

Is Epistemic Uncertainty Faithfully Represented by Evidential Deep Learning Methods?

ICML 2024

Mira Jürgens¹, Nis Meinert², Viktor Bengs³, Eyke Hüllermeier³, Willem Waegeman¹

¹Ghent University, ²German Aerospace Center (DLR), ³LMU Munich



What is Evidential Deep Learning (EDL)?

Learn a **second-order distribution** directly by which both *aleatoric* and *epistemic* uncertainty can be disentangled predicted.

Dirichlet-Categorical Model

- level 1: $y \sim \text{Cat}(\boldsymbol{\theta})$ with $\boldsymbol{\theta} \in \Delta_K$
- level 2: $\boldsymbol{\theta} \sim \text{Dir}(\mathbf{m})$ with $\mathbf{m} \in \mathbb{R}_+^K$

Dirichlet-Categorical Model

- level 1: $y \sim \text{Cat}(\theta)$ with $\theta \in \Delta_K$
- level 2: $\theta \sim \text{Dir}(\mathbf{m})$ with $\mathbf{m} \in \mathbb{R}_+^K$

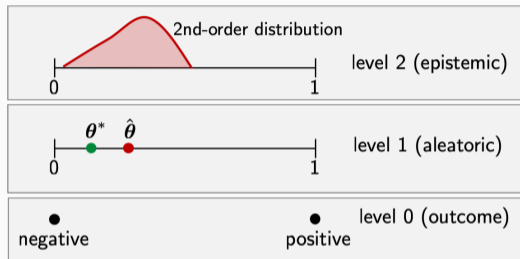


Fig. 1: Different levels of predictions^a (binary classification)

^aWimmer et al., "Quantifying aleatoric and epistemic uncertainty in machine learning: Are conditional entropy and mutual information appropriate measures?"

Normal-Inverse-Gamma Model

- level 1: $y \sim \mathcal{N}(\mu, \sigma^2)$ with $\mu \in \mathbb{R}$, $\sigma^2 \in \mathbb{R}_+$
- level 2: $(\mu, \sigma^2) \sim N\text{-}\Gamma^{-1}(\mathbf{m})$ with $\mathbf{m} = (\gamma, \nu, \alpha, \beta) \in \mathbb{R} \times \mathbb{R}_+^3$

¹Amini et al., “Deep evidential regression”.

²Meinert, Gawlikowski, and Lavin, “The unreasonable effectiveness of deep evidential regression”.

Evidential Deep Learning for Regression

Normal-Inverse-Gamma Model

- level 1: $y \sim \mathcal{N}(\mu, \sigma^2)$ with $\mu \in \mathbb{R}$, $\sigma^2 \in \mathbb{R}_+$
- level 2: $(\mu, \sigma^2) \sim N\text{-}\Gamma^{-1}(\mathbf{m})$ with $\mathbf{m} = (\gamma, \nu, \alpha, \beta) \in \mathbb{R} \times \mathbb{R}_+^3$

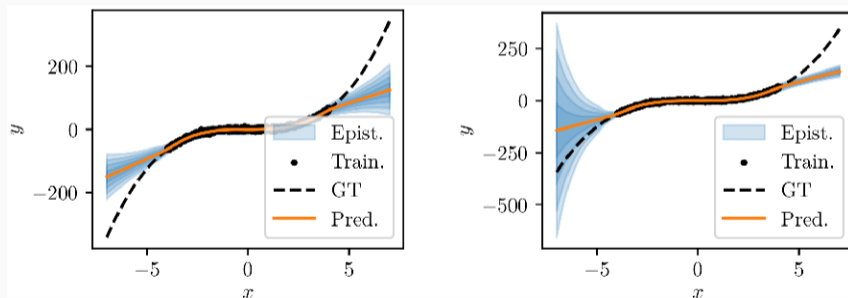


Fig. 2: Results of Deep Evidential Regression¹ reproduced for two different runs²

¹Amini et al., "Deep evidential regression".

²Meinert, Gawlikowski, and Lavin, "The unreasonable effectiveness of deep evidential regression".

Choose your Distribution!

