs. F.MNIST in our own experiments produces around 94% AUROC).

Table 2: AUROC scores for OOD detection reported in the literature. Various architectures, methods and datasets.

| Citation | $\mathcal{D}_{in}$ | $\mathcal{D}_{out}$ | Reported OOD score in format, Method:AUROC |
|---|---|---|---|
| **Softmax only** | | | |
| [Hendrycks and Gimpel, 2017b] | CIFAR10 | Various | Softmax:96.0 |
| [Hendrycks and Gimpel, 2017b] | CIFAR100 | Various | Softmax:90.0 |
| [Hendrycks and Gimpel, 2017b] | MNIST | Various | Softmax:91.0 |
| [Hendrycks and Gimpel, 2017b] | IMDB | Various | Softmax:94.0 |
| [Hendrycks and Gimpel, 2017b] | TIMIT | Various | Softmax:97.0 |
| **Various methods** | | | |
| [Malinin and Gales, 2018] | MNIST | Omniglot | Softmax:98.8, MC Dropout:99.2, PriorNetwork:100.0 |
| [Lee et al., 2018] | CIFAR10 | SVHN | Softmax:89.9, Mahalanobis(vanilla):93.9 |
| [Shafaei et al., 2019] | Various | Various | Softmax:72.6, MCDropout:73.5, Ensemble:72.8 |
| [Van Amersfoort et al., 2020] | CIFAR10 | SVHN | Softmax:90.6, Ensemble:93.3, DUQ:92.7 |
| [Van Amersfoort et al., 2020] | F.MNIST | MNIST | Softmax:84.3, Ensemble:86.1, DUQ:95.5 |
| [Mukhoti et al., 2021] | CIFAR10 | SVHN | Softmax:93.0, Ensemble:97.7, DUQ:93.7, DDU:97.9 |
| [Mukhoti et al., 2021] | CIFAR10 | CIFAR100 | Softmax:88.1, Ensemble:91.4, DUQ:85.9, DDU:91.8 |
| [Mukhoti et al., 2021] | 'Dirty MNIST' | F.MNIST | Softmax:84.0, Ensemble:97.7, DDU:99.8 |
| **Pre-training** | | | |
| [Hendrycks et al., 2019a] | CIFAR10 | SVHN | Softmax:91.8, Softmax(pre-train):95.7 |
| [Hendrycks et al., 2019a] | CIFAR10 | Various | Softmax:77.6, Softmax(pre-train):83.8 |

## B.2 Softmax Criticisms

Below we quote a selection of academic papers and talks on the topic of uncertainty in deep learning, evidencing a general belief that the softmax confidence of deep neural networks is unsuitable for capturing epistemic/model/distributional uncertainty.

**Academic Papers**

*"[The softmax output] is often erroneously interpreted as model confidence."*
[Gal and Ghahramani, 2015]

*"Deterministic models can capture aleatoric uncertainty but cannot capture epistemic uncertainty."*
[Gal et al., 2017]

*"NNs ... until recently have been unable to provide measures of uncertainty in their predictions."*
[Malinin and Gales, 2018]

*"When asked to predict on a data point unlike the training data, the NN should increase its uncertainty. There is no mechanism built into standard NNs to do this ... standard NNs cannot estimate epistemic uncertainty."*
[Pearce, 2021]

*"NNs are poor at quantifying predictive uncertainty."*
[Lakshminarayanan et al., 2017]

*"Deep neural networks with the softmax classifier are known to produce highly overconfident posterior distributions even for such abnormal samples."*
[Lee et al., 2018]

*"The output of the [softmax] classifier cannot identify these [far from the training data] inputs as out-of-distribution."*
[Hein et al., 2019]

*"The only uncertainty that can reliably be captured by looking at the softmax distribution is aleatoric uncertainty."*
[Van Amersfoort et al., 2020]

*"Softmax entropy is inherently inappropriate to capture epistemic uncertainty."*
[Mukhoti et al., 2021]

*"For [softmax] classifiers ... misclassification will occur with high confidence if the unknown is far from any known data."*
[Boult et al., 2019]

**Academic Talks**

*"Softmax is not telling you anything about ... model uncertainty."*
Elise Jennings, Training Program on Extreme-Scale Computing 2019
https://youtu.be/Puc_ujh5QZs?t=1323

*"The [softmax] network has no way of telling you 'I'm completely uncertain about the outcome and don't rely on my prediction'."*
Florian Wilhelm, PyData Berlin 2019
https://youtu.be/LCDIqL-8bHs?t=262

## B.3  Comprehensive Literature Review

This section provides a full literature review (summarised in section 2).

