



Uncertainty estimation via Prior Networks: A Tutorial

Andrey Malinin

e-mail - am969@yandex-team.ru

twitter - [@AndreyMalinin](https://twitter.com/AndreyMalinin)

22nd January 2021

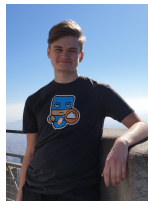
Joint work with...



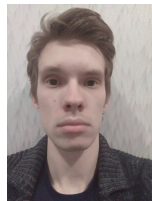
(a) Mark Gales



(b) Bruno Mlodozieniec



(c) Ivan
Provilkov



(d) Sergey
Chervontsev

Overview of the Talk

1. Motivation and Sources
2. Uncertainty Estimation via Ensembles
3. Uncertainty Estimation via Prior Networks
4. Ensemble Distribution Distillation

Overview of the Talk

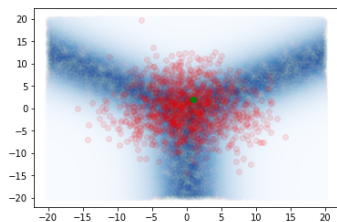
1. **Motivation and Sources**
2. Uncertainty Estimation via Ensembles
3. Uncertainty Estimation via Prior Networks
4. Ensemble Distribution Distillation

Why is Uncertainty important in practice?

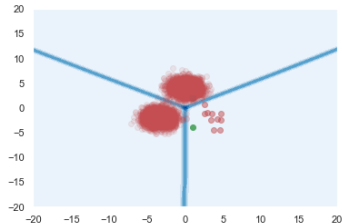
- Machine Learning (ML) systems are being deployed to many applications →
 - Image Classification, Speech Recognition, Machine Translation, etc...
- In some applications, the cost of a mistake is **high** or consequence **fatal** →
 - Medical applications, Financial applications and Autonomous vehicles
- Obtaining measures of uncertainty in predictions helps **avoid mistakes!**
 - Increases **safety** and **reliability** of ML system

- Given a **deployed** model and a **test input** \mathbf{x}^* we wish to:
 - Obtain a **prediction**
 - Obtain a measure of **uncertainty in prediction**
- Take **action** based estimate of uncertainty
 - Reject prediction / stop decoding sentence
 - Ask for human intervention
 - Use active learning
- Appropriate action depends on **source of uncertainty**

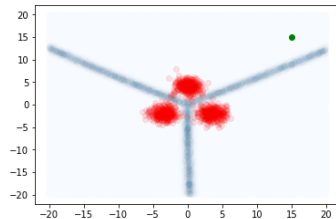
Sources of Uncertainty



(a) Data Uncertainty



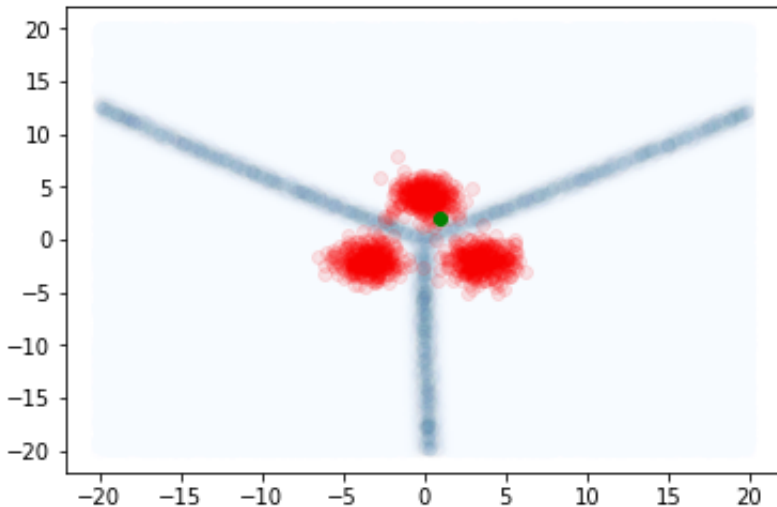
(b) Data Sparsity



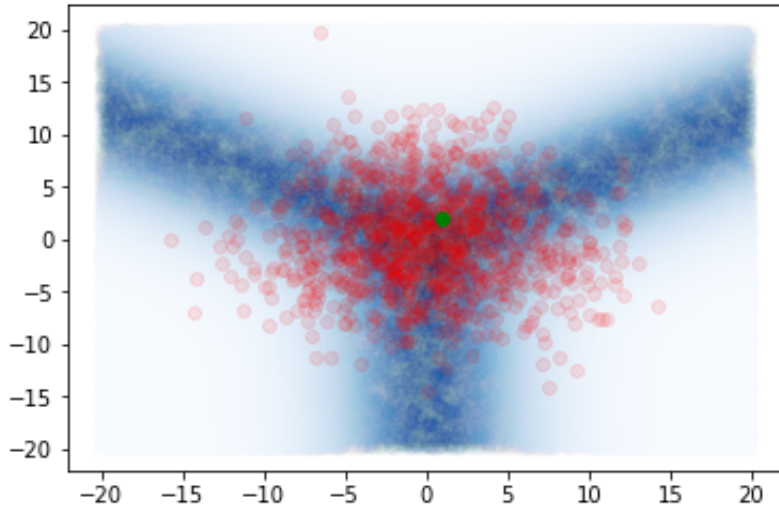
(c) Out-of-Distribution inputs

- Knowledge (epistemic) uncertainty refers to both:
 - Data Sparsity and Out-of-distribution inputs