Problem	Likelihood	Conjugate prior
classification ³	$y \sim \text{Cat}(\boldsymbol{\theta})$ $\boldsymbol{\theta} \in \Delta_K$	$\boldsymbol{\theta} \sim \text{Dir}(\mathbf{m})$ $\mathbf{m} \in \mathbb{R}_+^K$
regression (univariate) ⁴	$y \sim \mathcal{N}(\mu, \sigma^2)$ $\mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}_+$	$\boldsymbol{\theta} \sim N\text{-}\Gamma^{-1}(\mathbf{m})$ $\mathbf{m} = (\gamma, \nu, \alpha, \beta) \in \mathbb{R} \times \mathbb{R}_+^3$
regression (multivariate) ⁵	$\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ $\boldsymbol{\mu} \in \mathbb{R}^D, \boldsymbol{\Sigma} \in \mathbb{R}^{D \times D}$	$\boldsymbol{\theta} \sim \text{NIW}(\mathbf{m})$ $\mathbf{m} = (\boldsymbol{\mu}_0, \kappa, \nu, \boldsymbol{\Phi})$
point processes	$y \sim \text{Pois}(\boldsymbol{\theta})$ $\boldsymbol{\theta} \in \mathbb{R}_+$	$\boldsymbol{\theta} \sim \Gamma(\mathbf{m})$ $\mathbf{m} = (\alpha, \beta) \in \mathbb{R}_+^2$

$\boldsymbol{\theta}(\mathbf{x}, \phi) : \mathcal{X} \rightarrow \Theta$ (1st-order predictor), $\mathbf{m}(\mathbf{x}, \phi) : \mathcal{X} \rightarrow \mathcal{M}$ (2nd-order predictor)

³Sensoy, Kaplan, and Kandemir, "Evidential deep learning to quantify classification uncertainty".

⁴Amini et al., "Deep evidential regression".

⁵Meinert and Lavin, "Multivariate deep evidential regression".

While showing good results in downstream tasks, several critical analyses show theoretical flaws of EDL:

- Non-Convergence in the Classification case⁶
- Non-Convergence in the Regression Case⁷
- Non-Properness of its Loss functions⁸

→ **Our approach:** Does the learned **second order distribution** represent the (epistemic) uncertainty of the 1st-order parameters in a faithful way?

⁶Bengs, Hüllermeier, and Waegeman, “Pitfalls of epistemic uncertainty quantification through loss minimisation”.

⁷Meinert, Gawlikowski, and Lavin, “The unreasonable effectiveness of deep evidential regression”.

⁸Bengs, Hüllermeier, and Waegeman, “On second-order scoring rules for epistemic uncertainty quantification”.

Comparing 1st and 2nd-Order Risk Minimization

1st-Order vs 2nd-Order Risk Minimization

1st-Order Risk Minimization

- Loss function (e.g., NLL)

$$L_1 : \mathcal{Y} \times \mathbb{P}(\mathcal{Y}) \rightarrow \mathbb{R}$$

- Risk minimization

$$\min_{\Phi} \sum_{i=1}^N L_1(y_i, p(y | \theta(\mathbf{x}_i; \Phi))) + \lambda R(\Phi)$$

1st-Order vs 2nd-Order Risk Minimization

1st-Order Risk Minimization

- Loss function (e.g., NLL)

$$L_1 : \mathcal{Y} \times \mathbb{P}(\mathcal{Y}) \rightarrow \mathbb{R}$$

- Risk minimization

$$\min_{\Phi} \sum_{i=1}^N L_1(y_i, p(y | \theta(\mathbf{x}_i; \Phi))) + \lambda R(\Phi)$$

2nd-Order Risk Minimization

- Loss function (e.g., NLL)

$$L_1 : \mathcal{Y} \times \mathbb{P}(\mathcal{Y}) \rightarrow \mathbb{R}$$

- *Outer* expectation minimization

$$\min_{\Phi} \sum_{i=1}^N \mathbb{E}_{\theta \sim p(\theta | m(\mathbf{x}_i; \Phi))} [L_1(y_i, p(y | \theta))] + \lambda R(\Phi)$$