**Softmax origins:** Neural networks have long combined the cross-entropy loss with a sigmoid or softmax output for classification problems [Baum and Wilczek, 1987, Bridle, 1989, 1990]. The original attraction was compatibility with maximum likelihood and information theoretic objectives and the interpretation of outputs as posterior predictive probabilities [Richard and Lippmann, 1991].

**Softmax criticism:** Recent work has shown that interpreting softmax confidence as predictive probabilities can have several pitfalls; they are poorly calibrated and generally overconfident (i.e. 90% confidence does not mean correct 90% of the time) [Guo et al., 2017], and can be easily manipulated by adversarial examples [Nguyen et al., 2015, Ozbulak et al., 2018]. It's also been claimed there is no reason to trust them outside of the training distribution – *"[the softmax output] is often erroneously interpreted as model confidence."* [Gal and Ghahramani, 2015].

This paper investigates the last point. Whilst most work on uncertainty and neural networks criticises softmax uncertainties only informally, several works more rigorously demonstrate weaknesses. Hein et al. [2019] prove that any input can be magnified, $\tilde{\mathbf{x}} = \alpha\mathbf{x}$, with large $\alpha > 1$, to produce an arbitrarily high confidence prediction by ReLU networks. Whilst they define $\alpha\mathbf{x}$ as 'far from the training data', our work uses 'far' to mean unfamiliar image datasets with bounded pixel intensity (something trivial to check).

**More on Mukhoti et al.:** Mukhoti et al. [2021] released work concurrently with our own. Though they make claims that contradict our findings, these contradictions mostly arise from overly-strong wording and differences in task set up, and in our view the two papers are actually complimentary.

Their study is motivated by settings where there is significant overlap between classes, and hence high aleatoric uncertainty in the training data. Their experiments use a custom-created 'Ambiguous-MNIST' dataset, where a VAE generates inputs in-between MNIST classes, with multiple conflicting labels assigned to each input. 'Dirty-MNIST' refers to a mixture of Ambiguous-MNIST and standard MNIST data. This leads to justifiably low confidence softmax estimates for much of the training dataset, and hence softmax entropy is weakened when distinguishing training data from OOD data. However, even for this dataset, softmax confidence achieves a modest 84% AUROC in OOD detection.

There are two main arguments supporting their claim to *"prove we cannot infer epistemic uncertainty from a deterministic model's softmax entropy"*. The first hinges on the observation that by looking at softmax entropy, one cannot disentangle aleatoric from epistemic uncertainty. While we do not dispute this (section 6, cause 2b), according to our paper's findings, the same logic could be used to

claim the reverse; that we cannot infer aleatoric uncertainty from softmax entropy. Secondly, they argue that in an ensemble with epistemic uncertainty, some of the individual models must have lower softmax entropy than others, making them inherently unreliable. We do not dispute that there is increased variance in the softmax entropies of an ensemble for OOD data, but our findings would predict the *average* softmax entropy of individual models would also increase.

In our view Mukhoti et al. highlight an important failure of softmax conflating the two uncertainty types. The majority of benchmarks in the literature, e.g. those in table 2, use datasets for which high accuracy is achievable, and which do not have heavy aleatoric uncertainty – our analysis focuses on this set up, which has allowed us to investigate the relationship of softmax confidence and epistemic uncertainty in isolation. We do not comment on which of these set ups are more reflective to the real-world problems. In practise we note that it is simple for practitioners to know if there is high aleatoric uncertainty in their dataset by examining test and training accuracies (if high, there will be low aleatoric uncertainty).

**Softmax uses:** Two independent studies empirically assessed the ability of softmax confidence to detect OOD inputs and misclassified test inputs [Hendrycks and Gimpel, 2017a, Shafaei et al., 2019]. They found softmax confidence performs reasonably well on both tasks, though has room for improvement. This modest capability has been confirmed through work using softmax as a baseline (table 2).