Data (Aleatoric) Uncertainty



Data Uncertainty



Data Uncertainty

- Distinct Classes

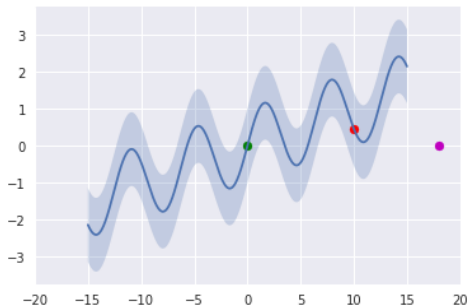
1 2 7

- Overlapping Classes

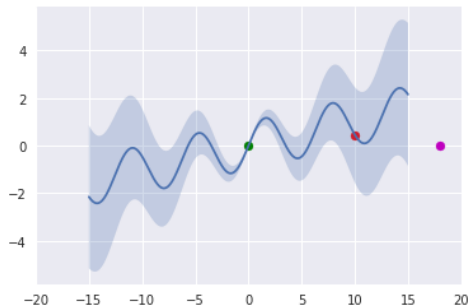
1 7 7

Data Uncertainty

- In regression tasks data uncertainty takes the form of additive noise



(a) Homoscedastic Noise



(b) Heteroscedastic Noise

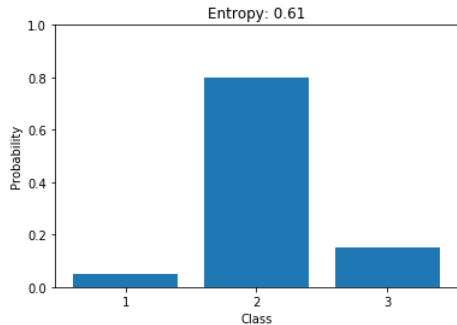
- Data Uncertainty is the *entropy* of the *true data distribution* \rightarrow

$$\mathcal{H}[\mathbf{P}_{\text{tr}}(y|\mathbf{x}^*)] = - \sum_{c=1}^K \mathbf{P}_{\text{tr}}(y = \omega_c|\mathbf{x}^*) \ln \mathbf{P}_{\text{tr}}(y = \omega_c|\mathbf{x}^*)$$

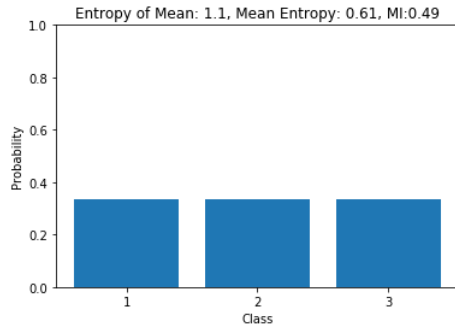
- Captured by the entropy of a model's posterior over classes \rightarrow

$$\mathcal{H}[\mathbf{P}(y|\mathbf{x}^*, \hat{\boldsymbol{\theta}})] = - \sum_{c=1}^K \mathbf{P}(y = \omega_c|\mathbf{x}^*, \hat{\boldsymbol{\theta}}) \ln \mathbf{P}(y = \omega_c|\mathbf{x}^*, \hat{\boldsymbol{\theta}})$$

Reminder - Entropy



(a) Low Entropy



(b) High Entropy

- Data Uncertainty is the *differential entropy* of the *true data distribution* \rightarrow

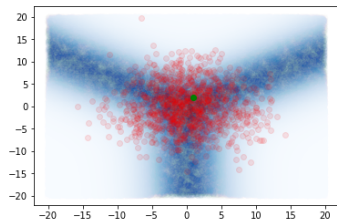
$$\mathcal{H}[p_{\text{tr}}(\mathbf{y}|\mathbf{x}^*)] = - \int p_{\text{tr}}(\mathbf{y}|\mathbf{x}^*) \ln p_{\text{tr}}(\mathbf{y}|\mathbf{x}^*) d\mathbf{y}$$

- Captured by the entropy of a model's posterior over classes \rightarrow

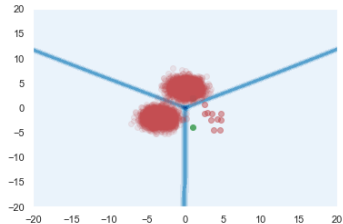
$$\mathcal{H}[p(\mathbf{y}|\mathbf{x}^*, \hat{\boldsymbol{\theta}})] = - \int p(\mathbf{y}|\mathbf{x}^*, \hat{\boldsymbol{\theta}}) \ln p(\mathbf{y}|\mathbf{x}^*, \hat{\boldsymbol{\theta}}) d\mathbf{y}$$

- Data Uncertainty is captured via Maximum Likelihood Estimation
 - Given sufficient training data, model flexibility and correct output distribution

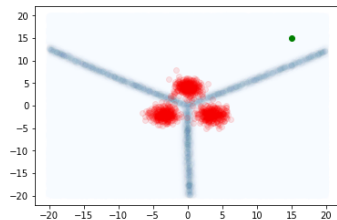
Sources of Uncertainty



(a) Data Uncertainty



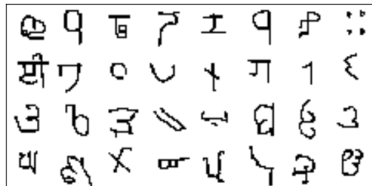
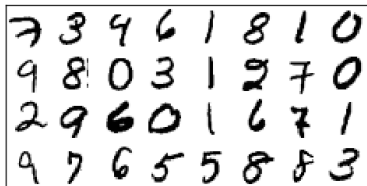
(b) Data Sparsity



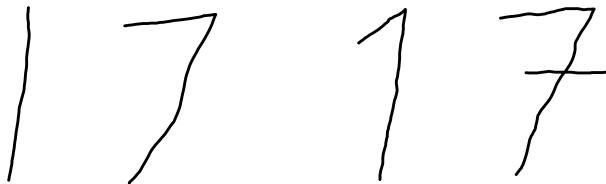
(c) Out-of-Distribution inputs

Knowledge Uncertainty - Out-of-Distribution

- Unseen classes



- Unseen variations of seen classes

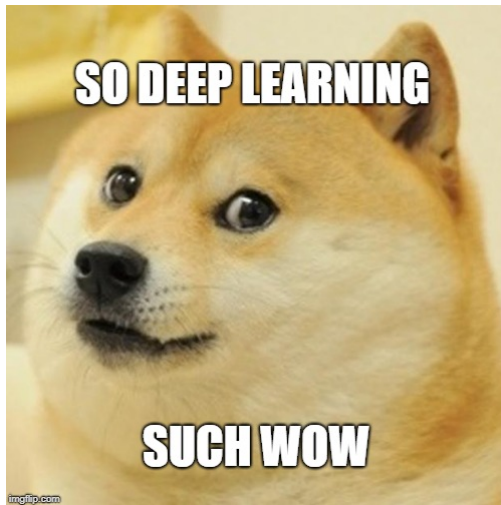


- Data Uncertainty → **Known-Unknown**
 - Class overlap (complexity of decision boundaries)
 - Homoscedastic and Heteroscedastic noise
- Knowledge Uncertainty → **Unknown-Unknown**
 - Test input in out-of-distribution region far from training data
- Appropriate **action** depends on **source** of uncertainty
 - Separating sources of uncertainty requires **Ensemble approaches**

