# Uncertainty Quantification in Deep Learning

Murat Şensoy
Senior Research Scientist
Blue Prism AI Labs
London, UK

1

# Outline

- Motivation
- Methods for Uncertainty Quantification in Deep Learning
- Evidential Deep Learning
- Real-World Applications
- Conclusions

2

2

# Overconfident Mistakes of Classifiers



**Classified as:**
**Typewriter keyboard**
83.14%

**Classified as:**
**Stone wall**
87.63%

Figure 1: Predictions by EfficientNet (Tan and Le, 2019) on test images from ImageNet: For the left image, the neural network predicts "typewriter keyboard" with certainty 83.14 %, for the right image "stone wall" with certainty 87.63 %.

Tan M, Le Q (2019) EfficientNet: Rethinking model scaling for convolutional neural networks. In: Proc. ICML, 36th Int. Conference on Machine Learning, Long Beach, California

Courtesy of: Hüllermeier, Eyke, and Willem Waegeman. "Aleatoric and epistemic uncertainty in machine learning: A tutorial introduction." *arXiv preprint arXiv:1910.09457* (2019).

3

---

# Overconfident Mistakes of Classifiers

There is really but one thing to say about **this** sorry movie It should never have been made The first one one of my favourites An American Werewolf in London is a great movie with a good plot good actors and good FX But this one It stinks to heaven with a cry of helplessness

There is really but one thing to say about **that** sorry movie It should never have been made The first one one of my favourites An American Werewolf in London is a great movie with a good plot good actors and good FX But this one It stinks to heaven with a cry of helplessness

Figure 2: Adversarial example (right) misclassified by a machine learning model trained on textual data: Changing only a single — and apparently not very important — word (highlighted in bold font) is enough to turn the correct prediction "negative sentiment" into the incorrect prediction "positive sentiment" (Sato et al., 2018).

Courtesy of: Hüllermeier, Eyke, and Willem Waegeman. "Aleatoric and epistemic uncertainty in machine learning: A tutorial introduction." *arXiv preprint arXiv:1910.09457* (2019).

4

"Being able to assess the reliability of a probability score for each instance is much more powerful than assigning an aggregate reliability score [...] independent of the instance to be classified."

Kull and Flach (2014). Reliability maps: A tool to enhance probability estimates and improve classification accuracy. In: Proc. of ECML'14.



5

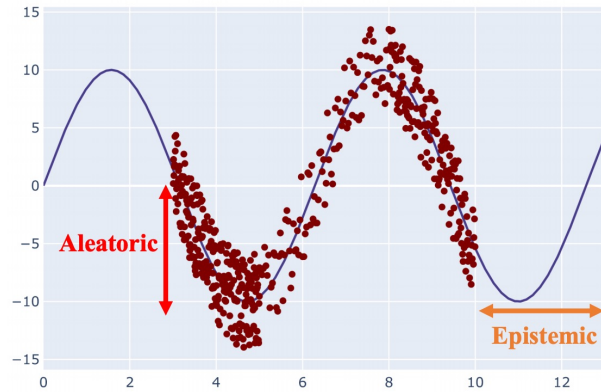# Not all mistakes are equals!




On 7th May 2016, a car operating with automated vehicle control systems crashed with a truck near Williston, Florida, USA. Unfortunately, the car driver died due to the severe injury. The car manufacturer reported that *the car's vision system classified the white side of the truck as the sky*.

6
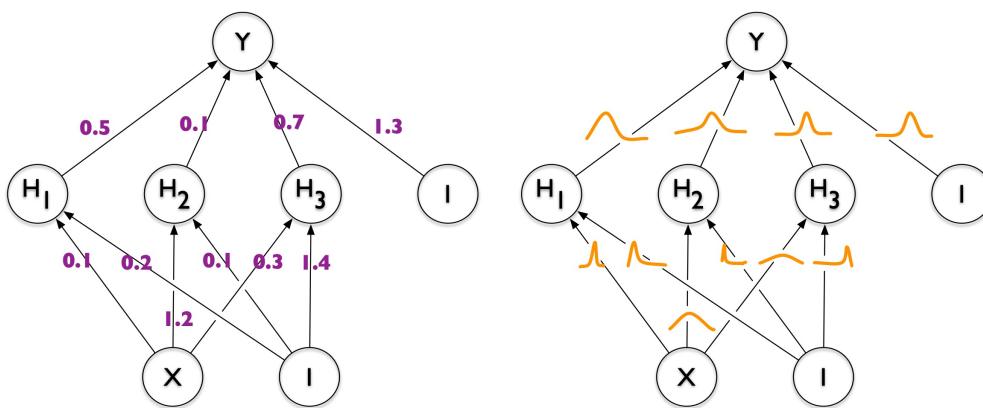
# Uncertainty Quantification in DL

- Two types of uncertainties: Aleatoric and Epistemic
- Wisdom of ignorance: knowing what you do not know
- Sampling-based methods:
  - Bayesian Networks
    - Multiplicative Normalizing Flow
    - Bayes by Backprop
  - Monte-Carlo Dropout
    - Variational dropout
  - Deep Ensembles
- Evidential Deep Learning



*** Abdar, Moloud, et al. "A review of uncertainty quantification in deep learning: Techniques, applications and challenges." *Information Fusion* (2021).