**Softmax manipulation:** Several prominent works have shown that OOD inputs can be intentionally crafted to create high softmax confidence predictions, for example through gradient-based perturbations to the input [Szegedy et al., 2014], pattern generators optimised via evolutionary algorithms [Nguyen et al., 2015], or synthetic objects rendered in unusual poses [Alcorn et al., 2018]. We test a more benign setup in our work, commonly used in OOD benchmarking, where OOD data is selected arbitrarily as some unrelated dataset.

**Improvements:** Various approaches improve over uncertainties obtained directly from the softmax, e.g. ensembling [Heskes, 1996, Tibshirani, 1996, Lakshminarayanan et al., 2017] and Bayesian neural networks [MacKay, 1992]. But these always come with an increased implementation and computational burden, and have so far failed to be widely adopted [Wenzel et al., 2020]. As a result softmax confidence is the sole measure of uncertainty available in most of today's neural networks.

**Pre-training:** Zeiler et al. [2014] showed networks pre-trained on large-scale generic vision tasks can be fine-tuned on other vision tasks with improved accuracy and reduced training times. Recent work has shown this also leads to improved robustness, tested on OOD detection in vision and NLP tasks [Orhan, 2019, Hendrycks et al., 2019a, 2020].

# C  Experimental Details and Further Results

This section lists the main hyperparameters for each experiment. For full details, we refer readers to the code – for each experiment and figure we point to the relevant script. These are hosted at `https://github.com/******`.

Smaller experiments in the paper were run on a machine with a single GPU (GTX 1080), while the benchmark results in table 1 were run on an internal cluster of GPUs (four Titan X's).

## C.1  Misc. Figures

- Fig. 1a script: `toy_classification_latent_05.py`

  Architecture:
  $\mathbf{x} \in \mathbb{R}^2 \to 16\,\text{fc} \to H = 16 \to \text{softmax}\ K = 3$, ReLU activations, trained for 50 epochs

- Fig. 1b script: `bottleneck_mnist_08_intro.py`

  Architecture:
  $\mathbf{x} \in \mathbb{R}^{28,28,1} \to 32\times(5,5)\,\text{conv} \to 64\times(3,3)\,\text{conv} \to (2,2)\,\text{pooling} \to 128\times(3,3)\,\text{conv} \to (2,2)\,\text{pooling} \to H = 16\,\text{fc} \to \text{softmax}\ K = 3$, ReLU activations, trained for 2 epochs

  PCA visualisation:
  We select the first two principle components over the final-layer space (on a mixture of training and OOD data points). We then create a grid on this 2D plane, plotting $U_{\text{entropy}}$ over the grid. 100 data points from the test set, and 100 data points from the OOD dataset are plotted.

- Fig. 2 script: `region_03_2class.py`, `region_04_3class.py`

- Fig. 8 script: `vector_field_03.py`

## C.2  Decision Boundary Analysis

Relating to experiment: Section 4.2, fig. 4

Run using script: `symmetry_search_04.py`

LeNet models are trained from scratch with following architectures.

For $K = 2$:
$\mathbf{x} \in \mathbb{R}^{28,28,1} \to 32\times(5,5)\,\text{conv} \to 64\times(3,3)\,\text{conv} \to (2,2)\,\text{pooling} \to 128\times(3,3)\,\text{conv} \to (2,2)\,\text{pooling} \to H = 4\,\text{fc} \to \text{softmax}\ K = 2$, Tanh activations

For $K = 5$:
$\mathbf{x} \in \mathbb{R}^{28,28,1} \to 32\times(5,5)\,\text{conv} \to 64\times(3,3)\,\text{conv} \to (2,2)\,\text{pooling} \to 128\times(3,3)\,\text{conv} \to (2,2)\,\text{pooling} \to H = 32\,\text{fc} \to \text{softmax}\ K = 5$, Tanh activations

Trained ResNet18's (C10 and C100) are taken from experiments in section 6.

Trained EfficientNetB0 is taken from tensorflow's public checkpoint.