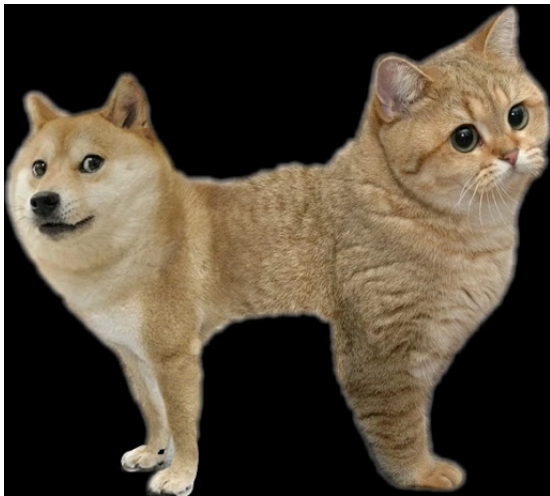
Overview of the Talk

1. Motivation and Sources
2. **Uncertainty Estimation via Ensembles**
3. Uncertainty Estimation via Prior Networks
4. Ensemble Distribution Distillation





Ensemble Approaches





- Uncertainty in θ captured by model posterior $p(\theta|\mathcal{D}) \rightarrow$

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}$$

- Can consider an **ensemble** of **probabilistic** models \rightarrow

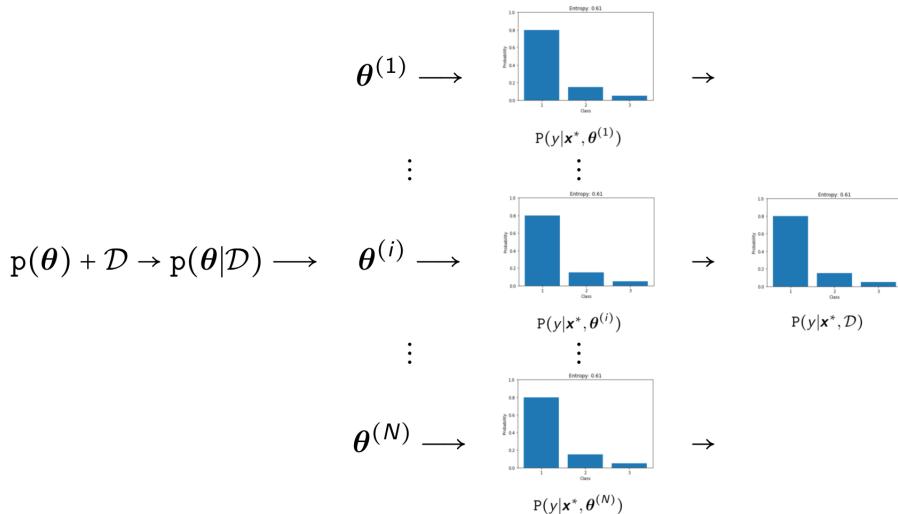
$$\{P(y|\mathbf{x}, \theta^{(m)})\}_{m=1}^M, \quad \theta^{(m)} \sim p(\theta|\mathcal{D}), \quad \{p(\mathbf{y}|\mathbf{x}, \theta^{(m)})\}_{m=1}^M, \quad \theta^{(m)} \sim p(\theta|\mathcal{D})$$
$$p(\mathbf{y}|\mathbf{x}, \theta) = \mathcal{N}(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Lambda}), \quad \{\boldsymbol{\mu}, \boldsymbol{\Lambda}\} = \mathbf{f}(\mathbf{x}; \theta)$$

- Bayesian inference of $P(y|\mathbf{x}^*, \theta) \rightarrow$

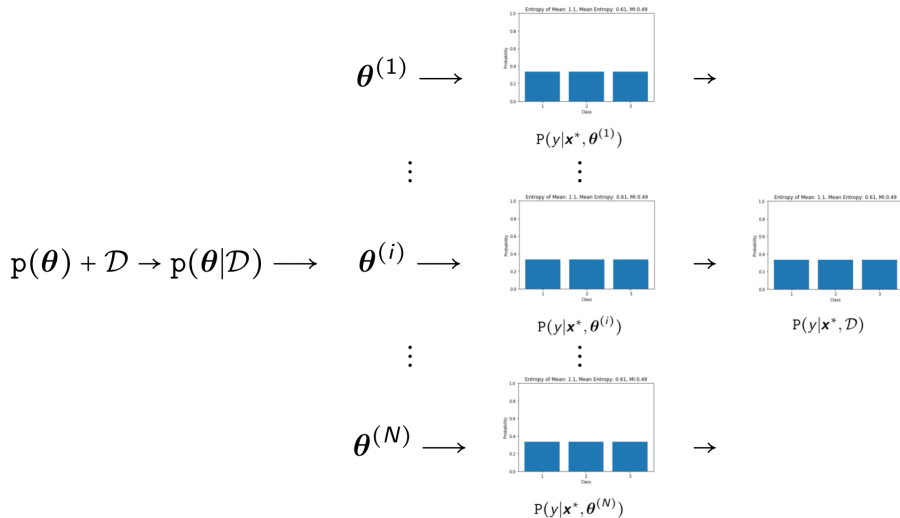
$$P(y|\mathbf{x}, \mathcal{D}) = \mathbb{E}_{p(\theta|\mathcal{D})}[P(y|\mathbf{x}, \theta)] \approx \frac{1}{M} \sum_{m=1}^M P(y|\mathbf{x}^*, \theta^{(m)}), \quad \theta^{(m)} \sim p(\theta|\mathcal{D})$$

- $P(y|\mathbf{x}^*, \mathcal{D})$ Is the **predictive posterior** or **ensemble mean**

