

Reverse KL-Divergence Training of Prior Networks: Improved Uncertainty and Adversarial Robustness

Andrey Malinin and Mark Gales

18 October 2019

Joint work with...



(a) Prof. Mark Gales

Overview of the Talk

1. Context: Why do we need Uncertainty Estimation?
2. RECAP: Sources of Uncertainty in Predictions
3. RECAP: Ensemble Approaches
4. Prior Networks
5. Adversarial Attack Detection

- 1. Context: Why do we need Uncertainty Estimation?**
2. Sources of Uncertainty in Predictions
3. Ensemble Approaches
4. Prior Networks
5. Adversarial Attack Detection

- Given a **deployed** model and a **test input 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

Applications of Uncertainty Estimation

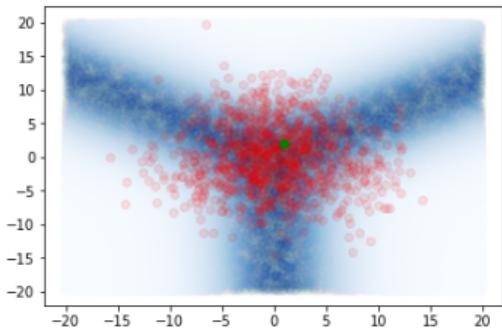
- Uncertainty should be assessed in the context of an **application**
- Threshold-based outlier detection →
 - Misclassification Detection [Hendrycks and Gimpel, 2016]
 - Out-of-distribution input Detection [Malinin and Gales, 2019]
 - Adversarial Attack Detection [Malinin and Gales, 2019]
- Active Learning [Gal, 2016]
- Reinforcement Learning uncertainty-driven exploration [Osband et al., 2016]
- Other...

Assessment of Uncertainty Quality

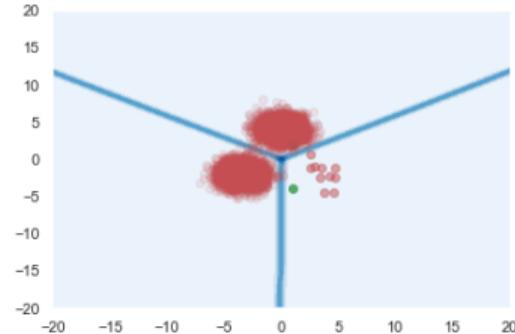
- Uncertainty should be assessed in the context of an **application**
- Threshold-based outlier detection →
 - Misclassification Detection [Hendrycks and Gimpel, 2016]
 - **Out-of-distribution input Detection** [Malinin and Gales, 2019]
 - **Adversarial Attack Detection** [Malinin and Gales, 2019]
- Active Learning [Gal, 2016]
- Reinforcement Learning uncertainty-driven exploration [Osband et al., 2016]
- Other...

1. Context: Why do we need Uncertainty Estimation?
2. **Sources of Uncertainty in Predictions**
3. Ensemble Approaches
4. Prior Networks
5. Adversarial Attack Detection

