

Predictive Uncertainty Estimation via **Prior Networks**

2018 NeurIPS

Xin Gao

2023.12.29

Introduction

- This work proposes a new framework for modeling predictive uncertainty called Prior Networks (PNs) which **explicitly models distributional uncertainty**.
- PNs do this by parameterizing a prior distribution over predictive distributions.

model uncertainty, data uncertainty and distributional uncertainty

- Model uncertainty, **Epistemic** uncertainty (reducible)
- Data uncertainty, **Aleatoric** uncertainty (irreducible)
- Distributional uncertainty arises due to mismatch between the training and test distributions (also called **dataset shift**)

Bayesian uncertainty

Distribution $p(\mathbf{x}, y)$ over input features \mathbf{x} and labels y

A classification model $P(\omega_c | \mathbf{x}^*, \mathcal{D})$ trained on a finite dataset $\mathcal{D} = \{\mathbf{x}_j, y_j\}_{j=1}^N \sim p(\mathbf{x}, y)$

$$P(\omega_c | \mathbf{x}^*, \mathcal{D}) = \int \underbrace{P(\omega_c | \mathbf{x}^*, \boldsymbol{\theta})}_{\text{Data}} \underbrace{p(\boldsymbol{\theta} | \mathcal{D})}_{\text{Model}} d\boldsymbol{\theta}$$

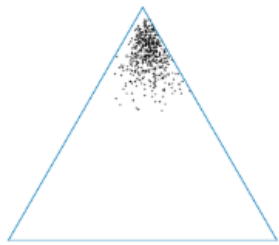
Approximation $q(\boldsymbol{\theta})$ $p(\boldsymbol{\theta} | \mathcal{D}) \approx q(\boldsymbol{\theta})$

Sampling $P(\omega_c | \mathbf{x}^*, \mathcal{D}) \approx \frac{1}{M} \sum_{i=1}^M P(\omega_c | \mathbf{x}^*, \boldsymbol{\theta}^{(i)}), \boldsymbol{\theta}^{(i)} \sim q(\boldsymbol{\theta})$

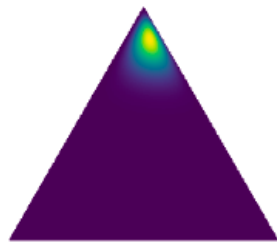
Distribution of distribution

A categorical distribution μ over class labels y

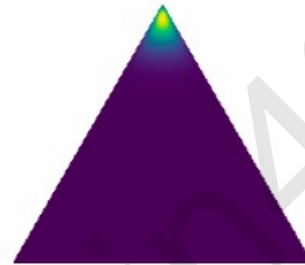
$$\left\{ P\left(\omega_c \mid \mathbf{x}^*, \boldsymbol{\theta}^{(i)}\right) \right\}_{i=1}^M$$



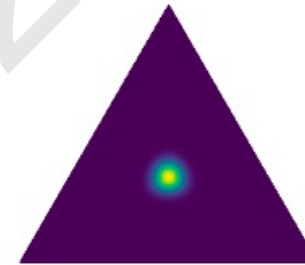
(a) Ensemble



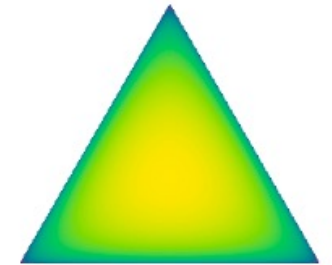
(b) Distribution



(a) Confident Prediction



(b) High data uncertainty



(c) Out-of-distribution

Figure 1: Distributions on a Simplex

Figure 2: Desired behaviors of a distribution over distributions

Explicitly parameterize a distribution over distributions on a simplex

- **Confident in-distribution data:** a sharp distribution centered on one of the corners of the simplex
- **Noise or class overlap (data uncertainty):** a sharp distribution focused on the center of the simplex
- **Confident out-of-distribution:** a flat distribution, large uncertainty

Prior Network

In Prior Networks **data uncertainty** is described by **the point-estimate categorical distribution μ** and **distributional uncertainty** is described by **the distribution over predictive categoricals $p(\mu|\mathbf{x}^*, \theta)$**

$$P(\omega_c | \mathbf{x}^*, \mathcal{D}) = \underbrace{\int \int}_{\text{Data}} \underbrace{p(\omega_c | \boldsymbol{\mu})}_{\text{Distributional}} \underbrace{p(\boldsymbol{\mu} | \mathbf{x}^*, \boldsymbol{\theta})}_{\text{Model}} p(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\mu} d\boldsymbol{\theta}$$

Distributions over a simplex: a Dirichlet, Mixture of Dirichlet distributions or the Logistic-Normal distribution

$$\text{Dir}(\boldsymbol{\mu} | \boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\prod_{c=1}^K \Gamma(\alpha_c)} \prod_{c=1}^K \mu_c^{\alpha_c-1}, \quad \alpha_c > 0, \alpha_0 = \sum_{c=1}^K \alpha_c$$

Higher values of α_0 lead to sharper distributions

Uncertainty $\frac{K}{\alpha_0}$

Dirichlet Prior Network

A Prior Network which parametrizes a Dirichlet will be referred to as a Dirichlet Prior Network (DPN). A DPN will generate the concentration parameters α of the Dirichlet distribution.

$$p(\boldsymbol{\mu} \mid \mathbf{x}^*; \hat{\boldsymbol{\theta}}) = \text{Dir}(\boldsymbol{\mu} \mid \boldsymbol{\alpha}), \quad \boldsymbol{\alpha} = \mathbf{f}(\mathbf{x}^*; \hat{\boldsymbol{\theta}})$$

The posterior over class labels will be given by the mean of the Dirichlet:

$$P(\omega_c \mid \mathbf{x}^*; \hat{\boldsymbol{\theta}}) = \int p(\omega_c \mid \boldsymbol{\mu}) p(\boldsymbol{\mu} \mid \mathbf{x}^*; \hat{\boldsymbol{\theta}}) d\boldsymbol{\mu} = \frac{\alpha_c}{\alpha_0}$$