Objective: Enforce $p(\theta | \mathbf{m}(\mathbf{x}_i; \Phi))$ to look similar to a predefined distribution $p(\theta | \mathbf{m}_0)$ by minimizing per-instance KL-divergence:

$$R(\Phi) = \sum_{i=1}^N d_{\text{KL}}(p(\theta | \mathbf{m}(\mathbf{x}_i; \Phi)), p(\theta | \mathbf{m}_0))$$

Objective: Enforce $p(\theta | \mathbf{m}(x_i; \Phi))$ to look similar to a predefined distribution $p(\theta | \mathbf{m}_0)$ by minimizing per-instance KL-divergence:

$$R(\Phi) = \sum_{i=1}^N d_{\text{KL}}(p(\theta | \mathbf{m}(x_i; \Phi)), p(\theta | \mathbf{m}_0))$$

- choose $p(\theta | \mathbf{m}_0)$ to parametrize uniform distribution (on 1st-order distributions)

Objective: Enforce $p(\theta | \mathbf{m}(x_i; \Phi))$ to look similar to a predefined distribution $p(\theta | \mathbf{m}_0)$ by minimizing per-instance KL-divergence:

$$R(\Phi) = \sum_{i=1}^N d_{\text{KL}}(p(\theta | \mathbf{m}(x_i; \Phi)), p(\theta | \mathbf{m}_0))$$

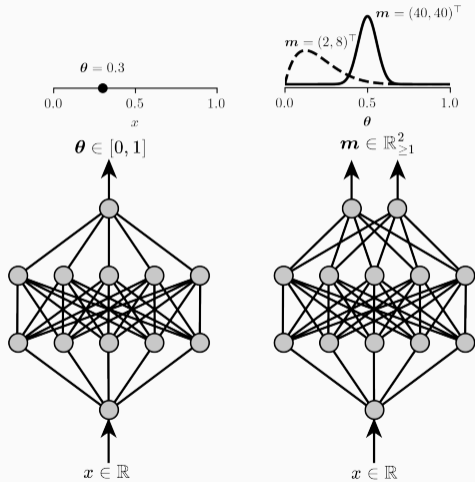
- choose $p(\theta | \mathbf{m}_0)$ to parametrize uniform distribution (on 1st-order distributions)
- minimizing **KL-divergence** \leftrightarrow maximizing **entropy**

Objective: Enforce $p(\theta | \mathbf{m}(x_i; \Phi))$ to look similar to a predefined distribution $p(\theta | \mathbf{m}_0)$ by minimizing per-instance KL-divergence:

$$R(\Phi) = \sum_{i=1}^N d_{\text{KL}}(p(\theta | \mathbf{m}(x_i; \Phi)), p(\theta | \mathbf{m}_0))$$

- choose $p(\theta | \mathbf{m}_0)$ to parametrize uniform distribution (on 1st-order distributions)
- minimizing **KL-divergence** \leftrightarrow maximizing **entropy**

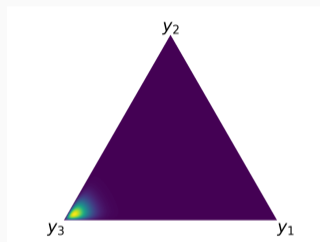
Example: 1st vs 2nd Order Risk Minimization



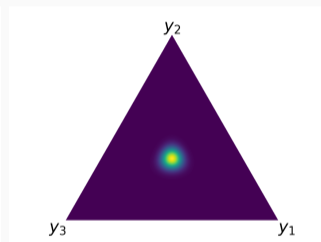
Do the resulting 2nd-order distributions represent the underlying epistemic uncertainty of the 1st-order predictor?

What is a Faithful Representation of Epistemic Uncertainty?

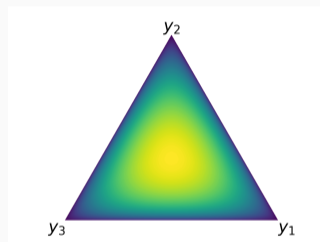
Spread of the distribution should yield a valid estimate of EU:



(a) low uncertainty



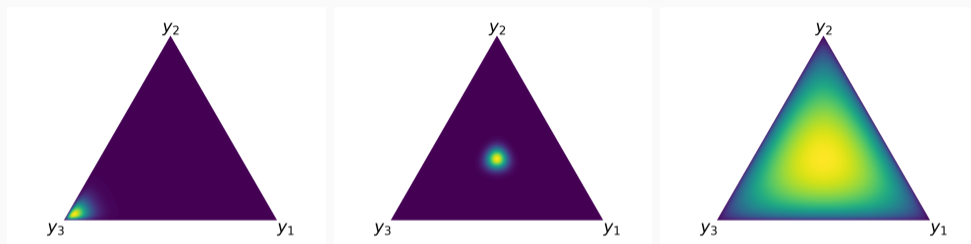
(b) high aleatoric uncertainty



(c) high epistemic uncertainty

What is a Faithful Representation of Epistemic Uncertainty?