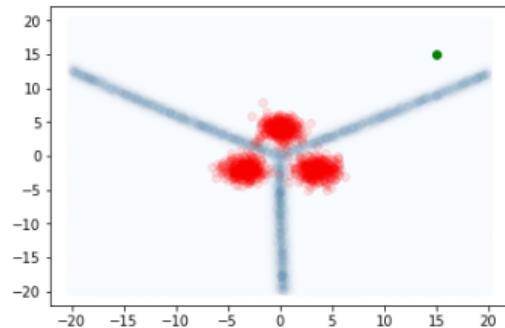
Sources of Uncertainty



(a) Data Uncertainty



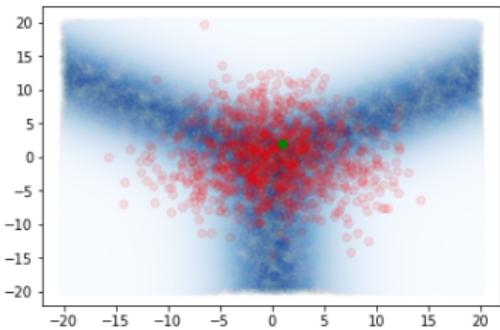
(b) Data Sparsity



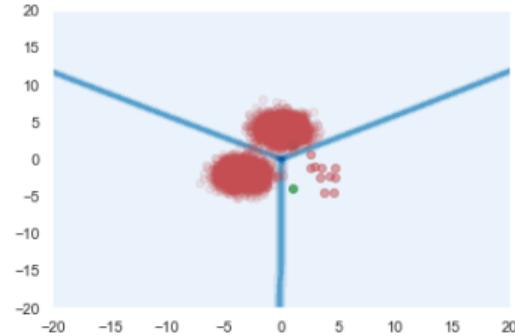
(c) Out-of-Distribution inputs

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

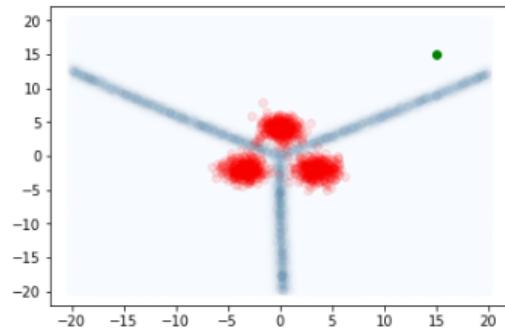
Sources of Uncertainty



(a) Data Uncertainty



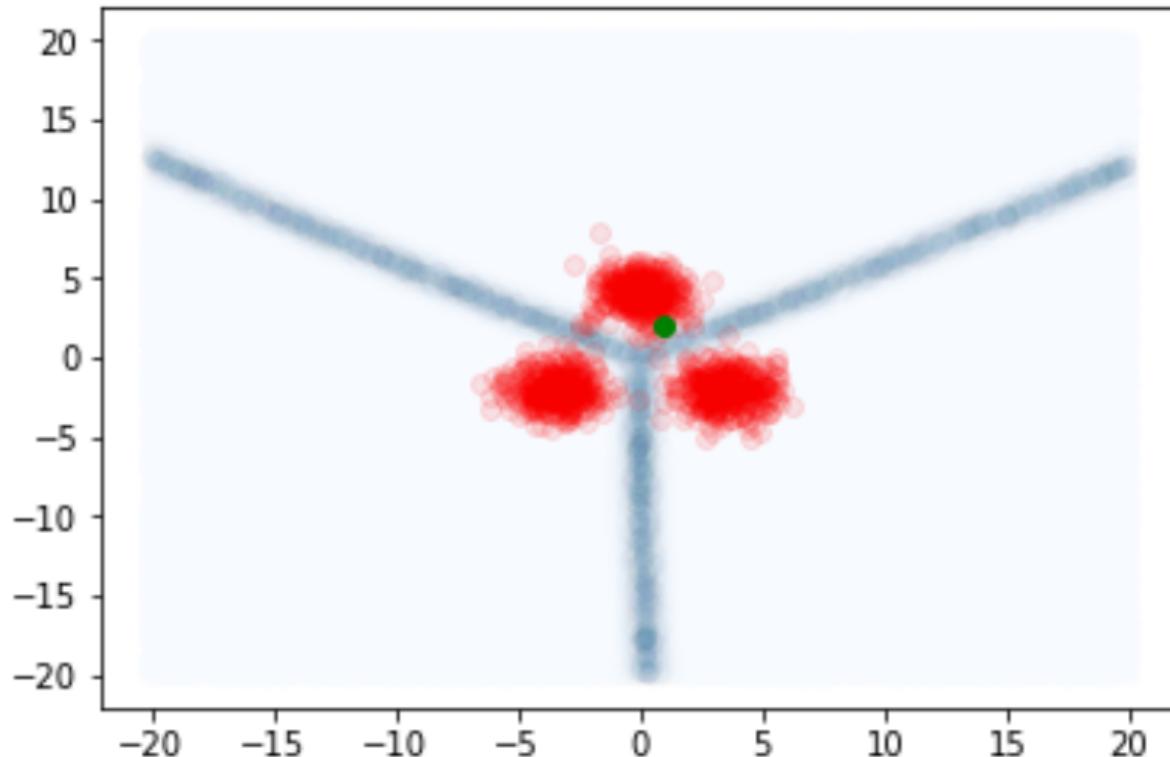
(b) Knowledge Uncertainty



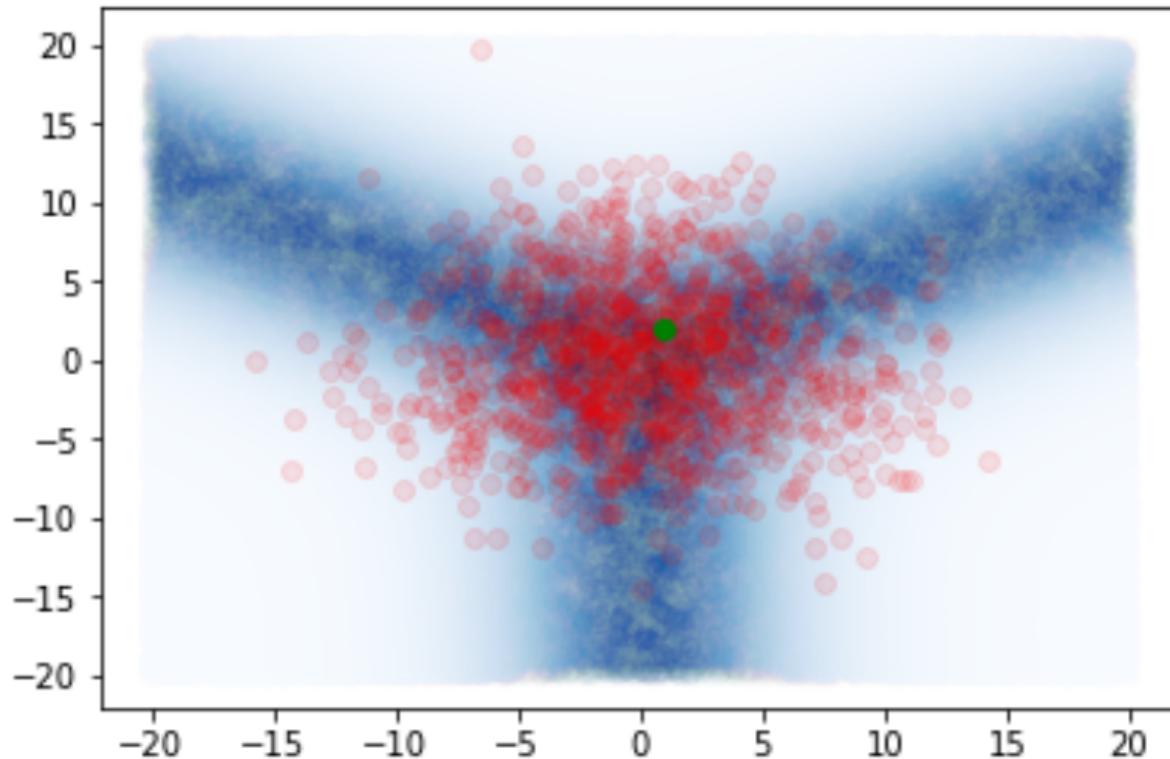
(c) Knowledge Uncertainty

- Knowledge (epistemic) uncertainty refers to both:
 - Data Sparsity **and** Knowledge Uncertainty

Data (Aleatoric) Uncertainty



Data Uncertainty



Data Uncertainty

- Distinct Classes