If an exponential output function is used for the DPN, where $\alpha_c = e^{z_c}$, then the expected posterior probability of a label ω_c is given by the output of the softmax

$$P(\omega_c \mid \mathbf{x}^*; \hat{\boldsymbol{\theta}}) = \frac{e^{z_c(\mathbf{x}^*)}}{\sum_{k=1}^K e^{z_k(\mathbf{x}^*)}}$$

Training loss

Minimize the KL divergence between

- the model and a sharp Dirichlet distribution focused on the appropriate class for in-distribution data
- the model and a flat Dirichlet distribution for out-of-distribution data

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{p}_{\text{in}}(\mathbf{x})}[KL[\text{Dir}(\boldsymbol{\mu} \mid \hat{\boldsymbol{\alpha}}) \parallel \mathbf{p}(\boldsymbol{\mu} \mid \mathbf{x}; \boldsymbol{\theta})]] + \mathbb{E}_{\mathbf{p}_{\text{out}}(\mathbf{x})}[KL[\text{Dir}(\boldsymbol{\mu} \mid \tilde{\boldsymbol{\alpha}}) \parallel \mathbf{p}(\boldsymbol{\mu} \mid \mathbf{x}; \boldsymbol{\theta})]]$$

It is simple to specify a **flat Dirichlet distribution** by setting all $\tilde{\alpha}_c = 1$

The **in-distribution** target $\hat{\alpha}_c, \hat{\mu}_c = \frac{\hat{\alpha}_c}{\hat{\alpha}_0}$

$$\hat{\mu}_c = \begin{cases} 1 - (K - 1)\epsilon & \text{if } \delta(y = \omega_c) = 1 \\ \epsilon & \text{if } \delta(y = \omega_c) = 0 \end{cases}$$

Measures

Expected predictive categorical $P(\omega_c \mid \mathbf{x}^*; \mathcal{D})$

Max probability: measure of confidence in the prediction

$$\mathcal{P} = \max_c P(\omega_c \mid \mathbf{x}^*; \mathcal{D})$$

Entropy: entropy of the predictive distribution, behaves similar to max probability, represents the uncertainty encapsulated in the entire distribution

$$\mathcal{H}[P(y \mid \mathbf{x}^*; \mathcal{D})] = - \sum_{c=1}^K P(\omega_c \mid \mathbf{x}^*; \mathcal{D}) \ln(P(\omega_c \mid \mathbf{x}^*; \mathcal{D}))$$

Measures: MI

Mutual Information (MI) between the categorical label y and the parameters of the model θ is a measure of the spread of an ensemble $\{P(\omega_c | \mathbf{x}^*, \theta^{(i)})\}_{i=1}^M$ which assess uncertainty in predictions due to model uncertainty.

Here, MI implicitly captures elements of distributional uncertainty.

$$\underbrace{\mathcal{I}[y, \theta | \mathbf{x}^*, \mathcal{D}]}_{\text{Model Uncertainty}} = \underbrace{\mathcal{H}[\mathbb{E}_{p(\theta|\mathcal{D})}[P(y|\mathbf{x}^*, \theta)]]}_{\text{Total Uncertainty}} - \underbrace{\mathbb{E}_{p(\theta|\mathcal{D})}[\mathcal{H}[P(y|\mathbf{x}^*, \theta)]]}_{\text{Expected Data Uncertainty}}$$

MI between y and μ , the spread is now explicitly due to distributional uncertainty

$$\underbrace{\mathcal{I}[y, \mu | \mathbf{x}^*; \mathcal{D}]}_{\text{Distributional Uncertainty}} = \underbrace{\mathcal{H}[\mathbb{E}_{p(\mu|\mathbf{x}^*; \mathcal{D})}[P(y | \mu)]]}_{\text{Total Uncertainty}} - \underbrace{\mathbb{E}_{p(\mu|\mathbf{x}^*; \mathcal{D})}[\mathcal{H}[P(y | \mu)]]}_{\text{Expected Data Uncertainty}}$$

Entropy, MI

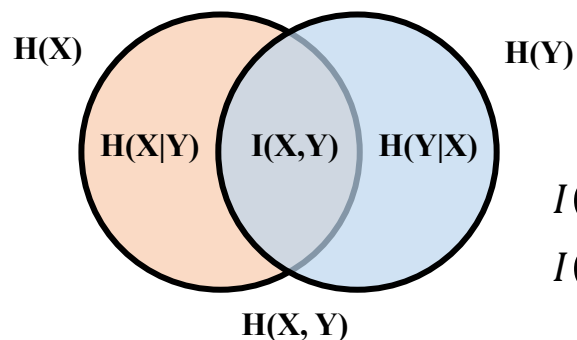
$$\mathcal{I}[y, \boldsymbol{\theta} \mid \mathbf{x}^*, \mathcal{D}] = \mathcal{H}[y \mid \mathbf{x}^*, \mathcal{D}] - \mathcal{H}[y \mid \boldsymbol{\theta}, \mathbf{x}^*]$$

Model
uncertainty

$$\begin{aligned} &= - \int P(y \mid \mathbf{x}^*, \mathcal{D}) \log P(y \mid \mathbf{x}^*, \mathcal{D}) dy + \int \int P(y, \boldsymbol{\theta} \mid \mathbf{x}^*, \mathcal{D}) \log P(y \mid \mathbf{x}^*, \boldsymbol{\theta}) dy d\boldsymbol{\theta} \\ &= - \int \left(\int P(y \mid \mathbf{x}^*, \boldsymbol{\theta}) P(\boldsymbol{\theta} \mid \mathcal{D}) d\boldsymbol{\theta} \right) \log \left(\int P(y \mid \mathbf{x}^*, \boldsymbol{\theta}) P(\boldsymbol{\theta} \mid \mathcal{D}) d\boldsymbol{\theta} \right) dy \\ &\quad + \int \int P(y \mid \mathbf{x}^*, \boldsymbol{\theta}) P(\boldsymbol{\theta} \mid \mathcal{D}) \log P(y \mid \mathbf{x}^*, \boldsymbol{\theta}) dy d\boldsymbol{\theta} \\ &= \mathcal{H}[\mathbb{E}_{P(\boldsymbol{\theta} \mid \mathcal{D})}[y \mid \mathbf{x}^*, \boldsymbol{\theta}]] + \int P(\boldsymbol{\theta} \mid \mathcal{D}) \left(\int P(y \mid \mathbf{x}^*, \boldsymbol{\theta}) \log P(y \mid \mathbf{x}^*, \boldsymbol{\theta}) dy \right) d\boldsymbol{\theta} \\ &= \mathcal{H}[\mathbb{E}_{P(\boldsymbol{\theta} \mid \mathcal{D})}[y \mid \mathbf{x}^*, \boldsymbol{\theta}]] - \mathbb{E}_{P(\boldsymbol{\theta} \mid \mathcal{D})}[\mathcal{H}[y \mid \mathbf{x}^*, \boldsymbol{\theta}]] \end{aligned}$$