## C.3  Counterfactual Experiments

Relating to experiment: Section 4.3, fig. 5

Run using script: `counterfact_mnist_02.py`

Architecture:
$\mathbf{x} \in \mathbb{R}^{28,28,1} \to 32\times(5,5)\,\text{conv} \to 64\times(3,3)\,\text{conv} \to (2,2)\,\text{pooling} \to 128\times(3,3)\,\text{conv} \to (2,2)\,\text{pooling} \to H = 16\,\text{fc} \to \text{softmax}\ K = 3$, ReLU activations, trained for 2 epochs

For the stack structure, we deviated slightly from this. To avoid the optimisation failing, we needed to evolve the weight vectors during training – see code for details. All other structures were fixed from initialisation.

Fig. 13 shows PCA visualisations (similar to those described in section C.1) relating to fig. 5.
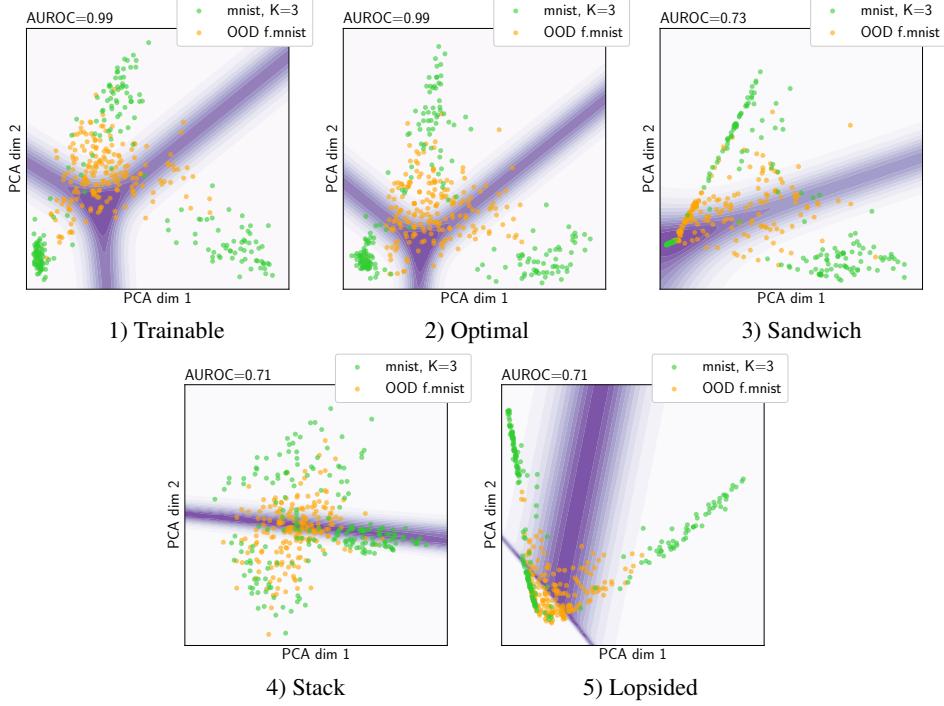
1) Trainable     2) Optimal     3) Sandwich

4) Stack     5) Lopsided

Figure 13: Relating to fig. 5, we provide PCA visualisations in the final-hidden layer for each structure tested.

## C.4 Filtering Experiments

### C.4.1 Visualising Activations

Relating to experiment: Section 5.1, fig. 6

Run using script: `bottleneck_mnist_10_interp_04.py`

Architecture:
$\mathbf{x} \in \mathbb{R}^{28,28,1} \rightarrow 32\times(5,5)$ conv $\rightarrow 64\times(3,3)$ conv $\rightarrow (2,2)$ pooling $\rightarrow 128\times(3,3)$ conv $\rightarrow (2,2)$ pooling $\rightarrow H = 64$ fc $\rightarrow$ softmax $K = 3$, ReLU activations, trained for 40 epochs

We used a simple algorithm to order the activations before displaying, grouping by which class each final-layer neurons was most activated for.