Spread of the distribution should yield a valid estimate of EU:



(a) low uncertainty

(b) high aleatoric uncertainty

(c) high epistemic uncertainty

Desirable **convergence properties**:⁹

1. Monotonicity: decreasing uncertainty with increasing sample size N
2. Convergence to Dirac delta distribution when $N \rightarrow \infty$

⁹Bengs, Hüllermeier, and Waegeman, "Pitfalls of epistemic uncertainty quantification through loss minimisation".

What is a Faithful Representation of Epistemic Uncertainty?

How can we directly evaluate the faithfulness of the resulting 2nd-order distribution?

What is a Faithful Representation of Epistemic Uncertainty?

How can we directly evaluate the faithfulness of the resulting 2nd-order distribution?

Def. 3.1: Reference Distribution

Let $\theta_{\mathcal{D}_N}(\mathbf{x}; \Phi_{\mathcal{D}_N})$ denote the minimizer of the 1st-order minimization problem for a training set \mathcal{D}_N of size N , where $\mathcal{D}_N \sim P^N$. Define the *reference 2nd-order distribution* as

$$q_N(\theta | \mathbf{x}) := \mathbb{P}(\theta_{\mathcal{D}_N}(\mathbf{x}; \Phi_{\mathcal{D}_N}) = \theta) \quad \text{for } x \in \mathcal{X}.$$

What is a Faithful Representation of Epistemic Uncertainty?

How can we directly evaluate the faithfulness of the resulting 2nd-order distribution?

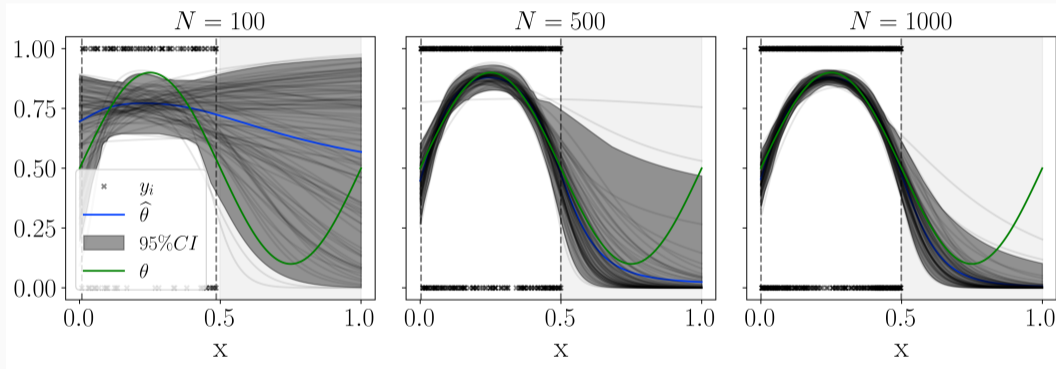
Def. 3.1: Reference Distribution

Let $\theta_{\mathcal{D}_N}(\mathbf{x}; \Phi_{\mathcal{D}_N})$ denote the minimizer of the 1st-order minimization problem for a training set \mathcal{D}_N of size N , where $\mathcal{D}_N \sim P^N$. Define the *reference 2nd-order distribution* as

$$q_N(\theta | \mathbf{x}) := \mathbb{P}(\theta_{\mathcal{D}_N}(\mathbf{x}; \Phi_{\mathcal{D}_N}) = \theta) \quad \text{for } x \in \mathcal{X}.$$

⇒ **Can be approximated empirically** by resampling and computing the empirical distribution function.