Total uncertainty

Expected Data Uncertainty



$$I(X, Y) = H(X) - H(X|Y)$$

$$I(X, Y) = H(Y) - H(Y|X)$$



$$\begin{aligned} \int \int P(x, y) \log \frac{P(x, y)}{P(x)P(y)} dx dy &= - \int \int P(x, y) \log P(x) dx dy - \int \int P(x, y) \log P(y) dx dy \\ &= - \int P(x) \log P(x) dx - \int P(y) \log P(y) dy \\ &= H(X) - H(X|Y) \end{aligned}$$

Measures: the differential entropy

The differential entropy: **maximized when the Dirichlet Distribution is flat**

$$\mathcal{H}[\mathbf{p}(\boldsymbol{\mu}|\mathbf{x}^*; \mathcal{D})] = - \int_{\mathcal{S}^{K-1}} \mathbf{p}(\boldsymbol{\mu}|\mathbf{x}^*; \mathcal{D}) \ln(\mathbf{p}(\boldsymbol{\mu}|\mathbf{x}^*; \mathcal{D})) d\boldsymbol{\mu}$$

Distribution Uncertainty

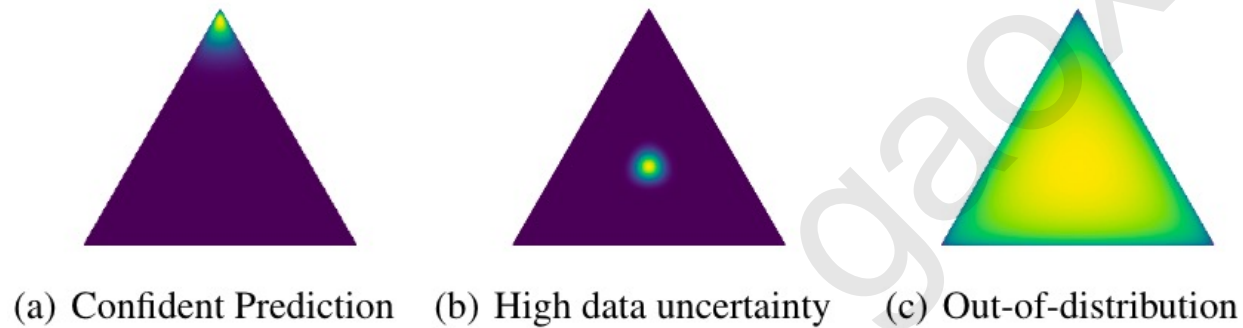
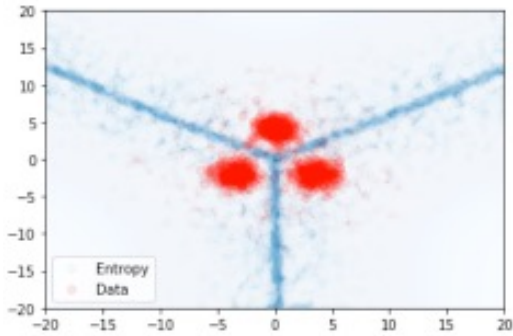


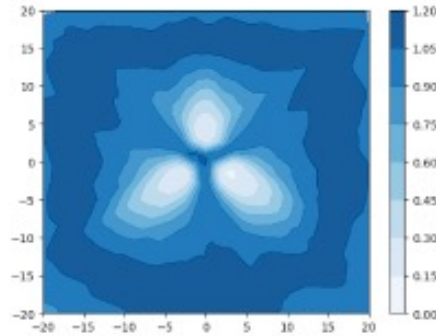
Figure 2: Desired behaviors of a distribution over distributions

- (a) Sharp distribution,
concentrated categorical prediction
- (b) Sharp distribution,
equiprobable categorical prediction
- (c) Flat distribution,**
equiprobable categorical prediction