7

7

# Bayesian Neural Networks



Predictive distribution $p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}^*|\mathbf{x}^*, \boldsymbol{\omega})p(\boldsymbol{\omega}|\mathbf{X}, \mathbf{Y})\mathrm{d}\boldsymbol{\omega}$
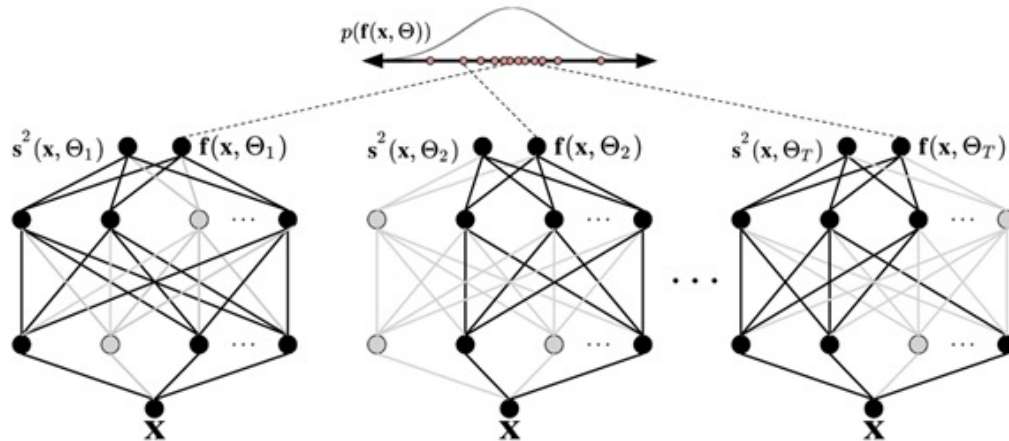
**Source:** Weight Uncertainty in Neural Networks, Blundell et al., ICML 2015

8

8

# Monte-Carlo Dropout (Bernoulli Approx. VI)



$$p(\mathbf{f}(\mathbf{x}, \Theta))$$

$$s^2(\mathbf{x}, \Theta_1) \quad \mathbf{f}(\mathbf{x}, \Theta_1) \qquad s^2(\mathbf{x}, \Theta_2) \quad \mathbf{f}(\mathbf{x}, \Theta_2) \qquad s^2(\mathbf{x}, \Theta_T) \quad \mathbf{f}(\mathbf{x}, \Theta_T)$$

Gal, Y., and Z. Ghahramani. 2016. "Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning." *33rd International Conference on Machine Learning (ICML 2016)*. https://arxiv.org/abs/1506.02142.

9

9

# Deep Ensembles

1. Let each neural network parametrize a distribution over the outputs, i.e. $p_\theta(y|\mathbf{x})$. Use a **proper scoring rule** as training criterion
   ‣ Classification: cross entropy loss
   ‣ Heteroscedastic Regression: net outputs mean $\mu_\theta(\mathbf{x})$ and variance $\sigma_\theta^2(\mathbf{x})$

$$\ell(\theta, \mathbf{x}_n, y_n) = \frac{1}{2}\log \sigma_\theta^2(\mathbf{x}) + \frac{(y - \mu_\theta(\mathbf{x}))^2}{2\sigma_\theta^2(\mathbf{x})} + \text{const.}$$

2. Augment with **adversarial training**
3. Train an **ensemble of $M$ networks in parallel** with random initialization
4. Combine predictions at test time

$$p(y|\mathbf{x}) = \frac{1}{M}\sum_{m}^{M} p_{\theta_m}(y|\mathbf{x}, \theta_m)$$

**Source**: https://www.gatsby.ucl.ac.uk/~balaji/deep-ensembles-poster.pdf

Lakshminarayanan, B., A. Pritzel, and C. Blundell. 2017. "Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles." Advances in Neural Information Processing Systems. https://arxiv.org/abs/1612.01474.

10

10

# Evidential Deep Learning

Murat Sensoy, Lance Kaplan, and Melih Kandemir, "Evidential deep learning to quantify classification uncertainty." *Advances in Neural Information Processing Systems (NuerIPS)*. pp. 3179-3189 , 2018.

Murat Sensoy, Lance Kaplan, Federico Cerutti, "Uncertainty-Aware Deep Classifiers using Generative Models", The 34rd Conference on Artificial Intelligence (AAAI), 2020.

Murat Sensoy, Maryam Saleki, Simon Julier, John Reid. "Not all Mistakes are Equal." Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems (AAMAS). 2020.

Murat Sensoy, Maryam Saleki, Simon Julier, John Reid. "Misclassification Risk and Uncertainty Quantification in Deep Classifiers." *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*. 2021.
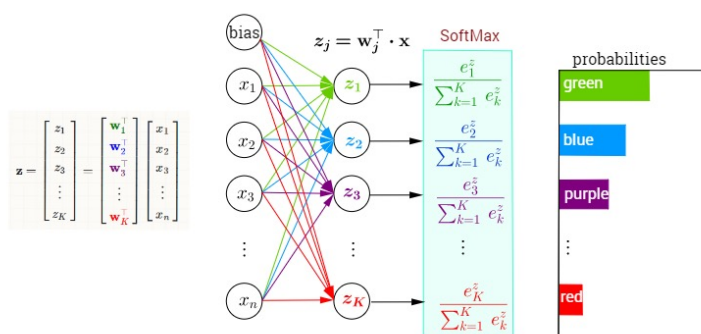
11

11

# Standard Approach for Classification

Usually **SoftMax** function is used to estimate class probabilities based on the output of a Deep Neural Network.

**Multi-Class Classification with NN and SoftMax Function**



**SoftMax Function**

$$\sigma(x_j) = \frac{e^{x_j}}{\sum_i e^{x_i}}$$

SoftMax leads to neural networks that are over-confident when they encounter samples which are highly different from examples in the training set.
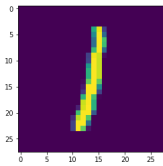
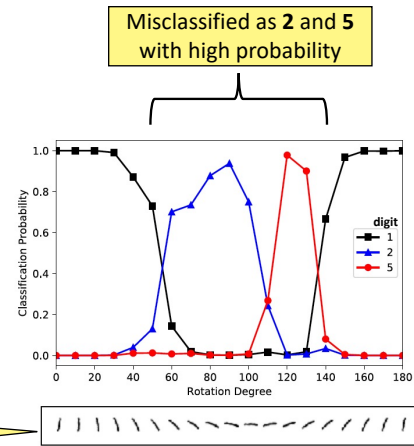*Lets check the example in the following slide*

12

12

# Predicting Class Labels with Softmax and Cross-Entropy

Consider the following image of digit 1 from the well-known MNIST dataset.