Reference Distribution: Example



Theoretical Results

Theoretical Results: Unregularized Case ($\lambda = 0$)

Theorem 3.2 and Theorem 3.1. in our paper show that

- **non-identifiability** issues arise for inner loss minimisation, leading to a wide range of possible values for the estimated uncertainty
- **convergence to a Dirac Delta distribution** in the case of outer loss minimisation (generalization of Theorem 1 of Bengs et al¹⁰ to all exponential family members)

¹⁰Bengs, Hüllermeier, and Waegeman, “Pitfalls of epistemic uncertainty quantification through loss minimisation”.

Theoretical Results: Regularized Case ($\lambda > 0$)

Theorem 3.3 shows that for inner and outer loss minimisation with entropy regularization

- there exists $\lambda \geq 0$, $x \in \mathcal{X}$ for which the 2nd-order distribution differs from the reference distribution.

Theoretical Results: Regularized Case ($\lambda > 0$)

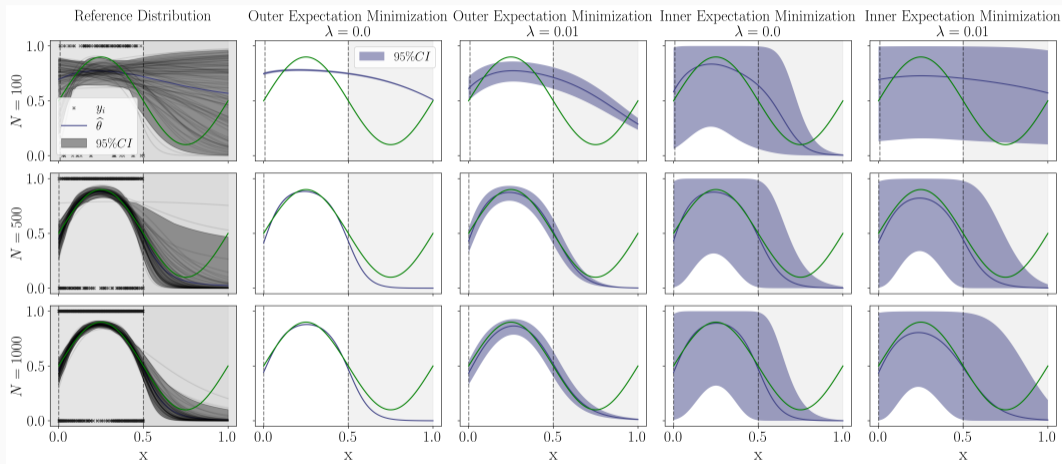
Theorem 3.3 shows that for inner and outer loss minimisation with entropy regularization

- there exists $\lambda \geq 0$, $x \in \mathcal{X}$ for which the 2nd-order distribution differs from the reference distribution.

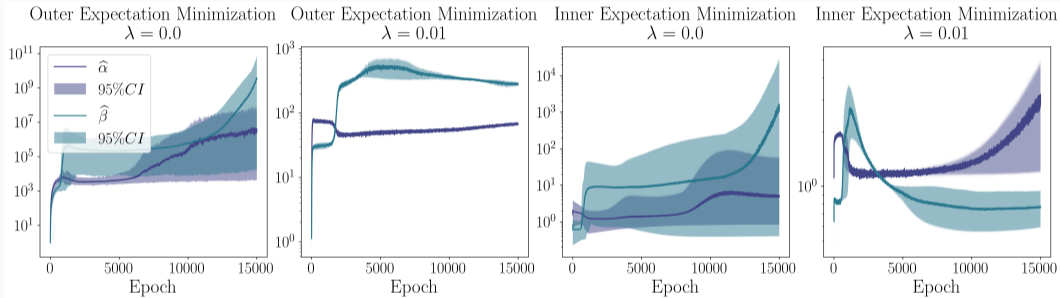
Common ways in EDL to optimize the parameter λ are mainly based on heuristics, defining an *uncertainty budget* that cannot be exceeded!