### C.4.2 Literal Filter Interpretation

Relating to experiment: Section 5.1, fig. 7

Run using script: `permute_test_04_mac_02.py`

Architecture:
$\mathbf{x} \in \mathbb{R}^{9,9} \rightarrow 16\times(3,3)$ conv $\rightarrow 32\times(3,3)$ conv $\rightarrow (2,2)$ pooling $\rightarrow 64\times(2,2)$ conv $\rightarrow$ softmax $K = 10$, ReLU activations, trained for 15 epochs

### C.4.3 Effect of Depth on Filtering

Relating to experiment: Section E.3, fig. 17

Run using script: `bottleneck_mnist_10_windows.py`

Example architecture for 8 convolutional layers (4x VGG blocks):
$\mathbf{x} \in \mathbb{R}^{28,28,1} \rightarrow 32\times(3,3)$ conv $\rightarrow 32\times(3,3)$ conv $\rightarrow (2,2)$ pooling $\rightarrow$ batchnorm $\rightarrow 64\times(3,3)$ conv $\rightarrow 64\times(3,3)$ conv $\rightarrow (2,2)$ pooling $\rightarrow$ batchnorm $\rightarrow 128\times(3,3)$ conv $\rightarrow 128\times(3,3)$ conv $\rightarrow (2,2)$

pooling $\rightarrow$ batchnorm $\rightarrow$ 128$\times$(3,3) conv $\rightarrow$ 128$\times$(3,3) conv $\rightarrow$ (2,2) pooling $\rightarrow$ batchnorm $\rightarrow$ dropout $\rightarrow$ 512 fc $\rightarrow$ softmax $K = 10$ , ReLU activations, trained for 12 epochs

## C.5   Standard OOD Detection Benchmarking

Relating to experiment: Section 6, table 1

Run using script: Models were trained for 100 epochs using script from https://github.com/huyvnphan/PyTorch_CIFAR10, which uses the official pytorch ResNet definitions https://github.com/pytorch/vision/blob/master/torchvision/models/resnet.py. We made only minor modifications to take different image datasets.

Comparison of GMM vs Softmax: `grid_run_05.py`

Results display: `grid_display_03.py`

Fig. 14 provides a breakdown of AUROC over individual datasets, the averages are given in table 1. We balance all datasets when computing AUROC.

We experimented extensively to tune the fitting of the GMM as far as possible, optimising for AUROC results using $U_{\text{density}}$. We performed random search over the first four variables below, and a full grid search over the last three.

- When appropriate, constrict one GMM component to model each class, or allow flexible fitting via expectation maximisation (EM)
- Number of final-layer neurons used (randomly selected) $\in [100, 256, 512]$ – note for ResNet18, $H = 512$
- Diagonal covariance regularisation $\in [1e-1, 1e-3, 1e-5, 1e-7]$
- Fit on $\mathbf{z}$ before or after non-linearity is applied
- Covariance type $\in [$full , diagonal , isotropic$]$
- Fit on $\mathbf{z}$ or $\log \mathbf{z}$
- Number components $\in [1, K, 10, 100]$

The final parameters providing the best fit are as follows.

- GMM fit with EM
- Number of final-layer neurons used $= 512$
- Diagonal covariance regularisation $= 1e-5$
- Fit on $\mathbf{z}$ after non-linearity is applied
- Covariance type $=$ full
- Fit on $\mathbf{z}$
- Number components $= K$

We further experimented using kernel density estimation (KDE) instead of GMM, but this gave consistently worse performance even after similar tuning efforts.

## C.6   Pre-training/Fine-tuning OOD Detection

Relating to experiment: Section 6, table 1

Run using script: `transfer_wind_12.py`, `read_results_transfer_02.py`

Details of datasets used are given in table 3.

We conducted a tuning process similar as in section C.5. The final parameters providing the best fit are as follows.