A neural network trained on the MNIST dataset can easily classify this image as the digit 1.

Misclassified as **2** and **5** with high probability

What happens if the image is rotated counter-clockwise? For some angles, it looks different from the images of digits

13

13

# Subjective Logic

Subjective logic is a calculus for subjective opinions, which in turn represent probabilities affected by degrees of uncertainty.

opinion owner

binary state variable

**b**elief + **d**isbelief + **u**ncertainty = 1

evidence for $x$    evidence for $\neg x$

$$\omega_x^s = (b, d, u) = Beta(\frac{2b}{u} + 1, \frac{2d}{u} + 1)$$

**Is the following true?**

$$digit(\ ,9)$$

Victor

**opinion**
(0.1, 0.8, 0.1)
**Most likely NOT**

Beta(3, 17)

Bob

**opinion**
(0.0, 0.0, 1.0)
**I do NOT know**

Beta(1, 1)

14

14

7

# Dirichlet Distribution



(a) $\alpha = [1, 1, 1]$  (b) $\alpha = [2, 5, 15]$  (c) $\alpha = [10, 10, 10]$

At the top, density plots (blue = low, yellow = high) for the Dirichlet distributions over the probability simplex in $\mathbb{R}^3$ for various values of the $\alpha$ parameters and, at the bottom, 500 categorical distributions sampled from each of these Dirichlet distributions.

The Dirichlet distribution is the conjugate prior of the categorical and multinomial distributions. It is a probability density function (pdf) for possible values of the probability mass function (pmf) $\pi = [\pi_1, \dots, \pi_K]$ over $K$ categories. It is characterized by parameters $\alpha = [\alpha_1, \cdots, \alpha_K]$ and is given by

$$\text{Dirichlet}(\pi | \alpha) = \begin{cases} \frac{1}{B(\alpha)} \prod_{i=1}^{K} \pi_i^{\alpha_i - 1} & \text{for } \pi \in \mathcal{S}_K, \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

where $\mathcal{S}_K$ is the $K$-dimensional unit simplex and $B(\alpha)$ is the $K$-dimensional multinomial beta function

$$b_i = \frac{\alpha_i - 1}{\sum_{k=1}^{K} \alpha_k} \qquad \pi_i = \frac{\alpha_i}{\sum_{k=1}^{K} \alpha_k}$$
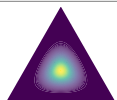
$$u = \frac{K}{\sum_{k=1}^{K} \alpha_k} \qquad \boxed{\text{Subjective Logic Interpretation}}$$

15

---

# Why is this useful?

Assume you have trained a classifier to distinguish *Cat*, *Fossa*, and *Fox* images



spoiler: this is from a fox cub picture

$c = [9, 9, 9]$
$\alpha = [10, 10, 10]$
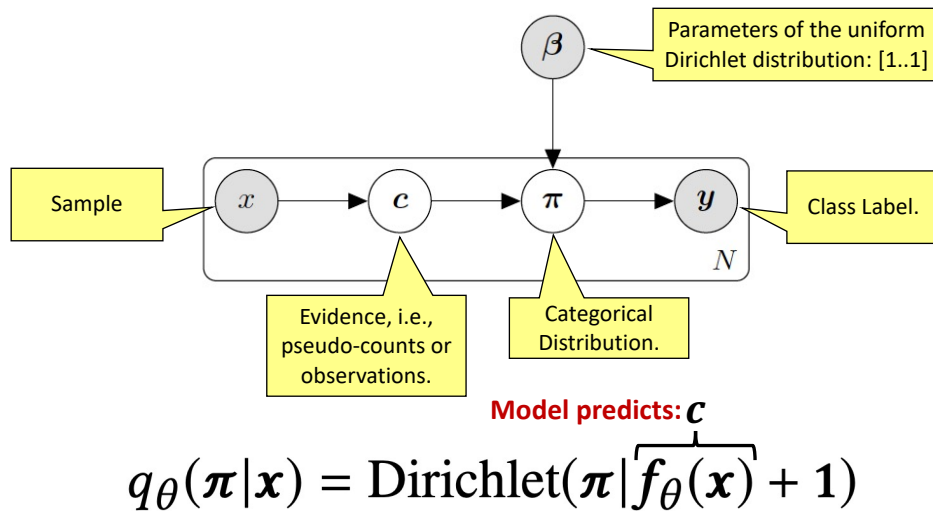
Aleatoric Uncertainty        Epistemic Uncertainty

$c = [0, 0, 0]$
$\alpha = [1, 1, 1]$

$\pi = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ ← Uniform Categorical Distribution → $\pi = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$

# Evidential Deep Classifiers



Parameters of the uniform Dirichlet distribution: [1..1]

Sample

Class Label.

Evidence, i.e., pseudo-counts or observations.

Categorical Distribution.

Model predicts: $c$

$$q_\theta(\pi|x) = \mathrm{Dirichlet}(\pi|\overbrace{f_\theta(x)} + 1)$$

17

17

# EDL loss functions

- For each input $x_i$, we define a *base loss function* parametrized by the latent categorical distribution $\pi_i$ and compute its expectation using the predicted Dirichlet distribution: $Dirichlet(\pi_i|f_\theta(x_i) + 1)$.