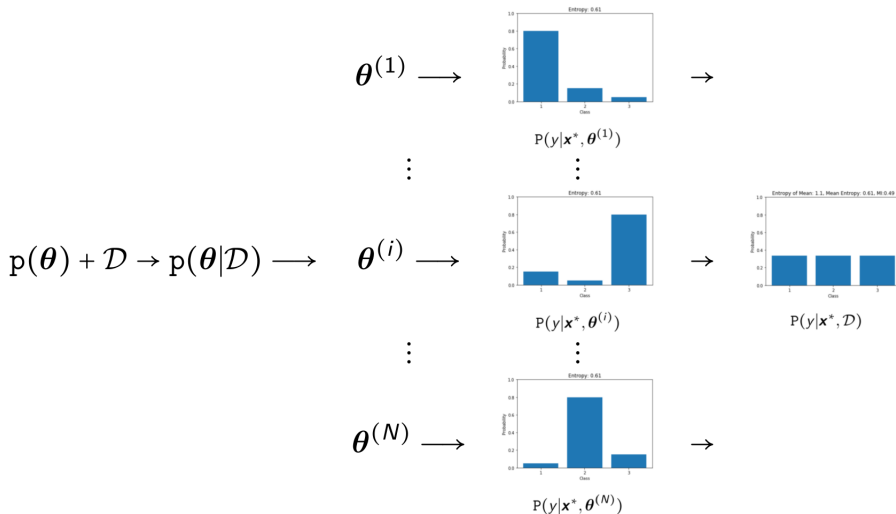
Ensemble for certain in-domain input



Ensemble for Out-of-Domain input



Ensemble for Out-of-Domain input



- Decompose sources of uncertainty via **Mutual Information** for classification:

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

- Mutual Information is a measure of **ensemble diversity**
- Intractable for regression, so use **Law of Total Variation**:

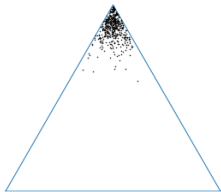
$$\underbrace{\mathbb{V}_{\mathbf{p}(\boldsymbol{\theta} | \mathcal{D})}[\boldsymbol{\mu}]}_{\text{Knowledge Uncertainty}} = \underbrace{\mathbb{V}_{\mathbf{p}(\mathbf{y} | \mathbf{x}, \mathcal{D})}[\mathbf{y}]}_{\text{Total Uncertainty}} - \underbrace{\mathbb{E}_{\mathbf{p}(\boldsymbol{\theta} | \mathcal{D})}[\boldsymbol{\Lambda}^{-1}]}_{\text{Data Uncertainty}} \quad (1)$$

Overview of the Talk

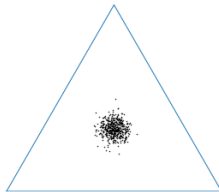
1. Motivation and Sources
2. Uncertainty Estimation via Ensembles
3. **Uncertainty Estimation via Prior Networks**
4. Ensemble Distribution Distillation

Distributions on a Simplex

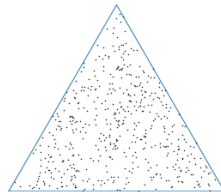
- Ensemble $\{P(y|\mathbf{x}^*, \boldsymbol{\theta}^{(m)})\}_{m=1}^M$ can be visualized on a [simplex](#)



(a) Confident



(b) Data Uncertainty



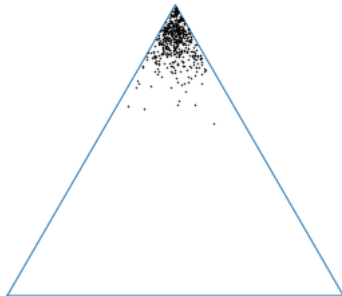
(c) Knowledge Uncertainty

- Same as sampling from **implicit** [Distribution over output Distributions](#)

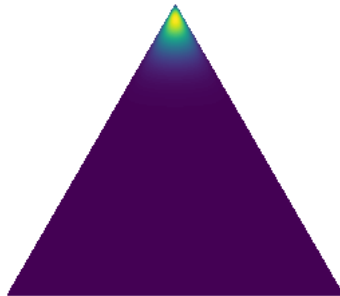
$$P(y|\mathbf{x}^*, \boldsymbol{\theta}^{(m)}) \sim p(\boldsymbol{\theta}|\mathcal{D}) \equiv \boldsymbol{\pi}^{(m)} \sim p(\boldsymbol{\pi}|\mathbf{x}^*, \mathcal{D})$$

- Expanding out $\pi^{(m)} = \begin{bmatrix} P(y = \omega_1) \\ P(y = \omega_2) \\ \vdots \\ P(y = \omega_K) \end{bmatrix}$, where each $\pi^{(m)}$ is a point on a simplex.