- Overlapping Classes



- Data Uncertainty → **Known-Unknown**
- Uncertainty due to properties of data
 - Class overlap (complexity of decision boundaries)
 - Human labelling error

Data Uncertainty

- Data Uncertainty is the *entropy* of the *true data distribution* →

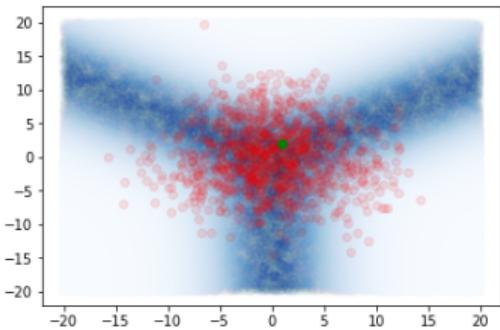
$$\mathcal{H}[P_{\text{tr}}(y|\mathbf{x}^*)] = - \sum_{c=1}^K P_{\text{tr}}(y = \omega_c | \mathbf{x}^*) \ln P_{\text{tr}}(y = \omega_c | \mathbf{x}^*)$$

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

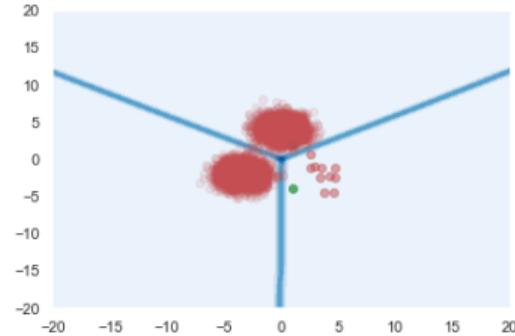
$$\mathcal{H}[P(y|\mathbf{x}^*, \hat{\theta})] = - \sum_{c=1}^K P(y = \omega_c | \mathbf{x}^*, \hat{\theta}) \ln P(y = \omega_c | \mathbf{x}^*, \hat{\theta})$$