- **Expected cross-entropy loss:**

$$\mathcal{L}_i(\theta) = \int \underbrace{\left[\sum_{j=1}^{K} -y_{ij}\log(\pi_{ij})\right]}_{\text{cross-entropy}} \underbrace{\frac{1}{B(\boldsymbol{\alpha}_i)}\prod_{j=1}^{K}\pi_{ij}^{\alpha_{ij}-1}}_{\text{Dirichlet}(\boldsymbol{\pi_i}|\boldsymbol{f_\theta}(\boldsymbol{x_i})+\mathbf{1})}\, d\boldsymbol{\pi}_i = \sum_{j=1}^{K} y_{ij}\left(\psi(\sum_{k=1}^{K}\alpha_{ik}) - \psi(\alpha_{ij})\right)$$

$\psi$ is the digamma function, i.e., the derivative of the log gamma function.
$y_i$ is the one-hot label vector for the sample $x_i$.

18

18

# EDL loss functions

**Type II maximum likelihood:**

We can treat $Dirichlet(\boldsymbol{\pi}_i | \boldsymbol{f}_{\boldsymbol{\theta}}(\boldsymbol{x}_i) + \mathbf{1})$ as a prior on the likelihood $Mult(\boldsymbol{y}_i | \boldsymbol{\pi}_i)$ and obtain the negated logarithm of the marginal likelihood by integrating out the class probabilities.

$$\mathcal{L}_i(\theta) = -\log\left(\int \underbrace{\left[\prod_{j=1}^{K} \pi_{ij}^{y_{ij}}\right]}_{Mult(\boldsymbol{y}_i|\boldsymbol{\pi}_i)} \underbrace{\frac{1}{B(\boldsymbol{\alpha}_i)} \prod_{j=1}^{K} \pi_{ij}^{\alpha_{ij}-1}}_{\text{Dirichlet}(\boldsymbol{\pi_i}|\boldsymbol{f_\theta}(\boldsymbol{x_i})+\mathbf{1})} d\boldsymbol{\pi}_i\right) = \sum_{j=1}^{K} y_{ij}\left(\log(\sum_{k=1}^{K} \alpha_{ik}) - \log(\alpha_{ij})\right)$$

19

19

# EDL loss functions

**The expected sum square error loss (Brier score):**

$$\mathcal{L}_i(\theta) = \int \underbrace{||\boldsymbol{y}_i - \boldsymbol{\pi}_i||_2^2}_{\text{SSE loss}} \underbrace{\frac{1}{B(\boldsymbol{\alpha}_i)} \prod_{j=1}^{K} \pi_{ij}^{\alpha_{ij}-1}}_{\text{Dirichlet}(\boldsymbol{\pi_i}|\boldsymbol{f_\theta}(\boldsymbol{x_i})+\mathbf{1})} d\boldsymbol{\pi}_i$$

$$= \sum_{j=1}^{K} \mathbb{E}\left[y_{ij}^2 - 2y_{ij}\pi_{ij} + \pi_{ij}^2\right] = \sum_{j=1}^{K} \left(y_{ij}^2 - 2y_{ij}\mathbb{E}[\pi_{ij}] + \mathbb{E}[\pi_{ij}^2]\right)$$

20

20

An advantage of this loss is that using the identity

$$\mathbb{E}[\pi_{ij}^2] = \mathbb{E}[\pi_{ij}]^2 + \mathrm{Var}(\pi_{ij}),$$

we get the following easily interpretable form

$$\mathcal{L}_i(\theta) = \sum_{j=1}^K (y_{ij} - \mathbb{E}[\pi_{ij}])^2 + \mathrm{Var}(\pi_{ij}) = \sum_{j=1}^K \underbrace{(y_{ij} - \bar{\pi}_{ij})^2}_{\mathcal{L}_{ij}^{err}} + \underbrace{\frac{\bar{\pi}_{ij}(1 - \bar{\pi}_{ij})}{(1 + \sum_k \alpha_{ik})}}_{\mathcal{L}_{ij}^{var}}.$$

▶ Proposition 1. For any $\alpha_{ij} \geq 1$, the inequality $\mathcal{L}_{ij}^{var} < \mathcal{L}_{ij}^{err}$ is satisfied.
i.e. The loss prioritizes data fit over variance estimation.

▶ Proposition 2. For a given sample $i$ with the correct label $j$, $L_i^{err}$ decreases when new evidence is added to $\alpha_{ij}$ and increases when evidence is removed from $\alpha_{ij}$.
i.e. The loss has a tendency to fit to the data.

▶ Proposition 3. For a given sample $i$ with the correct class label $j$, $L_i^{err}$ decreases when some evidence is removed from the biggest Dirichlet parameter $\alpha_{il}$ such that $l \neq j$.
i.e. The loss performs learned loss attenuation.

21

21

# Overall loss function

Evidence for the correct class is removed.

$$\mathcal{L}(\theta) = \sum_{i=1}^N \mathcal{L}_i(\theta) + \lambda_t \sum_{i=1}^N \mathrm{KL}[\mathrm{Dirichlet}(\boldsymbol{\pi}_i|\tilde{\boldsymbol{\alpha}}_i)\|\mathrm{Dirichlet}(\boldsymbol{\pi}_i|\mathbf{1})]$$

**Maximize model fit**          **Minimize evidence on errors.**

$$\mathrm{KL}[\mathrm{Dirichlet}(\boldsymbol{\pi}_i\tilde{\boldsymbol{\alpha}}_i)\|\mathrm{Dirichlet}(\boldsymbol{\pi}_i|\mathbf{1})]$$

$$= \log\left(\frac{\Gamma(\sum_{k=1}^K \tilde{\alpha}_{ik})}{\Gamma(K)\prod_{k=1}^K \Gamma(\tilde{\alpha}_{ik})}\right) + \sum_{k=1}^K (\tilde{\alpha}_{ik} - 1)\left[\psi(\tilde{\alpha}_{ik}) - \psi\left(\sum_{j=1}^K \tilde{\alpha}_{ij}\right)\right]$$

