



香港中文大學
The Chinese University of Hong Kong



之江實驗室
ZHEJIANG LAB

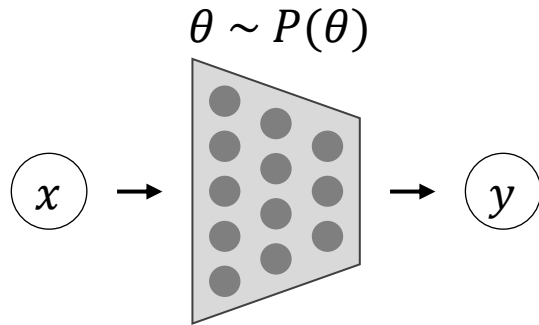
Uncertainty Estimation by Fisher Information-based Evidential Deep Learning

Danruo Deng¹, Guangyong Chen^{2*}, Yang Yu¹, Furui Liu², Pheng-Ann Heng¹

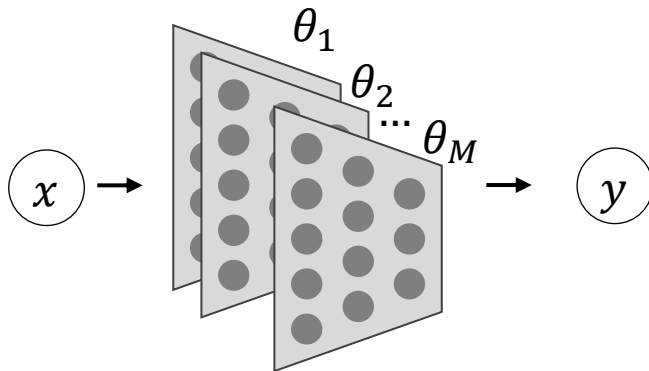
¹The Chinese University of Hong Kong, ²Zhejiang Lab



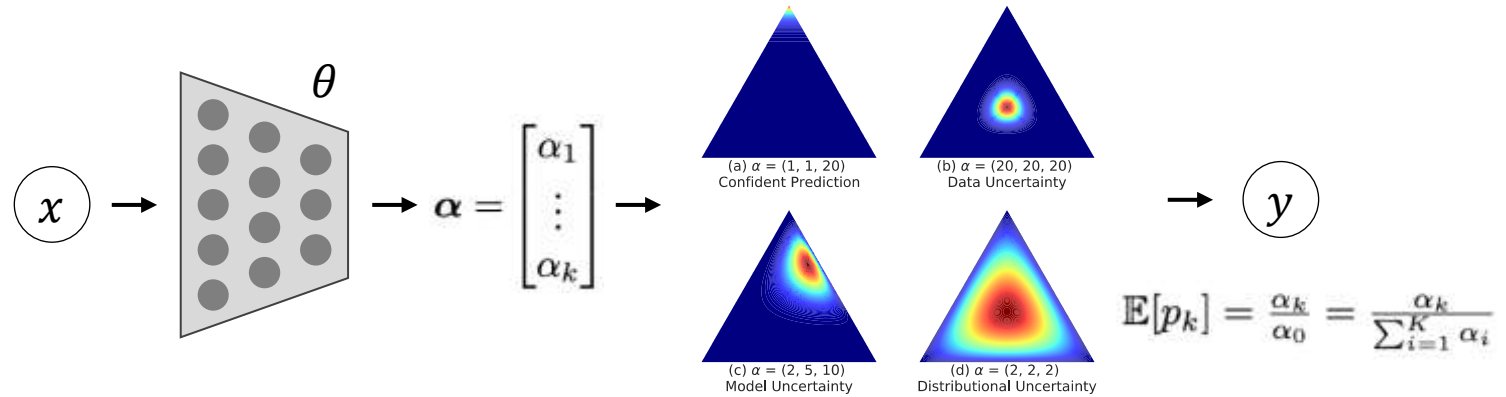
Bayesian Neural Networks (BNNs):



Ensemble Methods:



Evidential Neural Networks [1]:



$$p \sim \text{Dir}(\alpha)$$

$$\text{Dir}(p|\alpha) = \frac{\Gamma(\alpha_0)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K p_k^{\alpha_k - 1}, \alpha_0 = \sum_{k=1}^K \alpha_k$$

- Predictions with *softmax* is **over-confidence** and **cannot distinguish between different uncertainties**.
- BNNs and Ensemble methods are computationally expensive, and also **cannot distinguish between distributional uncertainty and other uncertainties**.
- Evidential neural networks **quantify different types of uncertainty** by modeling the output as the evidence use to obtain concentration parameters of a Dirichlet distribution.

[1] Sensoy et al., Evidential deep learning to quantify classification uncertainty. NeurIPS 2018.