- Data Uncertainty is captured as a consequence of Maximum Likelihood Estimation

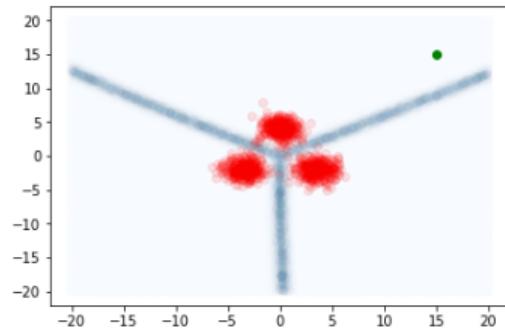
Sources of Uncertainty



(a) Data Uncertainty



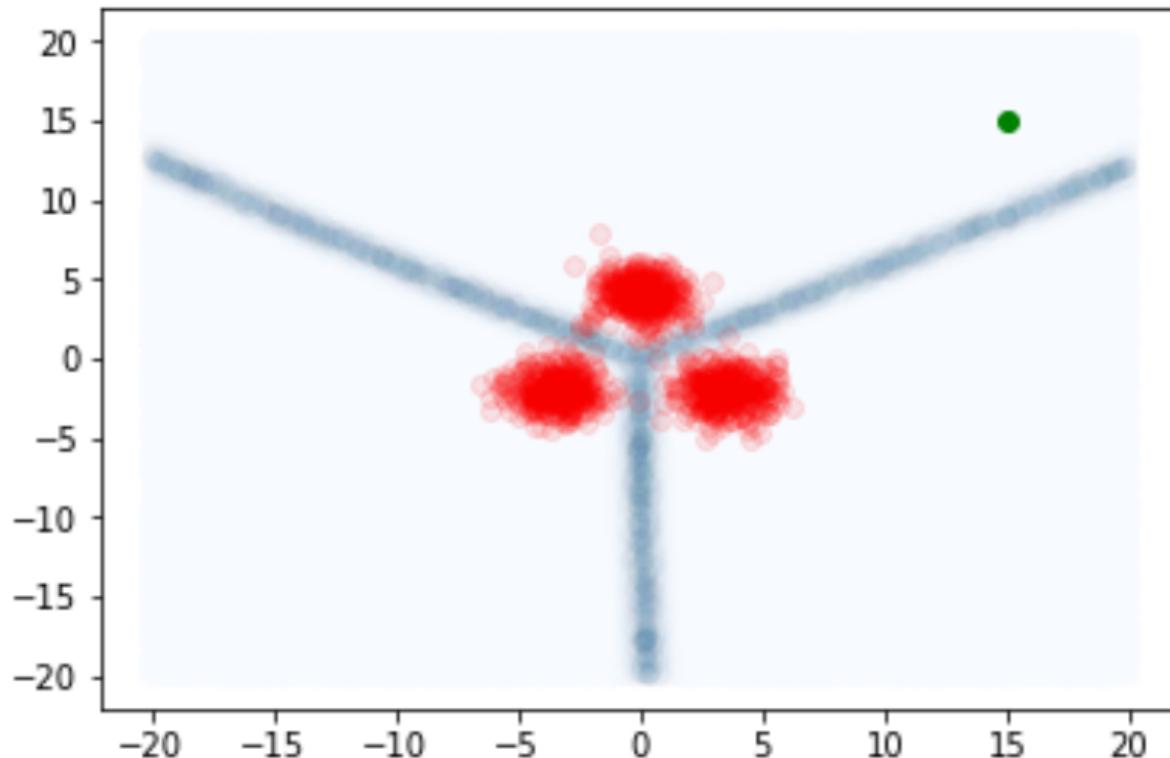
(b) Data Sparsity



(c) Out-of-Distribution inputs

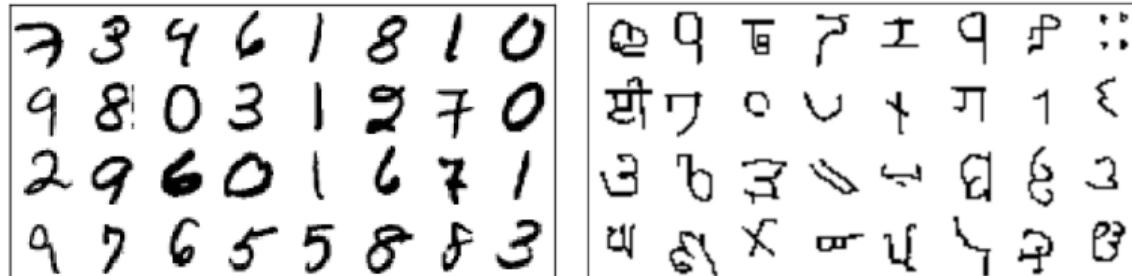
- Knowledge (epistemic) uncertainty refers to both:
 - Data Sparsity **and** Out-of-distribution inputs

Knowledge Uncertainty

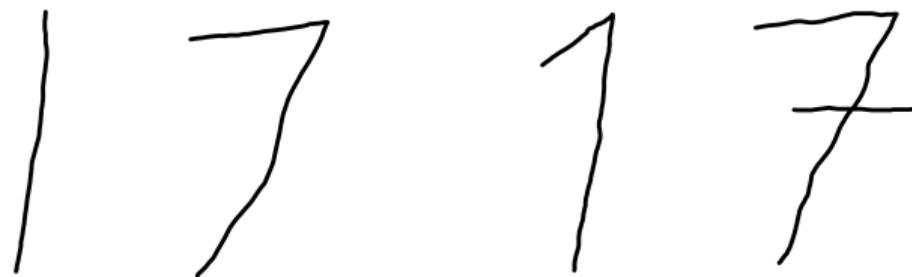


Knowledge Uncertainty - Out-of-Distribution

- Unseen classes



- Unseen variations of seen classes



Sources of Uncertainty

- Data Uncertainty → **Known-Unknown**
 - Class overlap (complexity of decision boundaries)
 - Human labelling error
- 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**
 - ... or **Prior Networks**

1. Context: Why do we need Uncertainty Estimation?
2. Sources of Uncertainty in Predictions
3. **Ensemble Approaches**
4. Prior Networks
5. Assessment of Uncertainty Quality

Ensemble Approaches

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

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

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

$$P(y|\mathbf{x}^*, \mathcal{D}) = \mathbb{E}_{p(\theta|\mathcal{D})}[P(y|\mathbf{x}^*, \theta)]$$