Distribution over Distributions

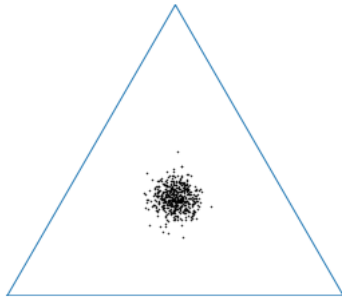


(a) $\{\pi^{(m)}\}_{m=1}^M$

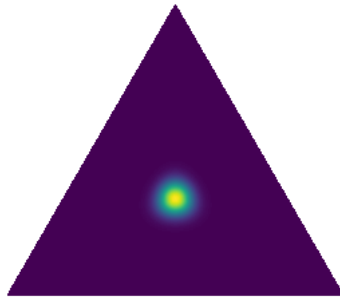


(b) $p(\pi|\mathbf{x}^*, \mathcal{D})$

Distribution over Distributions

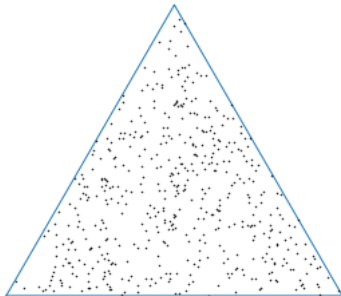


(a) $\{\pi^{(m)}\}_{m=1}^M$

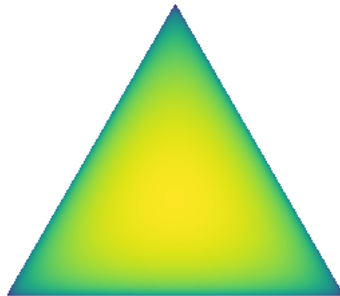


(b) $p(\pi|\mathbf{x}^*, \mathcal{D})$

Distribution over Distributions

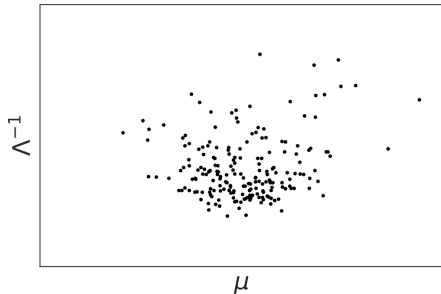


(a) $\{\pi^{(m)}\}_{m=1}^M$

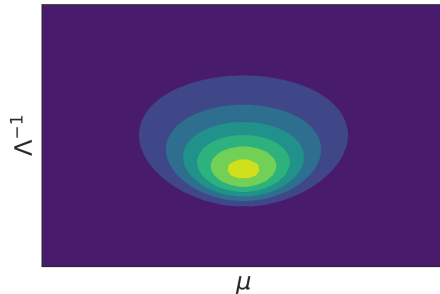


(b) $p(\pi|\mathbf{x}^*, \mathcal{D})$

Distribution over Distributions - Regression



(c) $\{\mu^{(m)}, \Lambda^{(m)}\}_{m=1}^M$



(d) $p(\mu, \Lambda | \mathbf{x}^*, \mathcal{D})$

- **Explicitly** model $p(\pi|\mathbf{x}^*, \mathcal{D})$ and $p(\mu, \Lambda|\mathbf{x}^*, \mathcal{D})$ using a **Prior Network**

$$\begin{aligned}p(\pi|\mathbf{x}^*; \hat{\theta}) &\approx p(\pi|\mathbf{x}^*, \mathcal{D}) \\p(\mu, \Lambda|\mathbf{x}^*, \theta) &\approx p(\mu, \Lambda|\mathbf{x}^*, \mathcal{D})\end{aligned}$$

- Predictive posterior distribution is given by expected categorical

$$\begin{aligned}P(y|\mathbf{x}^*; \hat{\theta}) &= \mathbb{E}_{p(\pi|\mathbf{x}^*; \hat{\theta})} [p(y|\pi)] = \hat{\pi} \\p(\mathbf{y}|\mathbf{x}^*; \hat{\theta}) &= \mathbb{E}_{p(\mu, \Lambda|\mathbf{x}^*; \hat{\theta})} [p(\mathbf{y}|\mu, \Lambda)]\end{aligned}$$

- A Classification Prior Network parametrizes the **Dirichlet Distribution**

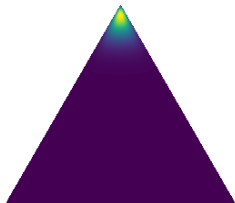
$$p(\pi|\mathbf{x}^*; \hat{\theta}) = \text{Dir}(\pi|\alpha), \quad \alpha = \mathbf{f}(\mathbf{x}^*; \hat{\theta})$$

- A Regression Prior Network parameterizes the **Normal-Wishart Distribution**

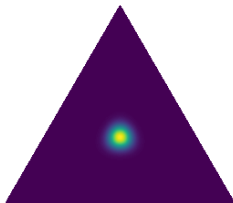
$$p(\mu, \mathbf{\Lambda}|\mathbf{x}^*, \theta) = \mathcal{NW}(\mu, \mathbf{\Lambda}|\mathbf{m}, \mathbf{L}, \kappa, \nu), \quad \{\mathbf{m}, \mathbf{L}, \kappa, \nu\} = \Omega = \mathbf{f}(\mathbf{x}^*; \theta)$$

- Dirichlet and Normal-Wishart Distributions \rightarrow
 - Conjugate priors to Categorical and Normal distributions, respectively.
 - Convenient properties \rightarrow analytically tractable!

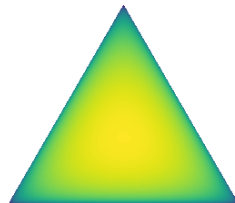
- Construct $p(\pi|\mathbf{x}^*, \hat{\theta})$ to emulate classification ensemble



(a) Low Uncertainty



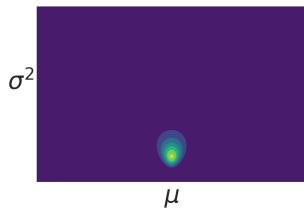
(b) Data Uncertainty



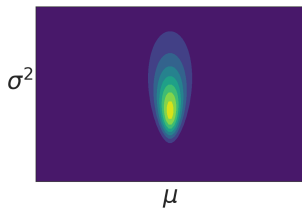
(c) Knowledge Uncertainty

Prior Networks vs Ensembles

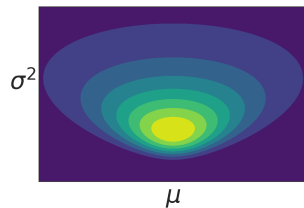
- Construct $p(\mu, \Lambda | \mathbf{x}^*, \hat{\theta})$ to emulate regression ensemble



(a) Low uncertainty



(b) Data uncertainty



(c) Knowledge Uncertainty

- Behaviour of Ensemble distribution over distributions
 - Controlled via **prior** $p(\theta)$ and **inference scheme**
- Behaviour of Prior Networks distribution over distributions
 - Controlled via **loss function** and **training data** \mathcal{D}

Uncertainty Measures for Prior Networks

- Ensemble uncertainty decomposition:

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

- Prior Network uncertainty decomposition

$$\underbrace{\mathcal{I}[y, \boldsymbol{\pi} | \mathbf{x}^*; \hat{\boldsymbol{\theta}}]}_{\text{Knowledge Uncertainty}} = \underbrace{\mathcal{H}[\mathbb{E}_{\mathbf{p}(\boldsymbol{\pi} | \mathbf{x}^*; \hat{\boldsymbol{\theta}})}[\mathbf{P}(y | \boldsymbol{\pi})]]}_{\text{Total Uncertainty}} - \underbrace{\mathbb{E}_{\mathbf{p}(\boldsymbol{\pi} | \mathbf{x}^*; \hat{\boldsymbol{\theta}})}[\mathcal{H}[\mathbf{P}(y | \boldsymbol{\pi})]]}_{\text{Data Uncertainty}}$$

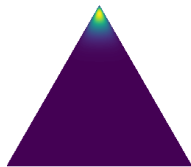
- **Ensemble** uncertainty decomposition (intractable!):

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

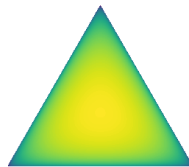
- **Prior Network** uncertainty decomposition (can be tractable!)