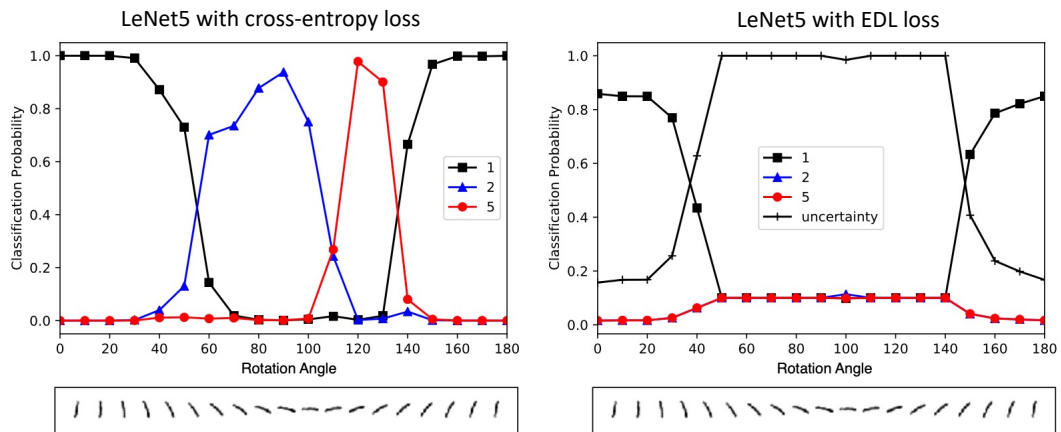$\lambda_t$ is the annealing coefficient; initially 0 and increased gradually to and 1 during training.
$\tilde{\alpha}$ refers to the predicted Dirichlet parameters after removing the evidence for the true category.

22

22

11

## Revisiting MNIST



LeNet5 with cross-entropy loss

LeNet5 with EDL loss

**The demo is available at https://muratsensoy.github.io/uncertainty.html**

23

23

## Evaluations

Recent research on uncertainty quantification is centred around Bayesian Neural Networks using
- LeNet5* architecture with ReLU nonlinearities
- MNIST and CIFAR10 datasets

and evaluated on tasks
- Out of distribution detection
- Adversarial Robustness

*LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, *86*(11), 2278-2324.

24

24

# Evaluations

We use the following approaches in our evaluations:

- Vanilla Neural Nets with weight decay (L2)
- Multiplicative Normalizing Flow (MNFG). Louizos and Welling, ICML, 2017.
- Deep Ensembles. Lakshminarayanan et al., NIPS, 2017.
- Monte Carlo Dropout. Gal and Ghahramani, ICML, 2016.
- Variational dropout and the local reparameterization trick (FFLU). Kingma et al. NIPS, 2015
- Bayes by Backprop (FFG). Blundell, ICML, 2015.
- Evidential Deep Learning (EDL). Sensoy et al., NIPS, 2018.

25

25



**MNIST Dataset**

Training Data

**CIFAR5 Dataset**

**notMNIST Dataset**

Testing Data

**CIFAR10 Dataset**

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

26

26

## Accuracy vs Uncertainty

First 5 categories of CIFAR10



| Method | MNIST | CIFAR 5 |
|---|---|---|
| *L2* | 99.4 | 76 |
| *Dropout* | 99.5 | 84 |
| *Deep Ensemble* | 99.3 | 79 |
| *FFG* | 99.1 | 78 |
| *FFLU* | 99.1 | 77 |
| *MNFG* | 99.3 | 84 |
| *EDL* | 99.3 | 83 |

Figure 2: The change of accuracy with respect to the uncertainty threshold for *EDL*.

Table 1: Test accuracies (%) for MNIST and CIFAR5 datasets.

27

27

## Out-of-Distribution Samples

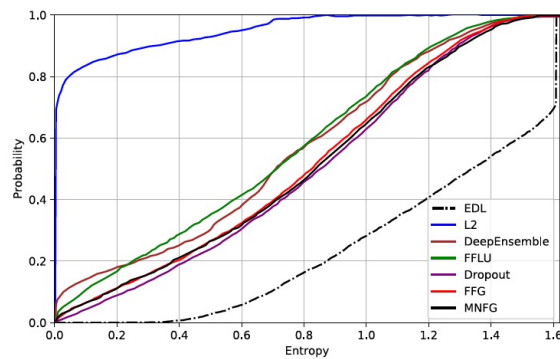*notMNIST dataset*

*Last 5 categories from CIFAR10*



Figure 3: Empirical CDF for the entropy of the predictive distributions on the notMNIST dataset (left) and samples from the last five categories of CIFAR10 dataset (right).