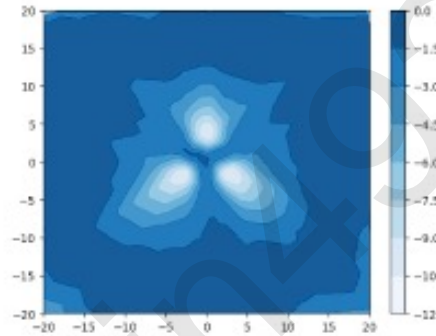
Experiments and results



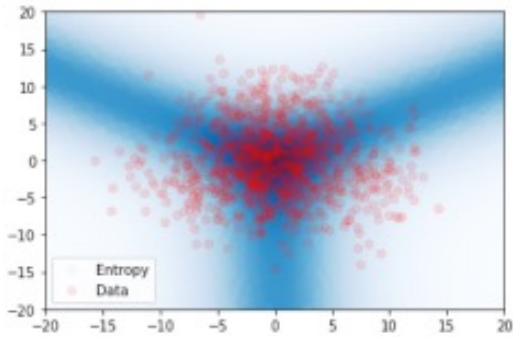
(a) $\sigma = 1$



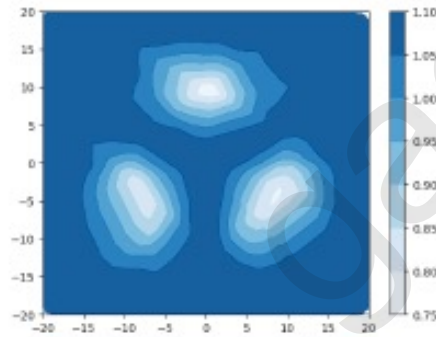
(b) Entropy $\sigma = 1$



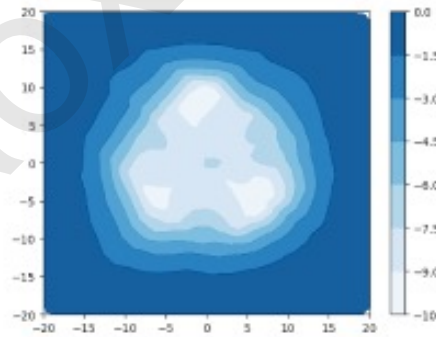
(c) Diff. Entropy $\sigma = 1$



(d) $\sigma = 4$



(e) Entropy $\sigma = 4$



(f) Diff. Entropy $\sigma = 4$

- **class overlap**

Entropy is high both in region of class overlap and far from training data

- difficult to distinguish out-of-distribution samples and in-distribution samples at a decision boundary

Differential entropy is low over the whole region of training data and high outside

- allowing the in-distribution region to be clearly distinguished from the out-of-distribution region

Experiments and results

MNIST and CIFAR-10 are low data uncertainty datasets - all classes are distinct

Differential entropy of the Dirichlet prior will be able to distinguish in-domain and out-of-distribution data better than entropy when the classes are less distinct.

OOD: positive class **ID: negative class**

Table 3: MNIST vs OMNIGLOT. Out-of-distribution detection AUROC on noisy data.

σ	Ent.		M.I.		D.Ent.	
	0.0	3.0	0.0	3.0	0.0	3.0
DNN	98.8	58.4	-	-	-	-
MCDP	98.8	58.4	99.3	79.1	-	-
DPN	100.0	51.8	99.5	22.3	100.0	99.8

total

model

distribution

zero mean isotropic
Gaussian noise
with a standard
deviation $\sigma=3$ noise

Uncertainty Estimation by Fisher Information-based Evidential Deep Learning

2023 ICML

Xin Gao

2024.1.5

Introduction

- It is not sensitive to arbitrary scaling of α_k classical EDL hinders the learning of evidence, especially for samples with high data uncertainty annotated with the one-hot label.
- We propose a simple and novel method, Fisher Information-based Evidential Deep Learning (I-EDL), to **weigh the importance of different classes for each training sample**.
- We introduce **PAC-Bayesian bound** to further improve the generalization ability.
- Our proposed method consistently outperforms traditional EDL-related algorithms in multiple uncertainty estimation tasks, in the confidence evaluation, OOD detection, and few-shot classification.

DUM and EDL

- **Dirichlet-based uncertainty models** quantify different types of uncertainty by modeling the output as the concentration parameters of a Dirichlet distribution.
- **Evidential deep learning (EDL)** adopts Dirichlet distribution and treats the output as evidence to quantify belief mass and uncertainty by jointly considering the Dempster–Shafer Theory of Evidence (DST) and subjective logic (SL).

State space: K mutually exclusive singletons (e.g., class labels)

=> belief mass, uncertainty mass $u + \sum_{k=1}^K b_k = 1$

=> Dirichlet prior, evidence $\text{Dir}(\mathbf{p} \mid \boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K p_k^{\alpha_k - 1}, \alpha_0 = \sum_{k=1}^K \alpha_k \quad \alpha_k = e_k + 1$

=> assign belief and uncertainty $b_k = \frac{\alpha_k - 1}{\alpha_0}, \quad u = \frac{K}{\alpha_0}$

=> point-estimated categorical prediction $\hat{p}_k = \frac{\alpha_k}{\alpha_0} = \frac{e_k + 1}{\sum_{c=1}^K e_c + K}$

Graphic Representation

- EDL supposes the observed labels y were drawn i.i.d. from an isotropic Gaussian distribution, i.e.

$$\mathbf{y} \sim \mathcal{N}(\mathbf{p}, \sigma^2 \mathbf{I})$$

where $p \sim \text{Dir}(f_\theta(x) + 1)$.

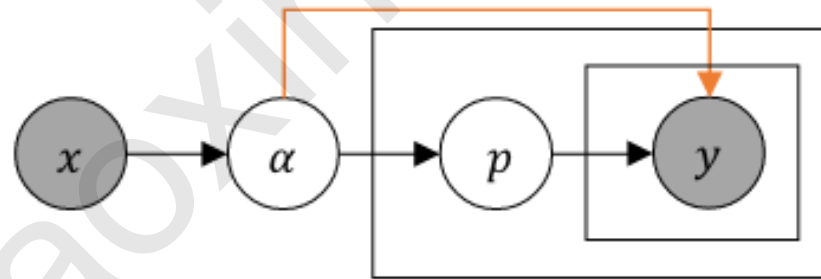
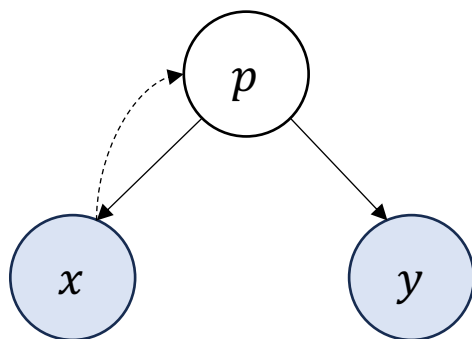


Figure 2. Graphical model representation of \mathcal{I} -EDL.

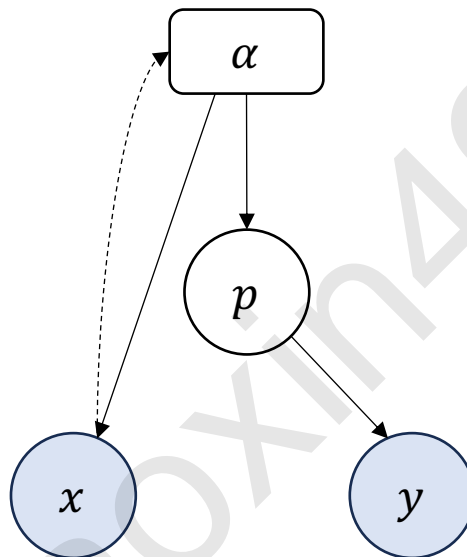
Training evidential neural networks by minimizing the expected MSE can be viewed as learning model parameters that maximize the expected likelihood of the observed labels.