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

$$\mathcal{L}(\theta, \mathcal{D}) = \underbrace{\mathcal{L}_{in}(\theta, \mathcal{D}_{trn})}_{\text{In Domain Loss}} + \gamma \cdot \underbrace{\mathcal{L}_{out}(\theta, \mathcal{D}_{out})}_{\text{OOD Loss}}$$



(a) In-Domain Target



(b) OOD Target

- How to train **Distribution over Distributions** using only $\{y^{(i)}, \mathbf{x}^{(i)}\}_{i=1}^N$?

Reverse KL-Divergence Loss and the ELBO [3, 2]

- Consider using Bayes' rule as follows:

$$p(\pi|\hat{\alpha}^{(i)}) \propto p(y^{(i)}|\pi)^{\hat{\beta}} p(\pi|\alpha_0), \quad p(\mu, \Lambda|\hat{\Omega}^{(i)}) \propto p(y^{(i)}|\mu, \Lambda)^{\hat{\beta}} p(\mu, \Lambda|\Omega_0)$$

- Minimizing **Reverse KL-Divergence** induces an ELBO-like loss:

$$\text{KL}[p(\pi|\mathbf{x}, \theta) \| p(\pi|\hat{\alpha}^{(i)})] = \underbrace{\hat{\beta} \cdot \mathbb{E}_{p(\pi|\mathbf{x}, \theta)} [-\ln p(y|\pi)]}_{\text{Reconstruction term}} + \underbrace{\text{KL}[p(\pi|\mathbf{x}, \theta) \| p(\pi|\alpha_0)]}_{\text{Prior}} + Z$$

$$\begin{aligned} \text{KL}[p(\mu, \Lambda|\mathbf{x}, \theta) \| p(\mu, \Lambda|\hat{\Omega}^{(i)})] &= \\ &= \hat{\beta} \cdot \mathbb{E}_{p(\mu, \Lambda|\mathbf{x}, \theta)} [-\ln p(y|\mu, \Lambda)] + \text{KL}[p(\mu, \Lambda|\mathbf{x}, \theta) \| p(\mu, \Lambda|\Omega_0)] + Z \end{aligned}$$

- Set $\hat{\beta} \gg 0$ in-domain and $\hat{\beta} = 0$ out-of-domain.

Reverse KL-Divergence Loss and the ELBO

- Prior parameters α_0 and $\Omega_0 = \{\mathbf{m}_0, \mathbf{L}_0, \kappa_0, \nu_0\}$ defined as follows:

$$\alpha_0 = \mathbf{1}$$

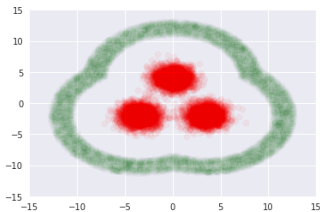
$$\mathbf{m}_0 = \sum_{i=1}^N \frac{\mathbf{y}^{(i)}}{N}, \quad \mathbf{L}_0^{-1} = \frac{\nu_0}{N} \sum_{i=1}^N (\mathbf{y}^{(i)} - \mathbf{m}_0)(\mathbf{y}^{(i)} - \mathbf{m}_0)^T, \quad \kappa_0 = \epsilon, \nu_0 = K + 1 + \epsilon$$

- Prior for classification - uninformative flat Dirichlet Prior
- Prior for regression - semi-informative Prior (uninformative would be improper)

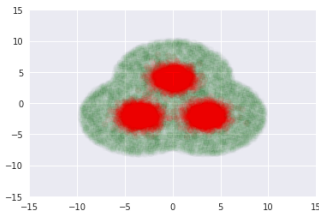
- But how to obtain out-of-domain training data $\mathcal{D}_{OOD} = \hat{p}_{out}(\mathbf{x})$?
 - Use a [different dataset](#), eg: CIFAR10 vs CIFAR100
 - [Synthesize](#) using generative model (VAE/GAN)
 - Generate using [adversarial attacks](#)
- Choice is highly non-trivial for many tasks (Depth Estimation) → main downside!

Reverse KL-Divergence Loss and the ELBO

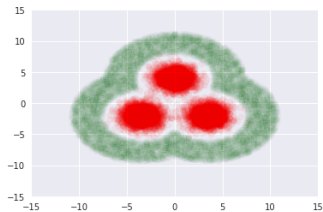
- Out-of-domain (OOD) training data must be on *boundary* on in-domain region →
 - Too loose → Some OOD might be considered in-domain
 - Too tight → Some in-domain might be considered OOD



(a) Too Loose

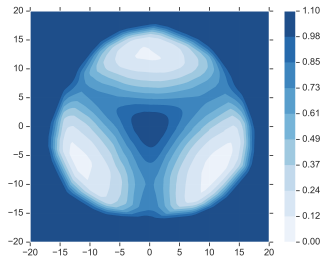


(b) Too Tight

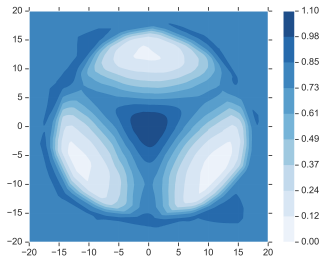


(c) Good

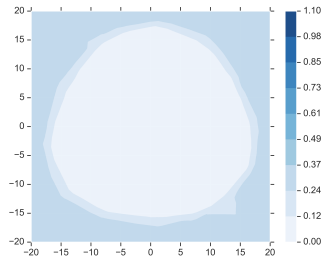
Prior Networks trained with RKL loss on Artificial Data



(a) Total Uncertainty



(b) Data Uncertainty



(c) Knowledge Uncertainty

Overview of the Talk

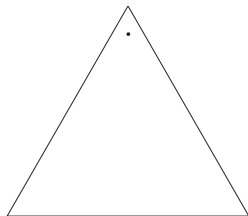
1. Motivation and Sources
2. Uncertainty Estimation via Ensembles
3. Uncertainty Estimation via Prior Networks
4. **Ensemble Distribution Distillation**

Ensemble Distribution Distillation [4, 2]

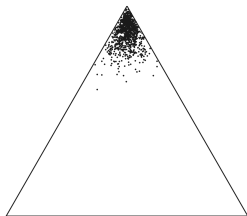
- Ensembles of multiple independently trained models $\{p(y|\mathbf{x}, \theta^{(m)})\}_{m=1}^M$
 - Improved performance
 - Robust uncertainty estimates derived from mean and diversity
 - **Computationally expensive!**
- Ensemble Distillation (EnD) \rightarrow distill ensemble **mean** into a single model
 - Improved performance and low computational cost
 - Lose information about **diversity** \rightarrow cannot separate data and knowledge uncertainty
- **Ensemble Distribution Distillation (EnD²)** \rightarrow
 - Distill mean and diversity of ensemble into single model
 - Improved performance and robust uncertainty at low computational cost