28

28

14

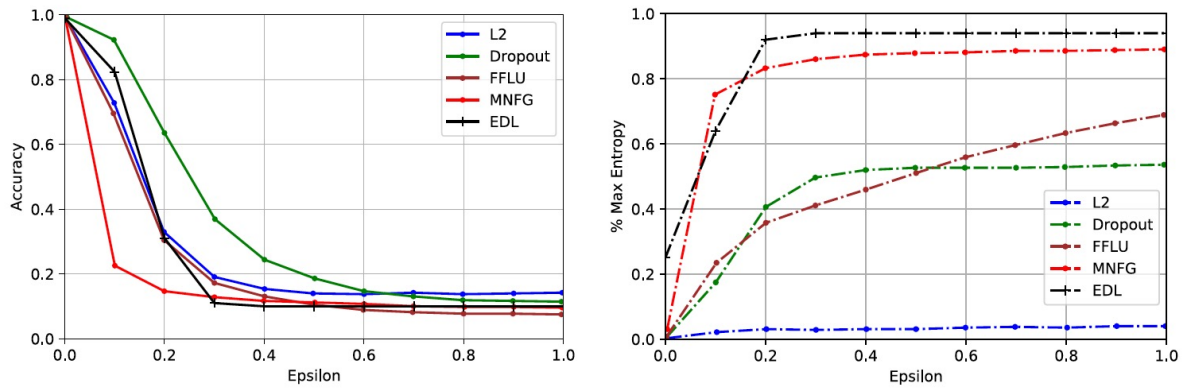## Adversarial Robustness (against FGSM)



Figure 4: Accuracy and entropy as a function of the adversarial perturbation $\epsilon$ on the MNIST dataset.

29

29

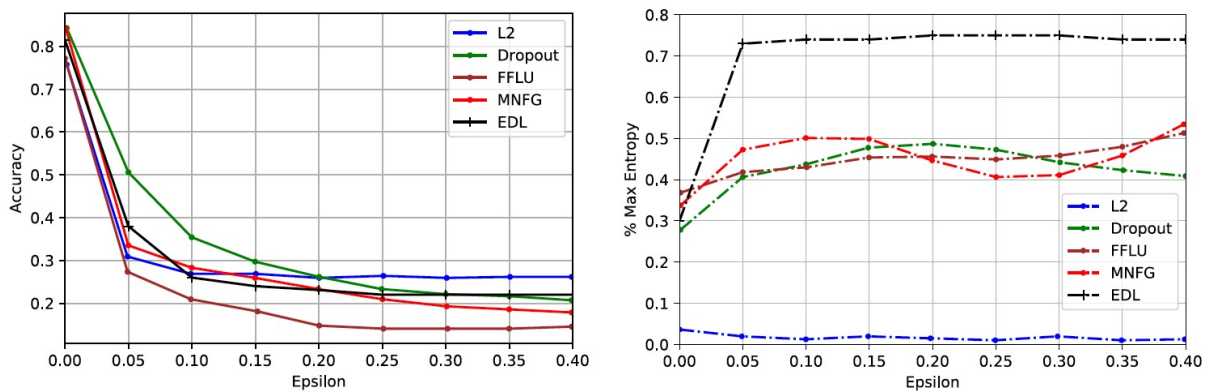## Adversarial Robustness (against FGSM)



Figure 5: Accuracy and entropy as a function of the adversarial perturbation $\epsilon$ on CIFAR5 dataset.

30

30

15

---

# Some Real-World Examples
# using
# Evidential Deep Learning

31

31

---

## Deep, spatially coherent Inverse Sensor Models with Uncertainty Incorporation using the evidential Framework

[1]Daniel Bauer and Lars Kuhnert are with the Ford Werke GmbH, Cologne, dbauer31@ford.de, lkuhnert@ford.de
[3]Lutz Eckstein is with the Institute for Automotive Engineering, RWTH Aachen University, office@ika.rwth-aachen.de

*Abstract*— To perform high speed tasks, sensors of autonomous cars have to provide as much information in as few time steps as possible. However, radars, one of the sensor modalities autonomous cars heavily rely on, often only provide sparse, noisy detections. These have to be accumulated over time to reach a high enough confidence about the static parts of the environment. For radars, the state is typically estimated by accumulating inverse detection models (IDMs). We employ the recently proposed evidential convolutional neural networks which, in contrast to IDMs, compute dense, spatially coherent inference of the environment state. Moreover, these networks are able to incorporate sensor noise in a principled way which we further extend to also incorporate model uncertainty. We present experimental results that show This makes it possible to obtain a denser environment perception in fewer time steps.

**U-Net** architecture is used in this study.

D. Bauer, L. Kuhnert and L. Eckstein, "Deep, spatially coherent Inverse Sensor Models with Uncertainty Incorporation using the evidential Framework," *2019 IEEE Intelligent Vehicles Symposium (IV)*, Paris, France, 2019, pp. 2490-2495, doi: 10.1109/IVS.2019.8813826.

32

32

## Slide 33

### Quantifying and Leveraging Classification Uncertainty for Chest Radiograph Assessment

Florin C. Ghesu[1], Bogdan Georgescu[1], Eli Gibson[1], Sebastian Guendel[1],
Mannudeep K. Kalra[2,3], Ramandeep Singh[2,3], Subba R. Digumarthy[2,3],
Sasa Grbic[1], and Dorin Comaniciu[1]

[1] Digital Technology and Innovation, Siemens Healthineers, Princeton, NJ, USA
[2] Department of Radiology, Massachusetts General Hospital, Boston, MA, USA
[3] Harvard Medical School, Boston, MA, USA
florin.ghesu@siemens-healthineers.com

Table 1: Comparison between the reference method [4] and several versions of our method calibrated at sample rejection rates of 0%, 10%, 25% and 50% (based on the PLCO dataset [2]). Lesion refers to lesions of the bones or soft tissue.

**ROC-AUC**

| Finding | Guendel et al. [4] | Ours [0%] | Ours [10%] | Ours [25%] | Ours [50%] |
|---|---|---|---|---|---|
| Granuloma | 0.83 | 0.85 | 0.87 | **0.90** | **0.92** |
| Fibrosis | 0.87 | 0.88 | 0.90 | **0.92** | **0.94** |
| Scaring | 0.82 | 0.81 | 0.84 | **0.89** | **0.93** |
| Lesion | 0.82 | 0.83 | 0.86 | **0.88** | **0.90** |
| Cardiac Ab. | 0.93 | 0.94 | 0.95 | **0.96** | **0.97** |
| **Average** | 0.85 | 0.86 | 0.89 | **0.91** | **0.93** |

**DenseNet121 is used in this study.**



Fig. 2: Evolution of the F1-scores for the positive (+) and negative (−) classes relative to the sample rejection threshold - determined using the estimated uncertainty. We show the performance for granuloma and fibrosis based on the PLCO dataset [2]. The baseline (horizontal dashed lines) is determined using the method from [4] (working point at max. average of per-class F1 scores). Decision threshold for our method is fixed at 0.5.