- GMM fit with EM
- Number of final-layer neurons used $= 1280$

AUROC



Softmax, $U_{\max}$ (mean=87.7)

| Train \ OOD | c100 | c10 | svhn | f.mnist | mnist |
|---|---|---|---|---|---|
| c100 | | 75.6 | 77.2 | 87.1 | 74.6 |
| c10 | 83.3 | | 87.0 | 89.0 | 82.2 |
| svhn | 91.9 | 92.3 | | 97.1 | 81.2 |
| f.mnist | 80.6 | 81.5 | 90.6 | | 96.7 |
| mnist | 95.7 | 96.7 | 99.1 | 93.8 | |

Softmax, $U_{\text{entropy}}$ (mean=88.4)

| Train \ OOD | c100 | c10 | svhn | f.mnist | mnist |
|---|---|---|---|---|---|
| c100 | | 76.3 | 79.0 | 89.4 | 77.5 |
| c10 | 83.6 | | 87.0 | 89.6 | 82.6 |
| svhn | 92.1 | 92.5 | | 97.5 | 81.4 |
| f.mnist | 81.9 | 82.8 | 92.1 | | 98.0 |
| mnist | 95.7 | 96.7 | 99.2 | 93.8 | |

Density, $U_{\text{density}}$ (mean=92.6)

| Train \ OOD | c100 | c10 | svhn | f.mnist | mnist |
|---|---|---|---|---|---|
| c100 | | 74.7 | 84.7 | 90.5 | 80.1 |
| c10 | 86.2 | | 88.2 | 92.4 | 90.5 |
| svhn | 95.7 | 95.6 | | 95.6 | 85.9 |
| f.mnist | 97.6 | 97.2 | 99.2 | | 99.8 |
| mnist | 99.5 | 99.6 | 99.7 | 99.2 | |

Softmax, $\max \cos \theta_{i,\mathbf{z}}$ (mean=89.8)

| Train \ OOD | c100 | c10 | svhn | f.mnist | mnist |
|---|---|---|---|---|---|
| c100 | | 76.2 | 80.1 | 89.4 | 77.6 |
| c10 | 83.6 | | 82.3 | 91.1 | 84.5 |
| svhn | 92.8 | 92.8 | | 95.3 | 80.6 |
| f.mnist | 90.4 | 90.6 | 94.5 | | 97.9 |
| mnist | 99.0 | 99.1 | 99.4 | 98.2 | |

Softmax, $||\mathbf{z}||$ (mean=86.3)

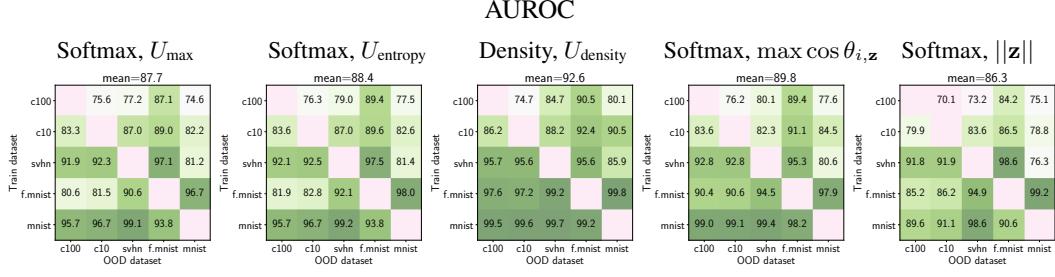| Train \ OOD | c100 | c10 | svhn | f.mnist | mnist |
|---|---|---|---|---|---|
| c100 | | 70.1 | 73.2 | 84.2 | 75.1 |
| c10 | 79.9 | | 83.6 | 86.5 | 78.8 |
| svhn | 91.8 | 91.9 | | 98.6 | 76.3 |
| f.mnist | 85.2 | 86.2 | 94.9 | | 99.2 |
| mnist | 89.6 | 91.1 | 98.6 | 90.6 | |

Figure 14: Full permutations of training datasets, with each dataset used as OOD. Numbers are AUROC for an evenly balanced combination of train and OOD datasets, averaged over five runs. E.g. for a ResNet18 trained on CIFAR 10, the mean AUROC against MNIST using $U_{\max}$ is 82.2%.

- Diagonal covariance regularisation $= 1e - 3$
- Fit on $\mathbf{z}$ after non-linearity is applied
- Covariance type = full
- Fit on $\log \mathbf{z}$
- Number components $= K$

Table 3: Transfer dataset information.