Motivation: EDL underestimates data uncertainty



Limitation: EDL **cannot** distinguish samples with **different data uncertainties**.

Analysis: For samples with **high data uncertainty** but annotated with **one-hot vectors**, the learning process of evidence for those mislabeled classes is **over-penalized and remains hindered**.

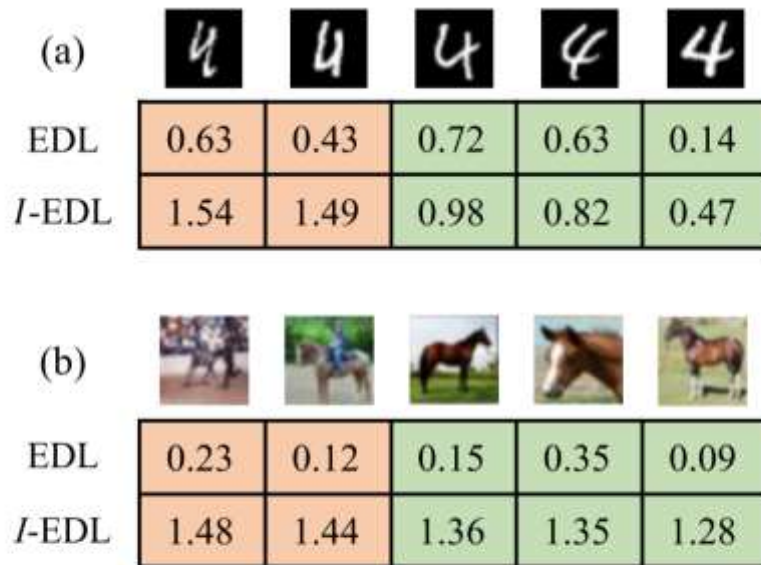
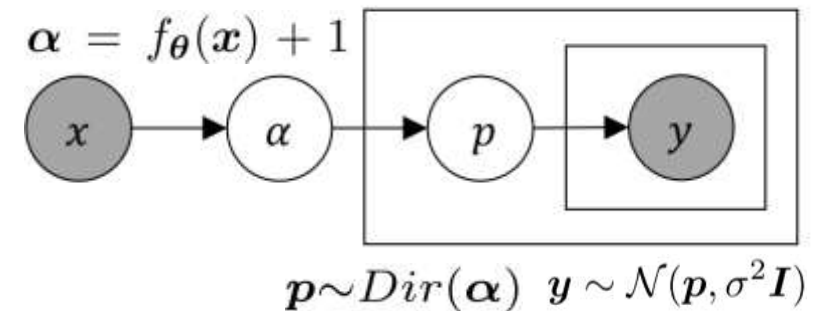


Figure 1. Data uncertainty for (a) digit “4” in MNIST, (b) “horse” in CIFAR10. \mathcal{I} -EDL has the ability to distinguish between hard samples (orange) and easy samples (green), but EDL cannot.

Graphical model representation of EDL:



Objective function:

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} \mathbb{E}_{p \sim \text{Dir}(\alpha)} [(y - p)^T (y - p)]$$

=

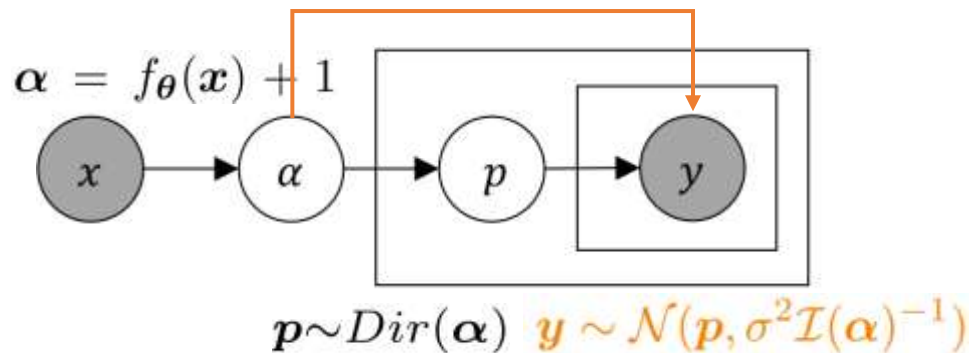
$$\max_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{P}} [\log \mathbb{E}_{p \sim \text{Dir}(\alpha)} [\mathcal{N}(y | p, \sigma^2 I)]]$$



Key Idea: A certain class label with **higher evidence** is allowed to have a **larger variance**, so that **the evidence for missing labels can be preserved** while maximizing the likelihood of the observed labels.

Method: Use Fisher Information Matrix (FIM) to measure the amount of information that **the observed class probabilities \mathbf{p}** carry about the **concentration parameters α** of a Dirichlet distribution that models \mathbf{p} .