- Can consider an ensemble of models \rightarrow

$$\{P(y|\mathbf{x}^*, \theta^{(m)})\}_{m=1}^M, \theta^{(m)} \sim p(\theta|\mathcal{D})$$

- Choose desired behaviour of ensemble via prior $p(\theta)$

Total Uncertainty

- Consider the entropy of the predictive posterior $P(y|\mathbf{x}^*, \mathcal{D}) \rightarrow$

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

- Measure of Total Uncertainty
 - Combination of Data uncertainty and Knowledge uncertainty

Expected Data Uncertainty

- Lets consider an ensemble of models $\{P(y|\mathbf{x}^*, \boldsymbol{\theta}^{(m)})\}_{m=1}^M$, $\boldsymbol{\theta}^{(m)} \sim p(\boldsymbol{\theta}|\mathcal{D})$
 - Each model $P(y|\mathbf{x}^*, \boldsymbol{\theta}^{(m)})$ captures an **different** estimate of data uncertainty.
- Ensemble estimate of data uncertainty → **Expected Data Uncertainty**

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

- **Not** the same as entropy of the predictive posterior $P(y|\mathbf{x}^*, \mathcal{D})$

Model Uncertainty

- If the predictions from the models are consistent

$$\underbrace{\mathcal{H}[\mathbb{E}_{p(\theta|\mathcal{D})}[P(y|\mathbf{x}^*, \theta)]]}_{\text{Total Uncertainty}} - \underbrace{\mathbb{E}_{p(\theta|\mathcal{D})}[\mathcal{H}[P(y|\mathbf{x}^*, \theta)]]}_{\text{Expected Data Uncertainty}} = 0$$

- If the predictions from the models are diverse

$$\underbrace{\mathcal{H}[\mathbb{E}_{p(\theta|\mathcal{D})}[P(y|\mathbf{x}^*, \theta)]]}_{\text{Total Uncertainty}} - \underbrace{\mathbb{E}_{p(\theta|\mathcal{D})}[\mathcal{H}[P(y|\mathbf{x}^*, \theta)]]}_{\text{Expected Data Uncertainty}} > 0$$

- Difference of the two is a measure of model uncertainty

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

Model Uncertainty → Knowledge Uncertainty

- If the predictions from the models are consistent

$$\underbrace{\mathcal{H}[\mathbb{E}_{p(\theta|\mathcal{D})}[P(y|x^*, \theta)]]}_{\text{Total Uncertainty}} - \underbrace{\mathbb{E}_{p(\theta|\mathcal{D})}[\mathcal{H}[P(y|x^*, \theta)]]}_{\text{Expected Data Uncertainty}} = 0$$

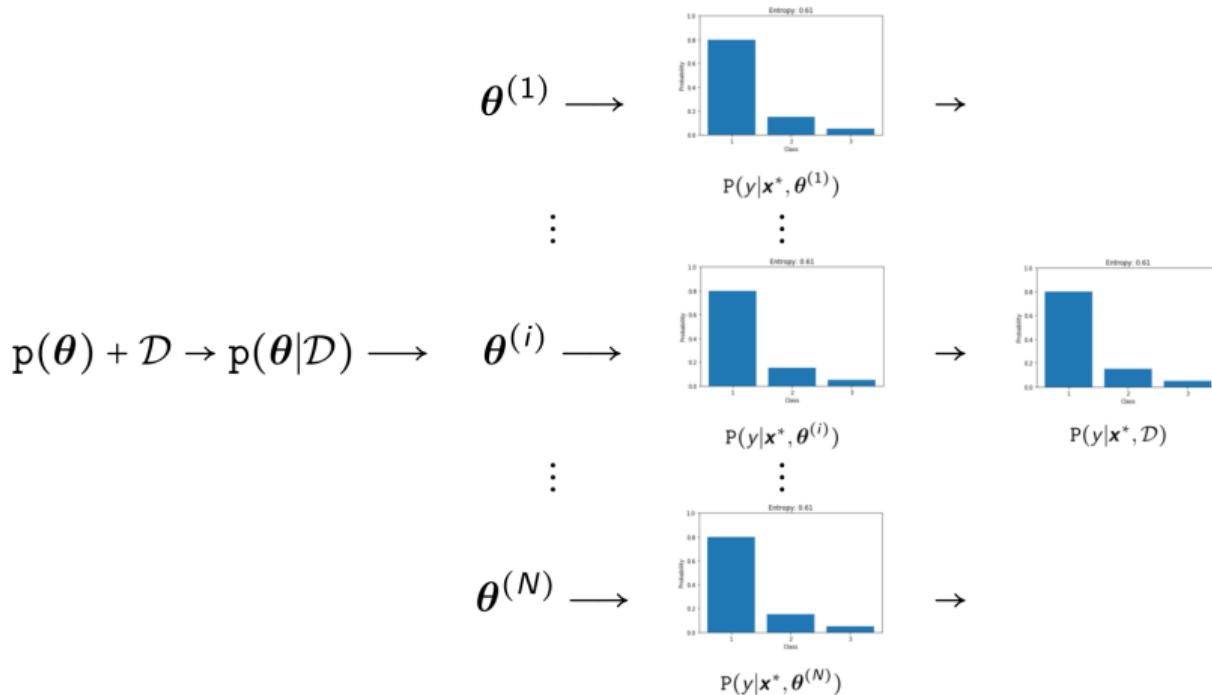
- If the predictions from the models are diverse