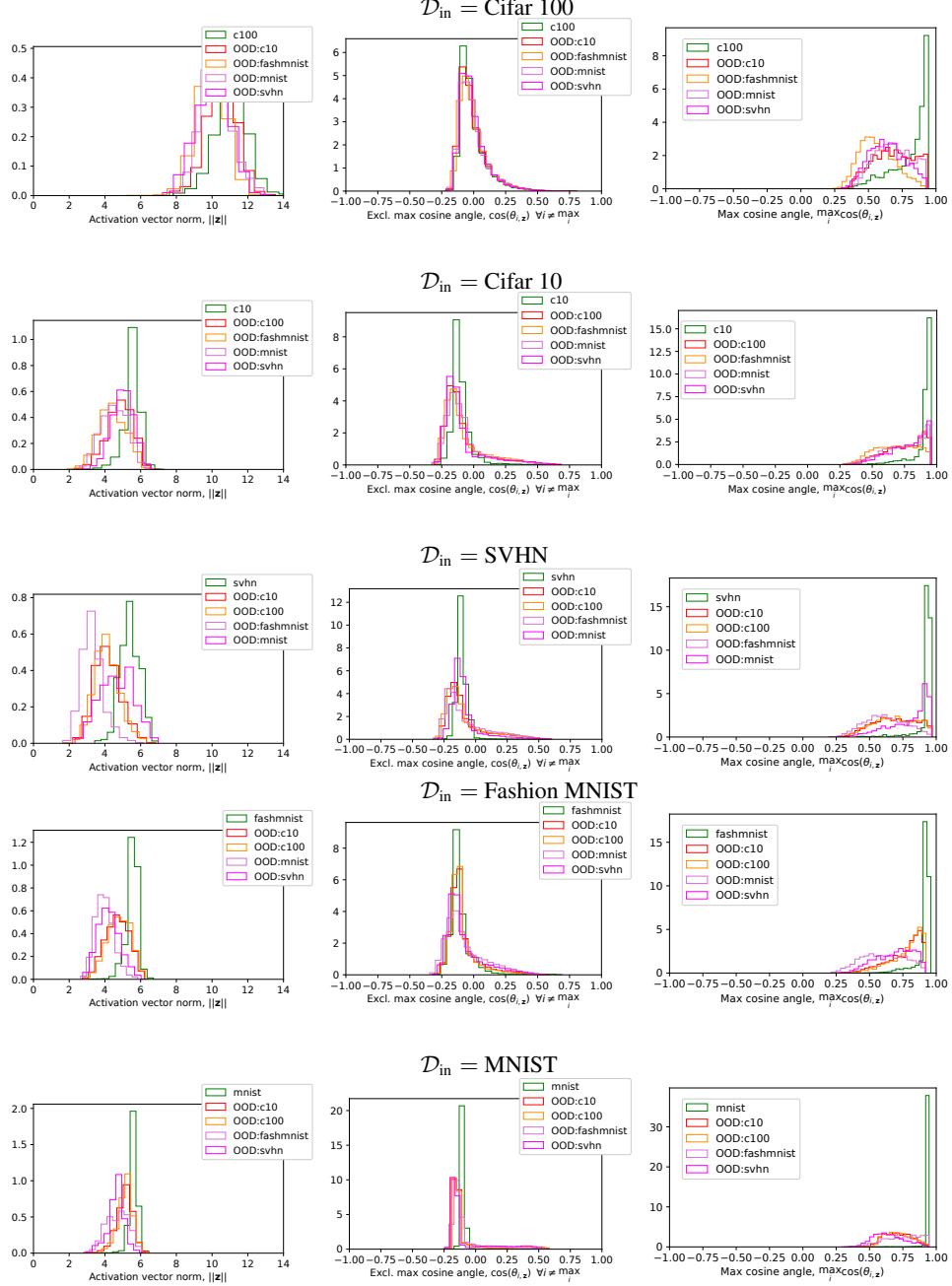| Dataset | Crop size | $N$ train data | $N$ test data | $K$ classes | Last-layer eps | Full fine-tune eps | Hyperlink |
|---|---|---|---|---|---|---|---|
| Satellite | (224, 224) | 1600 | 250 | 21 | 16 | 32 | uc merced |
| Cancer | (150, 150) | 4000 | 1000 | 8 | 16 | 32 | colorectal histology |
| Pets | (224, 224) | 3680 | 1000 | 37 | 16 | 32 | oxford iiit pet |
| Flowers | (224, 224) | 2040 | 1000 | 102 | 16 | 64 | oxford flowers102 |
| Beans | (224, 224) | 534 | 500 | 3 | 16 | 32 | beans |