Graphical model representation of **I-EDL**:



Objective function:

$$\max_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{P}} [\log \mathbb{E}_{\mathbf{p} \sim \text{Dir}(\boldsymbol{\alpha})} [\mathcal{N}(\mathbf{y} | \mathbf{p}, \sigma^2 \mathbf{I}(\boldsymbol{\alpha})^{-1})]]$$

$$\begin{aligned} \mathbf{I}(\boldsymbol{\alpha}) &= \mathbb{E}_{\text{Dir}(\mathbf{p} | \boldsymbol{\alpha})} \left[\frac{\partial \ell}{\partial \boldsymbol{\alpha}} \frac{\partial \ell}{\partial \boldsymbol{\alpha}^T} \right] \\ &= \text{diag}([\psi^{(1)}(\alpha_1), \dots, \psi^{(1)}(\alpha_K)]) - \psi^{(1)}(\alpha_0) \mathbf{1}\mathbf{1}^T \end{aligned}$$

- $\psi^{(1)}(\cdot)$ denotes the trigamma function, defined as $\psi^{(1)}(x) = d\psi(x)/dx = d^2 \ln \Gamma(x)/dx^2$.
- Since $\psi^{(1)}(x)$ is a monotonically decreasing function when $x > 0$, the class label with higher evidence corresponds to less Fisher information.



$$\max_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{P}} \left[\log \mathbb{E}_{\mathbf{p} \sim \text{Dir}(\boldsymbol{\alpha})} [\mathcal{N}(\mathbf{y} | \mathbf{p}, \sigma^2 \mathcal{I}(\boldsymbol{\alpha})^{-1})] \right]$$

↓ Jensen's inequality

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

↓ PAC-Bayesian Bound

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

where

$$\begin{aligned} \mathcal{L}_i(\boldsymbol{\theta}) &= \mathbb{E}_{\text{Dir}(\mathbf{p}_i | \boldsymbol{\alpha}_i)} [-\log \mathcal{N}(\mathbf{y}_i | \mathbf{p}_i, \sigma^2 \mathcal{I}(\boldsymbol{\alpha}_i)^{-1})] \\ &\propto \underbrace{\mathbb{E} [(\mathbf{y}_i - \mathbf{p}_i)^T \mathcal{I}(\boldsymbol{\alpha}_i) (\mathbf{y}_i - \mathbf{p}_i)]}_{\mathcal{L}_i^{\mathcal{I}\text{-MSE}}} - \underbrace{\sigma^2 \log |\mathcal{I}(\boldsymbol{\alpha}_i)|}_{\mathcal{L}_i^{|\mathcal{I}|}} \end{aligned}$$

Table 1. Given a sample $(\mathbf{x}_i, \mathbf{y}_i)$, the difference in loss function between \mathcal{I} -EDL and EDL are marked in blue.

	EDL	\mathcal{I} -EDL
MSE	$\sum_{j=1}^K (y_{ij} - \frac{\alpha_{ij}}{\alpha_{i0}})^2 + \sum_{j=1}^K \frac{\alpha_{ij}(\alpha_{i0} - \alpha_{ij})}{\alpha_{i0}^2(\alpha_{i0} + 1)}$	$\sum_{j=1}^K (y_{ij} - \frac{\alpha_{ij}}{\alpha_{i0}})^2 \psi^{(1)}(\alpha_{ij}) + \sum_{j=1}^K \frac{\alpha_{ij}(\alpha_{i0} - \alpha_{ij})}{\alpha_{i0}^2(\alpha_{i0} + 1)} \psi^{(1)}(\alpha_{ij})$
KL	$D_{\text{KL}}(\text{Dir}(\hat{\boldsymbol{\alpha}}_i) \ \text{Dir}(\mathbf{1}))$	$D_{\text{KL}}(\text{Dir}(\hat{\boldsymbol{\alpha}}_i) \ \text{Dir}(\mathbf{1}))$
\mathcal{I}	-	$-\log \mathcal{I}(\boldsymbol{\alpha}_i) $



- Standard neural network for classification with Softmax is **over-confidence** and **cannot distinguish between different uncertainties**.
- Although Evidential Deep Learning (EDL) models different types of uncertainties, it still **cannot distinguish between samples of different data uncertainties**.
- we propose a novel and simple method, ***Fisher Information-based Evidential Deep Learning (I-EDL)***, to alleviate the over-penalization of the mislabeled classes by considering importance weights with different classes.
- Extensive experiments on various **image classification, confidence evaluation and OOD detection tasks** demonstrate **the effectiveness of our approach** in achieving **high classification and uncertainty quantification**.



Thanks for your listening.