Graphic Representation

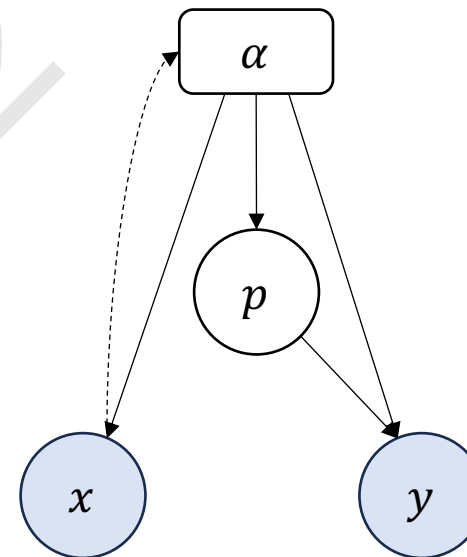
Classical DNN



EDL



I-EDL



x : Observed images

y : Observed labels

p : Probability map

α : Parameter of Dirichlet distribution

Solid arrows indicate generation while dashed ones refer to inference procedure from a neural network.

Higher evidence & Higher variance

- EDL supposes the observed labels \mathbf{y} were drawn i.i.d. from an isotropic Gaussian distribution, i.e.

$$\mathbf{y} \sim \mathcal{N}(\mathbf{p}, \sigma^2 \mathbf{I})$$

where $\mathbf{p} \sim \text{Dir}(f_\theta(x) + 1)$.

$$-\log \mathcal{N}(\mathbf{y}_i \mid \mathbf{p}_i, \mathbf{\Sigma}) = \frac{1}{2}(\mathbf{y}_i - \mathbf{p}_i)^T \mathbf{\Sigma}(\mathbf{y}_i - \mathbf{p}_i) + \frac{1}{2} \log |\mathbf{\Sigma}| + \text{const}$$

- The information of each class carried in categorical probabilities \mathbf{p} is different, thus the generation of each class for a specific sample should not be isotropic.

$$\mathbf{y} \sim \mathcal{N}(\mathbf{p}, \sigma^2 \mathcal{I}(\boldsymbol{\alpha})^{-1})$$

Fisher information matrix

- The Fisher information is a way of measuring the amount of information that an observable random variable X carries about an unknown parameter θ of a distribution that models X .


$$x_1, x_2, \dots, x_n \quad p(x; \theta)$$

- To assess the goodness of our estimate of θ we define a **score function** $s(\theta) = \nabla_{\theta} \log p(x | \theta)$
- The **expected value of score** wrt. our model is **zero**

$$\mathbb{E}_{p(x|\theta)} [s(\theta)] = \mathbb{E}_{p(x|\theta)} [\nabla \log p(x | \theta)] = 0$$

- The **covariance of score function** above is the definition of **Fisher Information Matrix**

$$\mathbf{F} = \mathbb{E}_{p(x|\theta)} [(s(\theta) - 0)(s(\theta) - 0)^T] = \mathbb{E}_{p(x|\theta)} [\nabla \log p(x | \theta) \nabla \log p(x | \theta)^T]$$

- The **negative expected Hessian of log likelihood** is equal to the **Fisher Information Matrix** \mathbf{F}

$$\mathbf{F} = -\mathbb{E}_{p(x|\theta)} [\mathbf{H}_{\log p(x|\theta)}]$$

Insights about FIM

In our context, **the Fisher information matrix (FIM)** is chosen to measure the amount of information that the categorical probabilities p carry about the concentration parameters α of a Dirichlet distribution that models p .

$$\ell = \log \text{Dir}(\mathbf{p} \mid \boldsymbol{\alpha})$$

$$\mathcal{I}(\boldsymbol{\alpha}) = \mathbb{E}_{\text{Dir}(\mathbf{p}|\boldsymbol{\alpha})} \left[\frac{\partial \ell}{\partial \boldsymbol{\alpha}} \frac{\partial \ell}{\partial \boldsymbol{\alpha}^T} \right] = \mathbb{E}_{\text{Dir}(\mathbf{p}|\boldsymbol{\alpha})} \left[-\frac{\partial^2 \ell}{\partial \boldsymbol{\alpha} \partial \boldsymbol{\alpha}^T} \right]$$

$$\mathcal{I}(\boldsymbol{\alpha}) = \text{diag} \left(\left[\psi^{(1)}(\alpha_1), \dots, \psi^{(1)}(\alpha_K) \right] \right) - \psi^{(1)}(\alpha_0) \mathbf{1}\mathbf{1}^T$$

$$\mathbf{y} \sim \mathcal{N}(\mathbf{p}, \sigma^2 \mathcal{I}(\boldsymbol{\alpha})^{-1})$$

$\alpha_k < \alpha_0$, trigamma function is a monotonically decreasing function when $x > 0$

MLE

- In MLE, we can learn model parameters θ by minimizing the expected negative log-likelihood loss function:

$$\begin{aligned} \min_{\theta} \quad & \mathbb{E}_{(x,y) \sim \mathcal{P}} \mathbb{E}_{p \sim \text{Dir}(\alpha)} \left[-\log p(\mathbf{y} \mid \mathbf{p}, \alpha, \sigma^2) \right] \\ \text{s.t.} \quad & \alpha = f_{\theta}(\mathbf{x}) + 1 \\ & \mathcal{I}(\alpha) = \mathbb{E}_{\text{Dir}(\mathbf{p} \mid \alpha)} \left[-\frac{\partial^2 \log \text{Dir}(\mathbf{p} \mid \alpha)}{\partial \alpha \alpha^T} \right] \\ & p(\mathbf{y} \mid \mathbf{p}, \alpha, \sigma^2) = \mathcal{N}(\mathbf{y} \mid \mathbf{p}, \sigma^2 \mathcal{I}(\alpha)^{-1}) \end{aligned}$$

- General loss can improve generalization but is intractable $(x, y) \sim P$
- We can find an upper bound of this optimization problem, converting general loss into empirical loss.