# D Plots



Figure 15: For ResNet18's trained on $\mathcal{D}_{\text{in}}$ from experiments in section 6, we plot distributions of $||\mathbf{z}||$ (left column), and cosine angles – both the maximum cosine angle (right column), and all except the maximum cosine angle (middle column).

```
from tensorflow import keras
from sklearn.mixture import GaussianMixture

# train model in usual way
model.fit(x_train, y_train)

# extract train data features from final hidden layer
model_z = keras.Model(inputs=model.input,
                      outputs=model.layers[-2])
z_train = model_z.predict(x_train)

# fit density model
gmm = GaussianMixture(n_components).fit(z_train)

# extract test features
z_test = model_z.predict(x_test)

# compute log probability score
uncertainty_estimate = -gmm.score_samples(z_test)
```

# E  Miscellaneous

## E.1  Code Snippet

We detail the procedure for computing $U_{\text{density}}$ through a code snippet.

## E.2  Illustrative Example of Softmax Failure

Fig. 16 plots final-layer activations as in fig. 6, but here ankle boots are included in the subset of three Fashion MNIST training classes. MNIST digits are consistently mapped to the same feature space. Note this is not a failure of softmax extrapolations, but due to 'feature overlap'.
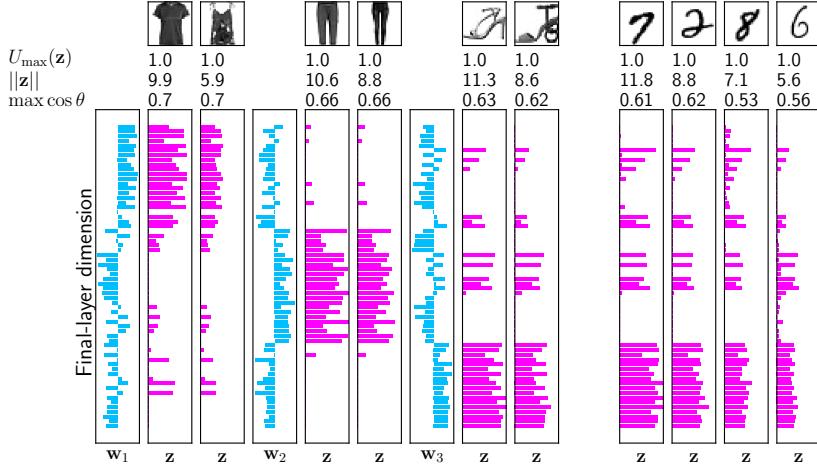


Figure 16: When ankle boots are included in the subset of classes trained on, OOD MNIST digits are consistently mapped into that class with high confidence, leading to poor OOD detection.

## E.3  The Importance of Depth on the Filtering Effect

We hypothesised that the feature filtering (section 5) effect might be stronger in deeper convolutional networks – the receptive field of convolutional filters increases with depth, meaning more holistic features are represented [Zeiler and Fergus, 2014]. OOD data may be less likely to contain these fuller features. To test this, we present a brief experiment; VGG-style networks with varying numbers of convolutional blocks are trained on all classes of either MNIST or Fashion MNIST, with the other dataset used as OOD data. Accuracy and AUROC are plotted for each dataset pair and network depth

in fig. 17. Whilst accuracy quickly plateaus in these simple tasks, AUROC signifcantly improves in the deeper networks, suggesting depth strenghthens the filtering effect.
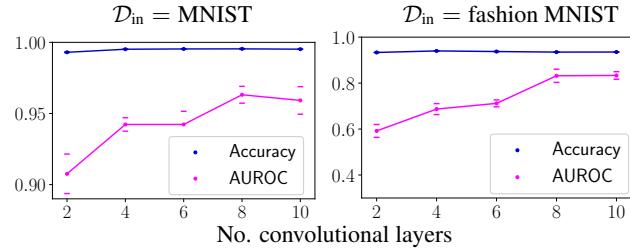


Figure 17: As depth increases, OOD detection improves. Mean±1 std. error over five runs.