Ensemble Distribution Distillation (EnD²) for Classification

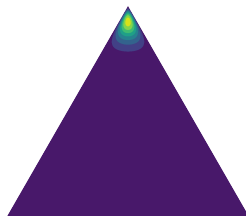
$$\frac{1}{M} \sum_{m=1}^M P(y|\mathbf{x}, \boldsymbol{\theta}^{(m)}) \quad \longleftarrow \{P(y|\mathbf{x}, \boldsymbol{\theta}^{(m)})\}_{m=1}^M \longrightarrow p(\boldsymbol{\pi}|\mathbf{x}; \boldsymbol{\phi})$$



(a) EnD



(b) Ensemble



(c) EnD²

- Distill **ensemble distribution** (mean and diversity) into a **single** model
 - Fully capture all information about the ensemble

Ensemble Distribution Distillation (EnD²) for Classification [4]

- Parameterize a Dirichlet distribution using neural network:

$$p(\boldsymbol{\pi}|\mathbf{x}; \boldsymbol{\phi}) = \text{Dir}(\boldsymbol{\pi}; \boldsymbol{\alpha}), \quad \boldsymbol{\alpha} = \mathbf{f}(\mathbf{x}; \boldsymbol{\phi}), \quad \alpha_c > 0$$

- Training data are ensemble predictions for every input:

$$\mathcal{D} = \left\{ \left\{ p(y|\mathbf{x}^{(i)}; \boldsymbol{\theta}^{(m)}), \mathbf{x}^{(i)} \right\}_{m=1}^M \right\}_{i=1}^N \sim \hat{\mathbf{p}}(\boldsymbol{\pi}, \mathbf{x})$$

- Train via Maximum Likelihood:

$$\mathcal{L}(\boldsymbol{\phi}, \mathcal{D}) = - \mathbb{E}_{\hat{\mathbf{p}}(\mathbf{x})} \left[\mathbb{E}_{\hat{\mathbf{p}}(\boldsymbol{\pi}|\mathbf{x})} [\ln p(\boldsymbol{\pi}|\mathbf{x}; \boldsymbol{\phi})] \right]$$

- Predict using **mean**, derive uncertainty from **mean** and **diversity**

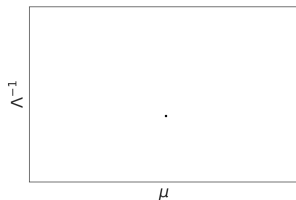
Classification results (% Error and % ROC-AUC)

| Method | CIFAR-10 | CIFAR-100 | TinyImageNet |
|------------------|----------|-----------|--------------|
| Single | 8.0 | 30.4 | 41.8 |
| Ensemble | 6.2 | 26.3 | 36.6 |
| EnD | 6.7 | 28.2 | 38.5 |
| EnD ² | 6.9 | 28.0 | 37.3 |

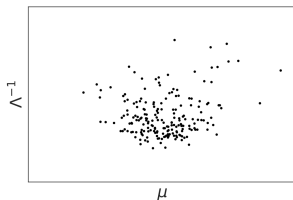
| Model | CIFAR100 vs. LSUN | | CIFAR100 vs. TinyImageNet | |
|------------------|-------------------|----------------|---------------------------|----------------|
| | Total Unc. | Knowledge Unc. | Total Unc. | Knowledge Unc. |
| Ensemble | 82.4 | 88.4 | 76.6 | 81.7 |
| EnD | 76.5 | - | 70.0 | - |
| EnD ² | 83.5 | 86.9 | 76.4 | 79.3 |

Ensemble Distribution Distillation (EnD²) for Regression [2]

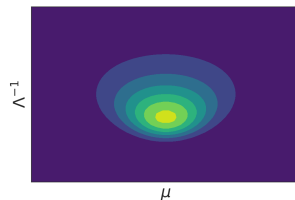
$$\frac{1}{M} \sum_{m=1}^M p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}^{(m)}) \quad \longleftarrow \{p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}^{(m)})\}_{m=1}^M \longrightarrow p(\boldsymbol{\mu}, \boldsymbol{\Lambda}|\mathbf{x}; \phi)$$



(a) EnD



(b) Ensemble



(c) EnD²

- Distill **ensemble distribution** (mean and diversity) into a **single** model
 - Fully capture all information about the ensemble

Ensemble Distribution Distillation (EnD²) for Regression

- Construct a *Regression Prior Network*

$$\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Lambda}), \quad \mathcal{NW}(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathbf{m}, \mathbf{L}, \kappa, \nu); \quad \{\mathbf{m}, \mathbf{L}, \kappa, \nu\} = f(\mathbf{x}; \phi)$$

- Training data are ensemble predictions for every input:

$$\mathcal{D} = \left\{ \left\{ p(\mathbf{y} | \mathbf{x}^{(i)}; \boldsymbol{\theta}^{(m)}), \mathbf{x}^{(i)} \right\}_{m=1}^M \right\}_{i=1}^N = \left\{ \left\{ \boldsymbol{\mu}^{(m,i)}, \boldsymbol{\Lambda}^{(m,i)} \right\}_{m=1}^M, \mathbf{x}^{(i)} \right\}_{i=1}^N \sim \hat{\mathbf{p}}(\boldsymbol{\mu}, \boldsymbol{\Lambda}, \mathbf{x})$$

- Train via Maximum Likelihood:

$$\mathcal{L}(\phi, \mathcal{D}) = - \mathbb{E}_{\hat{\mathbf{p}}(\mathbf{x})} \left[\mathbb{E}_{\hat{\mathbf{p}}(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathbf{x})} [\ln \mathbf{p}(\boldsymbol{\mu}, \boldsymbol{\Lambda} | \mathbf{x}; \phi)] \right]$$