$$\underbrace{\mathcal{H}[\mathbb{E}_{p(\theta|\mathcal{D})}[P(y|x^*, \theta)]]}_{\text{Total Uncertainty}} - \underbrace{\mathbb{E}_{p(\theta|\mathcal{D})}[\mathcal{H}[P(y|x^*, \theta)]]}_{\text{Expected Data Uncertainty}} > 0$$

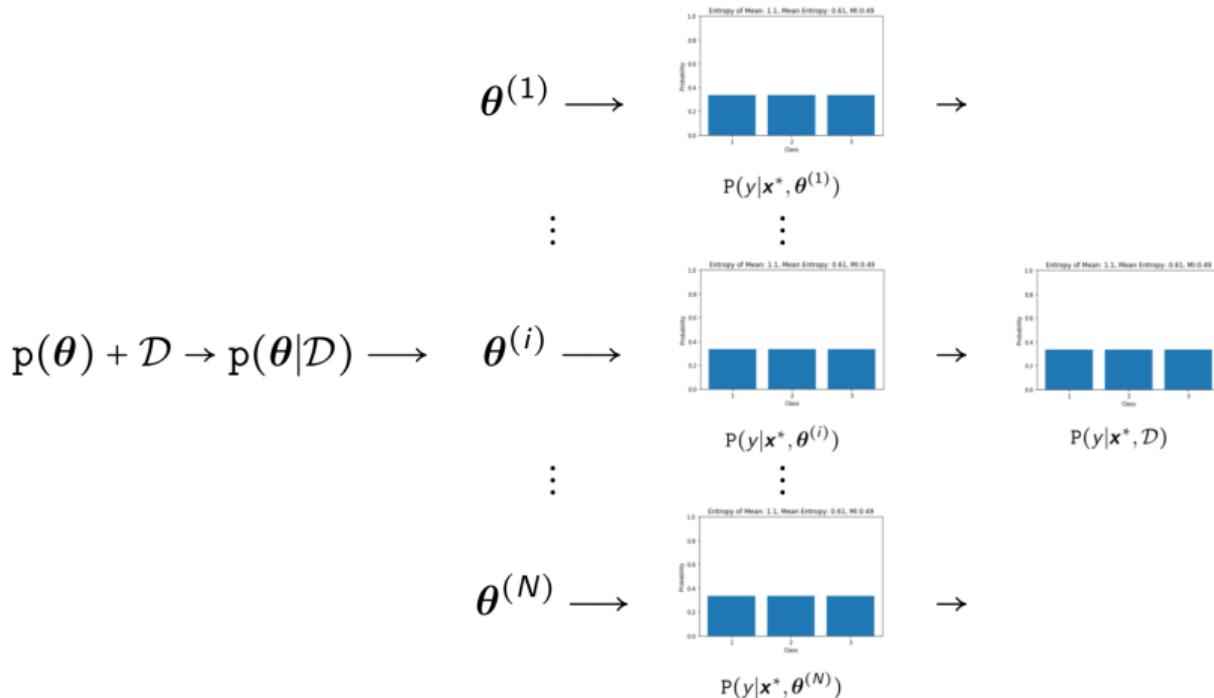
- Difference of the two is a measure of knowledge uncertainty