Empirical Results

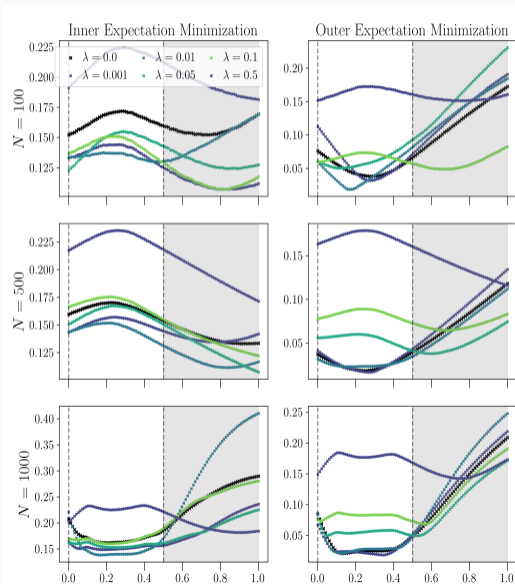
(Some) Experimental Results: Classification



(Some) Experimental Results: Convergence Analysis



(Some) Experimental Results: Distance Analysis



Conclusion

Summary

- DER does not result in distributions that are faithfully representing EU
- the regularization parameter λ yields an *uncertainty budget* that cannot be exceeded for a given amount of training data points

Outlook:

- further analysis on different regularization
 - different regularizers
 - more advanced second-order loss functions

is needed

For other exp. results, proofs, and more theoretical analysis, check out our paper.

THANK YOU

For other exp. results, proofs, and more theoretical analysis, check out our paper:



arXiv:2402.09056 [cs.AI]

(Is Epistemic Uncertainty Faithfully Represented by Evidential Deep Learning Methods?)

Empirical Results: Regression ($N = 100$)

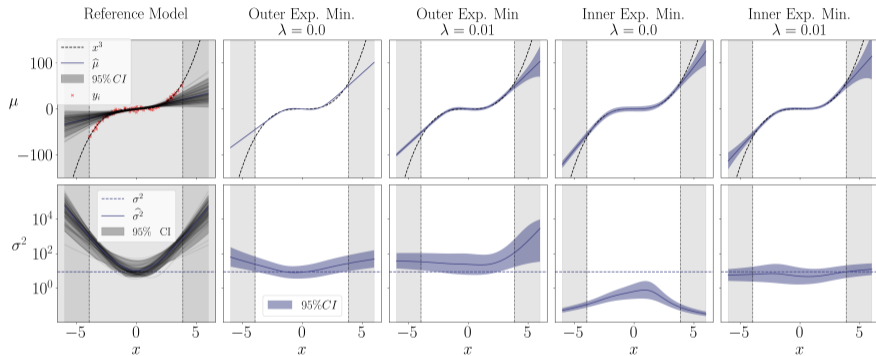
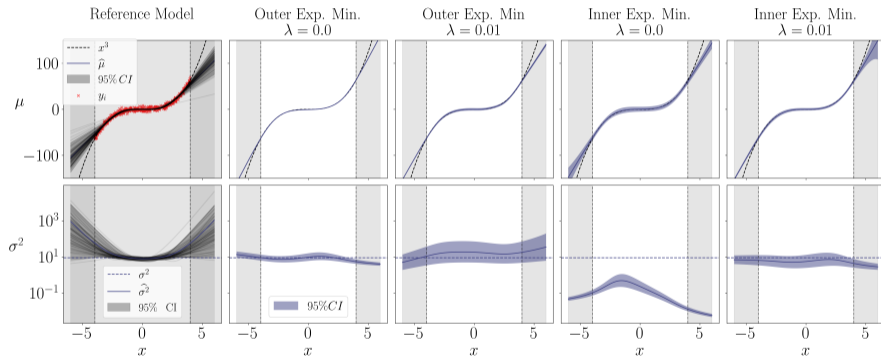


Fig. 4: Regression experiment with $\mathcal{D}_N = \{(x_i, x_i^3 + \epsilon)\}_{i=1}^N$ for $N \in \{100, 500, 1000\}$, where $x_i \in U([-4, 4])$, $\epsilon \sim \mathcal{N}(0, \sigma^2 = 9)$. The reference model learns the parameters $\theta = (\mu, \sigma)$ of the underlying normal distribution, by optimizing the negative log-likelihood.

Empirical Results: Regression ($N = 500$)



Empirical Results: Regression ($N = 1000$)