PAC-Bayesian Bound

- This theory focuses on **the upper bound** of the probability of **generalization error** for a **model output** by a learning algorithm, **given a certain data distribution**.

Theorem 3.1 ((Germain et al., 2009; Alquier et al., 2016; Masegosa, 2020)). Given a data distribution \mathcal{P} over $\mathcal{X} \times \mathcal{Y}$, a hypothesis set Θ , a prior distribution π over Θ , for any $\delta \in (0, 1]$, and $\lambda > 0$, with probability at least $1 - \delta$ over samples $\mathcal{D} \sim \mathcal{P}^n$, we have for all posterior ρ ,

$$\begin{aligned} \mathbb{E}_{\rho(\theta)}[\mathcal{L}(\theta)] &\leq \mathbb{E}_{\rho(\theta)}[\hat{\mathcal{L}}_{\mathcal{D}}(\theta)] \\ &\quad + \frac{1}{\lambda} \left[D_{\text{KL}}(\rho \parallel \pi) + \log \frac{1}{\delta} + \Psi_{\mathcal{P}, \pi}(\lambda, n) \right], \end{aligned}$$

RV: $\theta \Rightarrow p$
 $p \sim \text{Dir}(p|\alpha)$
 $\alpha = f_{\theta}(x) + 1$

where $\Psi_{\mathcal{P}, \pi}(\lambda, n) = \log \mathbb{E}_{\pi(\theta)} \mathbb{E}_{\mathcal{D} \sim \mathcal{P}^n} \left[e^{\lambda(\mathcal{L}(\theta) - \hat{\mathcal{L}}_{\mathcal{D}}(\theta))} \right]$

1. **Prior Distribution, π :** The distribution over the hypothesis set before observing any data. It reflects our initial beliefs about the parameters.
2. **Posterior, ρ :** After observing data, our beliefs about the hypothesis set are updated, leading to the posterior distribution.

Upper Bound

- In this paper, we treat $Dir(p|\alpha)$ as the **posterior distribution**, and the **prior** as $Dir(p|\mu)$, where μ is set to $\beta \gg 1$ for the corresponding class and 1 for all other class.

- The **upper bound** of the optimization problem in MLE can be expressed as

$$\frac{1}{N} \sum_{i=1}^N \mathcal{L}_i(\theta) + \frac{1}{\lambda} D_{\text{KL}}(\text{Dir}(\mathbf{p}_i | \alpha_i) \| \text{Dir}(\mathbf{p}_i | \mu_i))$$

where $\mathcal{L}_i(\theta) = \mathbb{E}_{\text{Dir}(\mathbf{p}_i|\alpha_i)} \left[-\log \mathcal{N}(\mathbf{y}_i | \mathbf{p}_i, \sigma^2 \mathcal{I}(\alpha_i)^{-1}) \right]$

- The **first term** is the expected FIM-weighted MSE subtract the negative log determinant of the FIM:

$$\mathcal{L}_i(\theta) \propto \underbrace{\mathbb{E} \left[(\mathbf{y}_i - \mathbf{p}_i)^T \mathcal{I}(\alpha_i) (\mathbf{y}_i - \mathbf{p}_i) \right]}_{\mathcal{L}_i^{\mathcal{I}\text{-MSE}}} - \underbrace{\sigma^2 \log |\mathcal{I}(\alpha_i)|}_{\mathcal{L}_i^{|\mathcal{I}|}}$$

- The **second term** can be simplified by setting $\hat{\alpha}_i = \alpha_i \odot (1 - \mathbf{y}_i) + \mathbf{y}_i$ as Sensoy et al.

$$\mathcal{L}_i^{\text{KL}} = D_{\text{KL}}(\text{Dir}(\mathbf{p}_i | \hat{\alpha}_i) \| \text{Dir}(\mathbf{p}_i | \mathbf{1}))$$

MLE & MSE & Cross-entropy

$$\max \log \mathcal{P}(y; \theta) = \max \sum_{i=1}^n \log \mathcal{P}(y_i; \theta)$$

- Gaussian

$$\sum_{i=1}^n \log \mathcal{N}(y_i \mid f_{\theta}(x_i), \Sigma) = -\frac{1}{2} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^T \Sigma (y_i - f_{\theta}(x_i)) - \frac{n}{2} \log |\Sigma| + \text{const}$$

$$\sum_{i=1}^n \log \mathcal{N}(y_i \mid f_{\theta}(x_i), I) = -\frac{1}{2} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^T (y_i - f_{\theta}(x_i)) + \text{const} \quad \Rightarrow \quad \text{MSE}$$

- Bernoulli

$$y \sim B(y, f_{\theta}(x))$$

$$\sum_{i=1}^n \log [f_{\theta}(x_i)]^{y_i} [1 - f_{\theta}(x_i)]^{(1-y_i)} = \sum_{i=1}^n y_i \log f_{\theta}(x_i) + (1 - y_i) \log(1 - f_{\theta}(x_i))$$

\Rightarrow Cross entropy

Objective function

- Finally, the objective function Eq.(2) can be reformulated as

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \mathcal{L}_i^{\mathcal{I}\text{-MSE}} - \lambda_1 \mathcal{L}_i^{|\mathcal{I}|} + \lambda_2 \mathcal{L}_i^{\text{KL}}$$

classical EDL can be viewed as
a degenerate version of I-EDL

$$\mathcal{L}_i^{\mathcal{I}\text{-MSE}} = \sum_{j=1}^K \left((y_{ij} - \frac{\alpha_{ij}}{\alpha_{i0}})^2 + \frac{\alpha_{ij}(\alpha_{i0} - \alpha_{ij})}{\alpha_{i0}^2(\alpha_{i0} + 1)} \right) \psi^{(1)}(\alpha_{ij}),$$

$$\mathcal{L}_i^{|\mathcal{I}|} = \sum_{j=1}^K \log \psi^{(1)}(\alpha_{ij}) + \log \left(1 - \sum_{j=1}^K \frac{\psi^{(1)}(\alpha_{i0})}{\psi^{(1)}(\alpha_{ij})} \right),$$

$$\begin{aligned} \mathcal{L}_i^{\text{KL}} = & \log \Gamma(\sum_{j=1}^K \hat{\alpha}_{ij}) - \log \Gamma(K) - \sum_{j=1}^K \log \Gamma(\hat{\alpha}_{ij}) \\ & + \sum_{j=1}^K (\hat{\alpha}_{ij} - 1) \left[\psi(\hat{\alpha}_{ij}) - \psi(\sum_{k=1}^K \hat{\alpha}_{ik}) \right], \end{aligned}$$

Table 1. Difference between \mathcal{I} -EDL and EDL are marked in blue.