In: Shen D. et al. (eds) Medical Image Computing and Computer Assisted Intervention – MICCAI 2019. MICCAI 2019. Lecture Notes in Computer Science, vol 11769. Springer, Cham
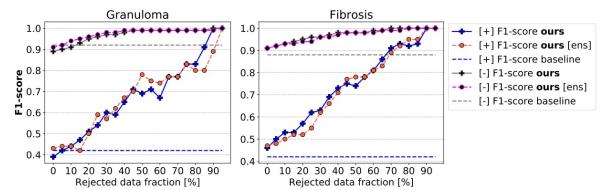
33

## Slide 34

### Quantifying and Leveraging Classification Uncertainty for Chest Radiograph Assessment

Florin C. Ghesu[1], Bogdan Georgescu[1], Eli Gibson[1], Sebastian Guendel[1],
Mannudeep K. Kalra[2,3], Ramandeep Singh[2,3], Subba R. Digumarthy[2,3],
Sasa Grbic[1], and Dorin Comaniciu[1]

[1] Digital Technology and Innovation, Siemens Healthineers, Princeton, NJ, USA
[2] Department of Radiology, Massachusetts General Hospital, Boston, MA, USA
[3] Harvard Medical School, Boston, MA, USA
florin.ghesu@siemens-healthineers.com

Labels in ChestX-Ray8 dataset are automatically generated by an NLP system from radiographic reports. A committee of 4 board certified experts analysed randomly selected test set of **689** cases; relabelled **120** samples are called *critical set*.



(a) $\hat{u}, \hat{p} = 0.90, 0.45$ (b) $\hat{u}, \hat{p} = 0.93, 0.48$ (c) $\hat{u}, \hat{p} = 0.54, 0.65$ (d) $\hat{u}, \hat{p} = 0.11, 0.05$
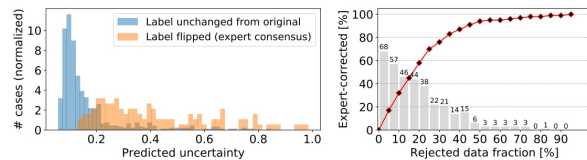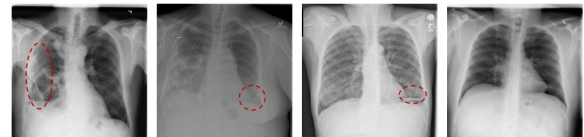
Fig. 3: **Left**: Predictive uncertainty distribution on 689 ChestX-Ray test images; a higher uncertainty is associated with cases of the critical set, which required a label correction according to expert committee. **Right**: Plot showing the capacity of the algorithm to eliminate cases from the critical set via sample rejection. Bars indicate the percentage of critical cases for each batch of 5% rejected cases.

Fig. 4: ChestX-Ray8 test images assessed for pleural effusion ($\hat{u}$: est. uncertainty, $\hat{p}$: output probability; with affected regions circled in red). Figures 4a, 4b and 4c show positive cases of the critical set with high predictive uncertainty – possibly explained by the atypical appearance of accumulated fluid in 4a, and poor quality of image 4b. Figure 4d shows a high confidence case with no pleural effusion.

In: Shen D. et al. (eds) Medical Image Computing and Computer Assisted Intervention – MICCAI 2019. MICCAI 2019. Lecture Notes in Computer Science, vol 11769. Springer, Cham

34

## Quantifying and Leveraging Classification Uncertainty for Chest Radiograph Assessment

Florin C. Ghesu[1], Bogdan Georgescu[1], Eli Gibson[1], Sebastian Guendel[1],
Mannudeep K. Kalra[2,3], Ramandeep Singh[2,3], Subba R. Digumarthy[2,3],
Sasa Grbic[1], and Dorin Comaniciu[1]

[1] Digital Technology and Innovation, Siemens Healthineers, Princeton, NJ, USA
[2] Department of Radiology, Massachusetts General Hospital, Boston, MA, USA
[3] Harvard Medical School, Boston, MA, USA
florin.ghesu@siemens-healthineers.com

**Uncertainty-driven Bootstrapping:** We also investigated the benefit of using bootstrapping based on the uncertainty measure on the example of plural effusion (ChestX-Ray8). We report performance as $[AUC; F1\text{-}score$ (pos. class); $F1\text{-}score$ (neg. class)]. After training our method, the baseline performance was measured at $[0.89; 0.60; 0.92]$ on testing. We then eliminated 5%, 10% and 15% of training samples with highest uncertainty, and retrained in each case on the remaining data. The metrics improved to $[0.90; 0.68; 0.92]_{5\%}$, $[0.91; 0.67; 0.94]_{10\%}$ and $[\mathbf{0.93}; \mathbf{0.69}; \mathbf{0.94}]_{15\%}$ on the test set. This is a significant increase, demonstrating the potential of this strategy to improve the robustness of the model to the label noise. We are currently focused on further exploring this method.