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

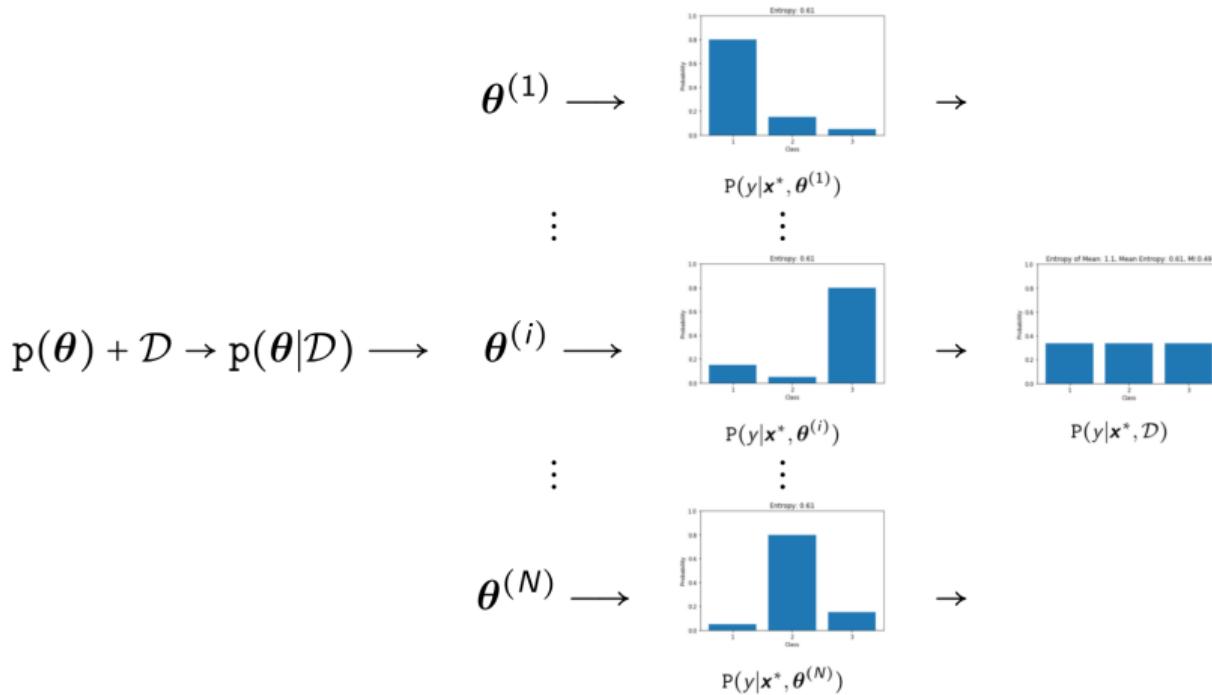
Ensemble for certain in-domain input



Ensemble for uncertain in-domain input

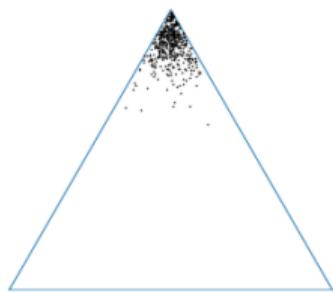


Ensemble for Out-of-Domain input

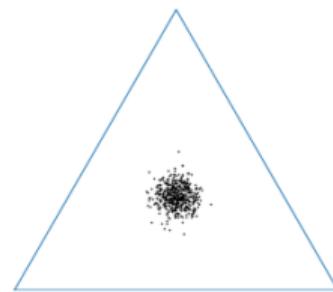


Distributions on a Simplex

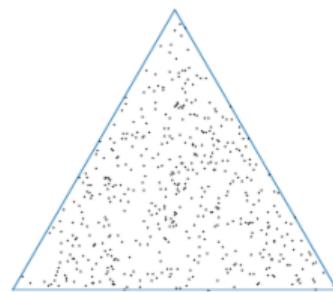
- Ensemble $\{P(y|\mathbf{x}^*, \theta^{(m)})\}_{m=1}^M$ can be visualized on a simplex for an input \mathbf{x}^*



(a) Confident



(b) Data Uncertainty



(c) Knowledge Uncertainty

- Ideally - compute all measures of uncertainty in **closed form**...
 - But **inference** is **intractable** for neural networks
 - Bayes' Rule is **intractable** for neural networks
- Solutions → use approximate inference
 - Compute approximate posterior $q(\theta) \approx p(\theta|\mathcal{D})$
 - Use **variational approximations** to measures of uncertainty
 - Use **Monte-Carlo approximations** to measures of uncertainty

- Variational Inference:
 - Bayes by Backprop [Blundell et al., 2015]
 - Probabalistic Backpropagation [Hernández-Lobato and Adams, 2015]
- Monte-Carlo Methods:
 - Monte-Carlo Dropout [Gal, 2016, Gal and Ghahramani, 2016]
 - Stochastic Gradient Langevin Dynamics [Welling and Teh, 2011]
 - Fast-Ensembling via Mode Connectivity [Garipov et al., 2018]
 - Stochastic Weight Averaging Gaussian (SWAG) [Maddox et al., 2019]
- Non-Bayesian Ensembles:
 - Bootstrap DQN [Osband et al., 2016]
 - Deep Ensembles [Lakshminarayanan et al., 2017]

Limitations

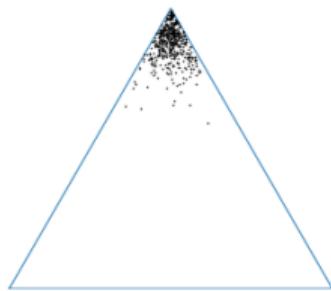
- Hard to guarantee diverse $\{P(y|\mathbf{x}^*, \theta^{(m)})\}_{m=1}^M$ for OOD \mathbf{x}^*
- Diversity of ensemble depends on:
 - Selection of prior
 - Nature of approximations
 - Architecture of network
 - Properties and size of data
- May be computationally expensive

Overview of the Talk

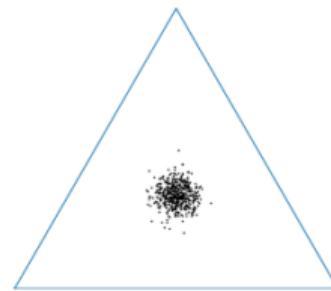
1. Context: Why do we need Uncertainty Estimation?
2. Sources of Uncertainty in Predictions
3. Ensemble Approaches
4. **Prior Networks**
5. Adversarial Attack Detection