	EDL	\mathcal{I} -EDL
MSE	$\sum_{i=1}^K (y_i - \frac{\alpha_i}{\alpha_0})^2$ $+ \sum_{i=1}^K \frac{\alpha_i(\alpha_0 - \alpha_i)}{\alpha_0^2(\alpha_0 + 1)}$	$\sum_{i=1}^K (y_i - \frac{\alpha_i}{\alpha_0})^2 \psi^{(1)}(\alpha_i)$ $+ \sum_{i=1}^K \frac{\alpha_i(\alpha_0 - \alpha_i)}{\alpha_0^2(\alpha_0 + 1)} \psi^{(1)}(\alpha_i)$
KL	$D_{\text{KL}}(\text{Dir}(\hat{\alpha}) \parallel \text{Dir}(\mathbf{1}))$	$D_{\text{KL}}(\text{Dir}(\hat{\alpha}) \parallel \text{Dir}(\mathbf{1}))$
\mathcal{I}	-	$-\log \mathcal{I}(\alpha) $

- For different labels in a sample
Though it has been correctly classified for a specific label, it still allows for more evidence for the overlapping labels.

Objective function

Uncertainty $\frac{K}{\alpha_0}$

```
def compute_fisher_mse(self, labels_1hot_, evi_alp_):
    evi_alp0_ = torch.sum(evi_alp_, dim=-1, keepdim=True)

    gamma1_alp = torch.polygamma(1, evi_alp_)
    gamma1_alp0 = torch.polygamma(1, evi_alp0_)

    gap = labels_1hot_ - evi_alp_ / evi_alp0_

    loss_mse_ = (gap.pow(2) * gamma1_alp).sum(-1).mean()

    loss_var_ = (evi_alp_ * (evi_alp0_ - evi_alp_) * gamma1_alp / (evi_alp0_ * evi_alp0_ * (evi_alp0_ + 1))).sum(-1).mean()

    loss_det_fisher_ = - (torch.log(gamma1_alp).sum(-1) + torch.log(1.0 - (gamma1_alp0 / gamma1_alp).sum(-1))).mean()

    return loss_mse_, loss_var_, loss_det_fisher_
```

Table 1. Difference between \mathcal{I} -EDL and EDL are marked in blue.

	EDL	\mathcal{I} -EDL
MSE	$\sum_{i=1}^K (y_i - \frac{\alpha_i}{\alpha_0})^2$ $+ \sum_{i=1}^K \frac{\alpha_i(\alpha_0 - \alpha_i)}{\alpha_0^2(\alpha_0 + 1)}$	$\sum_{i=1}^K (y_i - \frac{\alpha_i}{\alpha_0})^2 \psi^{(1)}(\alpha_i)$ $+ \sum_{i=1}^K \frac{\alpha_i(\alpha_0 - \alpha_i)}{\alpha_0^2(\alpha_0 + 1)} \psi^{(1)}(\alpha_i)$
KL	$D_{\text{KL}}(\text{Dir}(\hat{\alpha}) \parallel \text{Dir}(\mathbf{1}))$	$D_{\text{KL}}(\text{Dir}(\hat{\alpha}) \parallel \text{Dir}(\mathbf{1}))$
\mathcal{I}	-	$-\log \mathcal{I}(\alpha) $

<https://github.com/danruod/IEDL>

Experiments

- OOD detection

Table 3. AUPR scores of OOD detection (mean \pm standard deviation of 5 runs). [†] indicates that the first four lines are the results reported by Charpentier et al. (2020). Bold and underlined numbers indicate the best and runner-up scores, respectively.

Method	MNIST \rightarrow KMNIST [†]		MNIST \rightarrow FMNIST [†]		CIFAR10 \rightarrow SVHN [†]		CIFAR10 \rightarrow CIFAR100	
	Max.P	α_0	Max.P	α_0	Max.P	α_0	Max.P	α_0
Dropout	94.00 \pm 0.1	-	96.56 \pm 0.2	-	51.39 \pm 0.1	-	45.57 \pm 1.0	-
KL-PN	92.97 \pm 1.2	93.39 \pm 1.0	<u>98.44 \pm 0.1</u>	<u>98.16 \pm 0.0</u>	43.96 \pm 1.9	43.23 \pm 2.3	61.41 \pm 2.8	61.53 \pm 3.4
RKL-PN	60.76 \pm 2.9	53.76 \pm 3.4	<u>78.45 \pm 3.1</u>	<u>72.18 \pm 3.6</u>	53.61 \pm 1.1	49.37 \pm 0.8	55.42 \pm 2.6	54.74 \pm 2.8
PostN	95.75 \pm 0.2	94.59 \pm 0.3	97.78 \pm 0.2	97.24 \pm 0.3	<u>80.21 \pm 0.2</u>	77.71 \pm 0.3	81.96 \pm 0.8	82.06 \pm 0.8
EDL	97.02 \pm 0.8	<u>96.31 \pm 2.0</u>	98.10 \pm 0.4	98.08 \pm 0.4	78.87 \pm 3.5	<u>79.12 \pm 3.7</u>	<u>84.30 \pm 0.7</u>	<u>84.18 \pm 0.7</u>
\mathcal{I}-EDL	98.34 \pm 0.2	98.33 \pm 0.2	98.89 \pm 0.3	98.86 \pm 0.3	83.26 \pm 2.4	82.96 \pm 2.2	85.35 \pm 0.7	84.84 \pm 0.6

We mainly focus on the comparisons with DBU models, which solve OOD detection by distinguishing different types of uncertainty.