In: Shen D. et al. (eds) Medical Image Computing and Computer Assisted Intervention – MICCAI 2019. MICCAI 2019. Lecture Notes in Computer Science, vol 11769. Springer, Cham

35

35

---

# MedSpecSearch: Medical Specialty Search



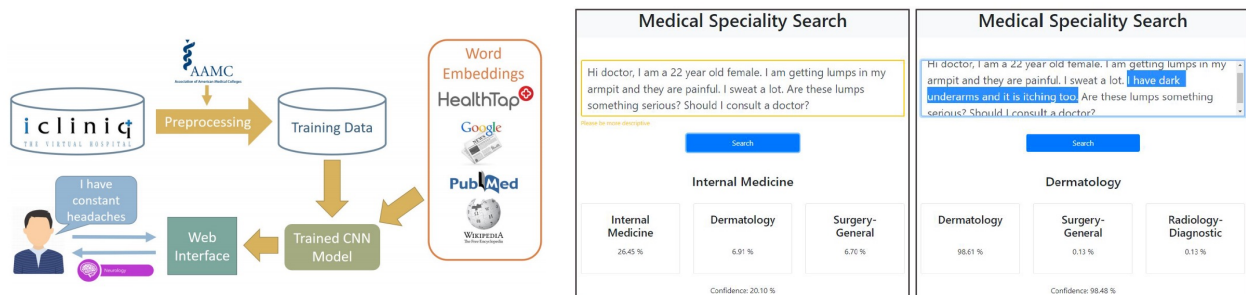**Fig. 2.** MedSpecSearch Front End with Two Example Queries

Using a confidence threshold like 90%, significantly reduces the amount of misclassifications of the system and increases the general prediction accuracy from 74% to 90.4%.

Şahin, M. U., Balatkan, E., Eran, C., Zeydan, E., & Yeniterzi, R. (2019, April). MedSpecSearch: Medical Specialty Search. In *European Conference on Information Retrieval* (pp. 225-229). Springer, Cham.
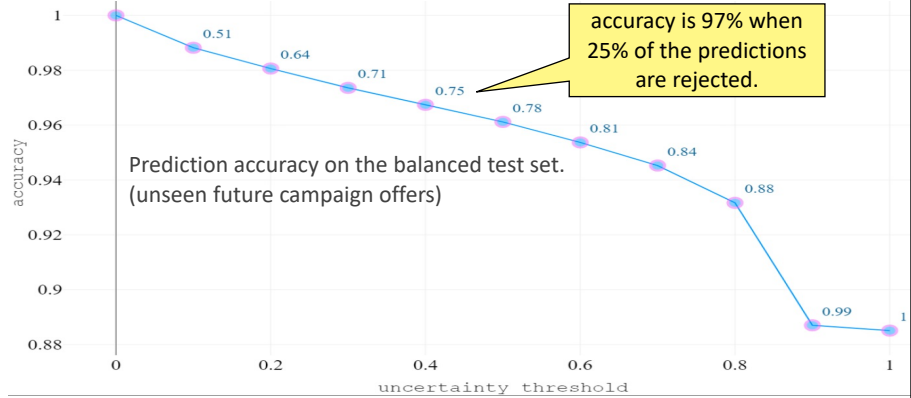
36

36

# Campaign Participation Prediction for GSM Users

- Campaign offers are sent to users as short text messages.

- Binary classification is used to predict user participation based on Google's *Wide&Deep* model.

Ayvaz, D., Aydoğan, R., Akçura, M. T., & Şensoy, M. (2021). Campaign participation prediction with deep learning. *Electronic Commerce Research and Applications*, *48*, 101058.

**An example offer sent by the GSM operator:**

"Get 500 minutes and 500MB only for 21$/month. You can join this campaign by texting 'MONTHLY500MIN' to 1111.".

accuracy is 97% when 25% of the predictions are rejected.

Prediction accuracy on the balanced test set. (unseen future campaign offers)

37

37

---

**EDL is extended by researchers in MIT for regression tasks.**

Evidential Deep Learning

Alexander Amini

MIT 6.S191

January 26, 2021

6.S191 Introduction to Deep Learning
introtodeeplearning.com  @MITDeepLearning

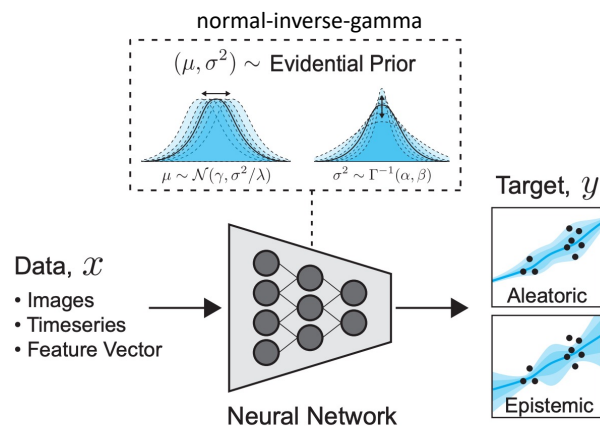MIT EECS

IntroToDeepLearning.com

MIT 6.S191: Evidential Deep Learning and Uncertainty

18,072 views • Premiered Mar 19, 2021

402    1    SHARE    SAVE    ...

38

38

# Deep Evidential Regression

**Source:** Alexander Amini, Wilko Schwarting, Ava Soleimany, Daniela Rus: Deep Evidential Regression. NeurIPS 2020



Figure 1: **Evidential regression** simultaneously learns a continuous target along with aleatoric (data) and epistemic (model) uncertainty. Given an input, the network is trained to predict the parameters of an evidential distribution, which models a higher-order probability distribution over the individual likelihood parameters, $(\mu, \sigma^2)$.

39

39

# Conclusions

- Overconfidence of deep learning models stands as an important problem and require intensive research for the applicability of these models to real-world problems.

- Existing research focus on Bayesian approaches and sampling-based methods, while some promising work, such as EDL, can achieve calibrated uncertainties without using costly sampling methods.

40

40

# Thank you!

41