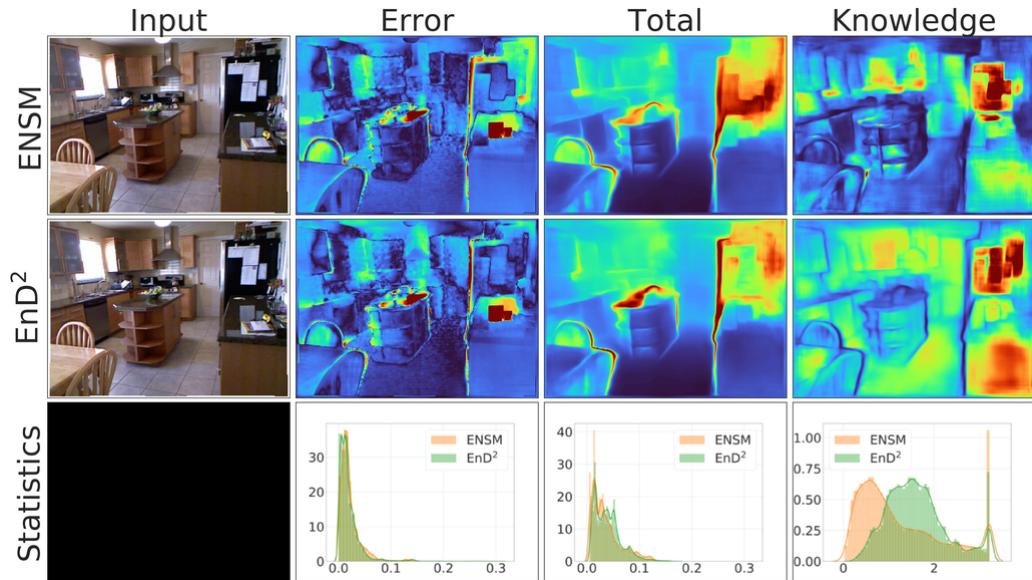
- Predict using **mean**, derive uncertainty from **mean** and **diversity**

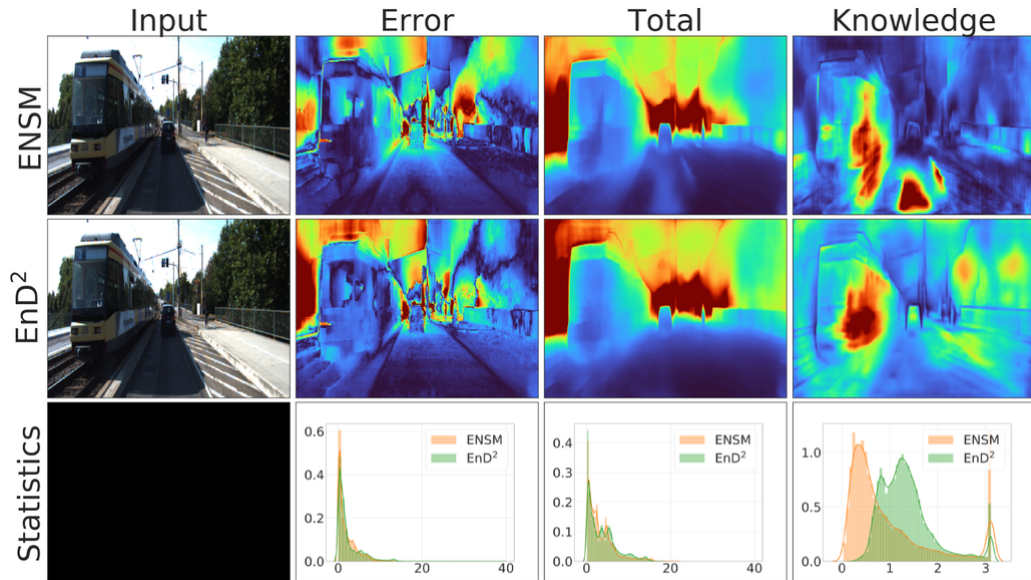
Monocular Depth Estimation - Predictive Performance

| Method | NYUv2 | | | KITTI | | |
|------------------|---------------------|----------------------|---------------------|---------------------|----------------------|---------------------|
| | rel(\downarrow) | rmse(\downarrow) | NLL(\downarrow) | rel(\downarrow) | rmse(\downarrow) | NLL(\downarrow) |
| ENSM 5 | 0.117 | 0.438 | 0.76 | 0.073 | 3.355 | 1.94 |
| EnD ² | 0.120 | 0.451 | -1.47 | 0.075 | 3.367 | 1.42 |
| MD-EnD | 0.121 | 0.451 | 8.48 | 0.079 | 3.446 | 2.30 |
| DER | 0.125 | 0.464 | -1.04 | 0.078 | 3.552 | 1.71 |

Monocular Depth Estimation - OOD Detection (ROC-AUC)

| Method | OOD | NYUv2 vs LSUN | | KITTI vs LSUN | |
|------------------------|-------|---------------|----------------|---------------|----------------|
| | | Total Unc. | Knowledge Unc. | Total Unc. | Knowledge Unc. |
| ENSM | LSN-B | 72.3 | 74.5 | 03.2 | 82.2 |
| EnD ² (Our) | | 73.3 | 81.7 | 1.7 | 88.7 |
| MD-EnD | | 63.0 | 50.2 | 0.4 | 44.8 |
| ENSM | LSN-C | 88.7 | 88.6 | 03.6 | 77.9 |
| EnD ² (Our) | | 89.3 | 96.4 | 2.0 | 83.4 |
| DER | | 87.7 | 87.8 | 03.5 | 04. |
| MD-EnD | | 69.8 | 42.2 | 01.2 | 50.6 |





- EnD^2 - still a single model \rightarrow evaluate robustness
- Are Dirichlet and Normal-Wishart appropriate?
 - Do we need to model ensemble in model detail?
 - Do we need to only capture bulk properties?
- Do we need auxiliary training data? Mixup?
- Can we use EnD^2 for analysis?
- Can we combine ensemble generation and EnD^2 ?

Thank you! Questions?

- [1] Andrey Malinin and Mark Gales,
“Predictive uncertainty estimation via prior networks,”
in Advances in Neural Information Processing Systems, 2018, pp. 7047–7058.
- [2] Andrey Malinin, Sergey Chervontsev, Ivan Provilkov, and Mark Gales,
“Regression prior networks,”
arXiv preprint arXiv:2006.11590, 2020.
- [3] Andrey Malinin and Mark JF Gales,
“Reverse kl-divergence training of prior networks: Improved uncertainty and adversarial robustness,”
in Advances in Neural Information Processing Systems, 2019.

- [4] Andrey Malinin, Bruno Mlodozienec, and Mark JF Gales,
“Ensemble distribution distillation,”
in *International Conference on Learning Representations*, 2020.