Experiments

- Few-shot Learning

Table 4. Classification accuracy (Acc.), AUPR scores for both confidence evaluation (Conf.) and OOD detection (OOD) under $\{5, 10\}$ -way $\{1, 5, 20\}$ -shot settings of mini-ImageNet. CUB is used for OOD detection. Each experiment is run for over 10,000 few-shot episodes.

Method	5-Way 1-shot			5-Way 5-shot			5-way 20-shot		
	Acc.	Conf. (Max. α)	OOD (α_0)	Acc.	Conf. (Max. α)	OOD (α_0)	Acc.	Conf. (Max. α)	OOD (α_0)
EDL	61.00 ± 0.22	80.59 ± 0.23	65.40 ± 0.26	80.38 ± 0.15	93.92 ± 0.09	76.53 ± 0.27	85.54 ± 0.12	97.51 ± 0.04	79.78 ± 0.23
\mathcal{I} -EDL	63.82 ± 0.20	82.00 ± 0.21	74.76 ± 0.25	82.00 ± 0.14	94.09 ± 0.09	82.48 ± 0.20	88.12 ± 0.09	97.54 ± 0.04	85.40 ± 0.19
Δ	2.82	1.41	9.36	1.62	0.17	5.95	2.58	0.04	5.62

Method	10-Way 1-shot			10-Way 5-shot			10-way 20-shot		
	Acc.	Conf. (Max. α)	OOD (α_0)	Acc.	Conf. (Max. α)	OOD (α_0)	Acc.	Conf. (Max. α)	OOD (α_0)
EDL	44.55 ± 0.15	65.97 ± 0.20	67.83 ± 0.24	62.52 ± 0.16	86.81 ± 0.10	76.34 ± 0.20	69.29 ± 0.17	94.21 ± 0.06	76.88 ± 0.17
\mathcal{I} -EDL	49.37 ± 0.13	68.29 ± 0.19	71.95 ± 0.20	67.89 ± 0.11	87.45 ± 0.09	82.29 ± 0.19	78.60 ± 0.08	94.40 ± 0.04	82.52 ± 0.14
Δ	4.82	2.32	4.12	5.37	0.64	5.95	9.31	0.19	5.64

Our method not only improves classification accuracy but also greatly improves the availability of uncertainty estimation in the more challenging few-shot scenarios.

Experiments

- Density plots of the predicted **differential entropy** and **mutual information** (Last paper, distributional uncertainty)
- **Lower** entropy or mutual information represents the model yields **a sharper distribution**, indicating that the sample has low uncertainty.
- Our method provides **more separable uncertainty estimates**, I-EDL produces sharper prediction peaks than EDL

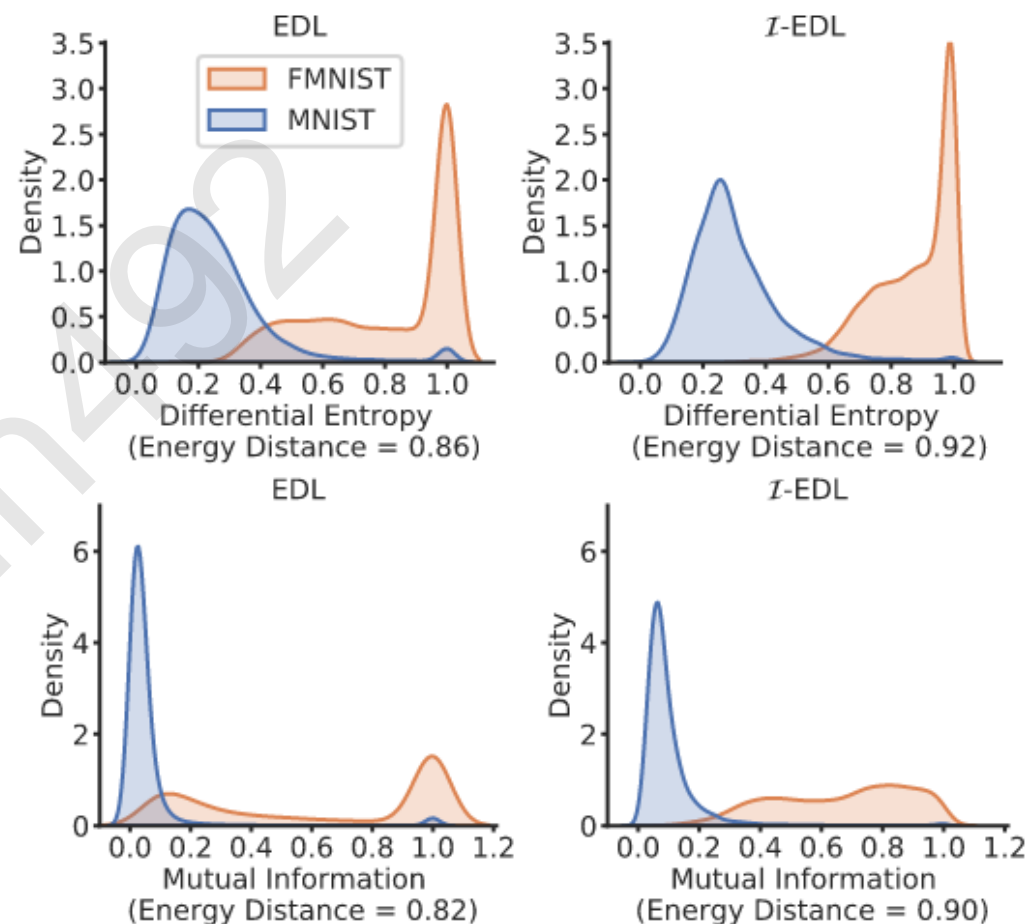


Figure 4. Uncertainty representation for ID (MNIST) and OOD (FMNIST). More results are shown in Appendix C.6.

Conclusion

- The observed label is jointly generated by the predicted categorical probability and the informativeness of each class contained in the sample.
- The informativeness is modeled by the uncertainty of the estimator of α (FIM), naturally including **data uncertainty**.