Distributions on a Simplex

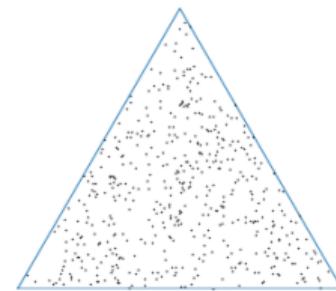
- Ensemble $\{P(y|\mathbf{x}^*, \boldsymbol{\theta}^{(m)})\}_{m=1}^M$ can be visualized on a simplex



(a) Confident



(b) Data Uncertainty



(c) Knowledge Uncertainty

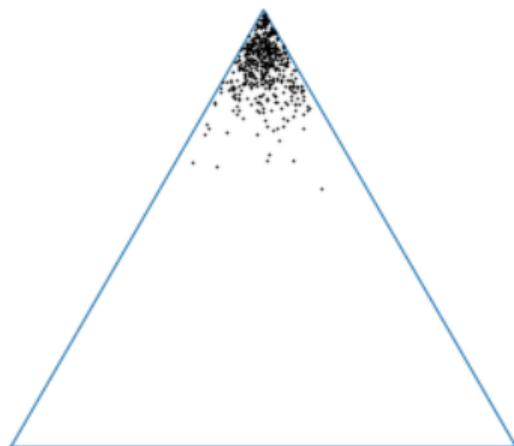
- Same as sampling from **implicit** Distribution over output Distributions

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

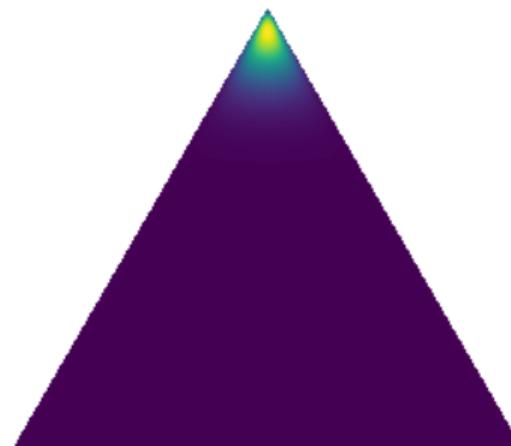
Distributions on a Simplex (cont)

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

Distribution over Distributions

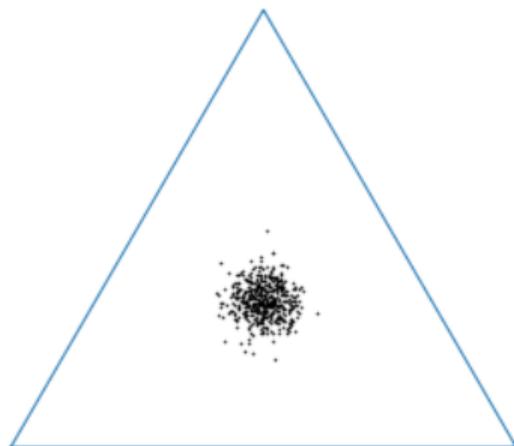


(a) $\{\mu^{(m)}\}_{m=1}^M$

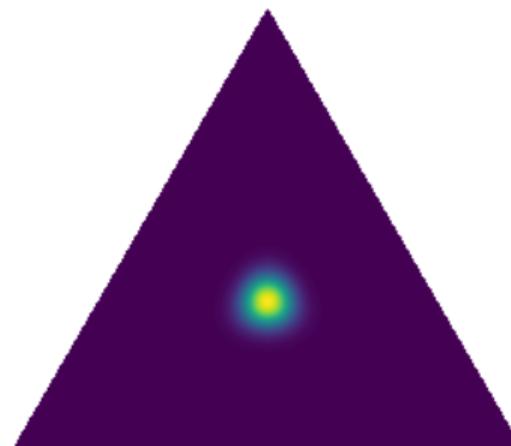


(b) $p(\mu|x^*, \mathcal{D})$

Distribution over Distributions

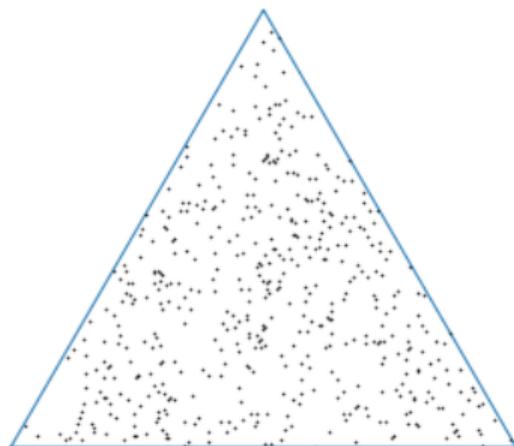


(a) $\{\mu^{(m)}\}_{m=1}^M$

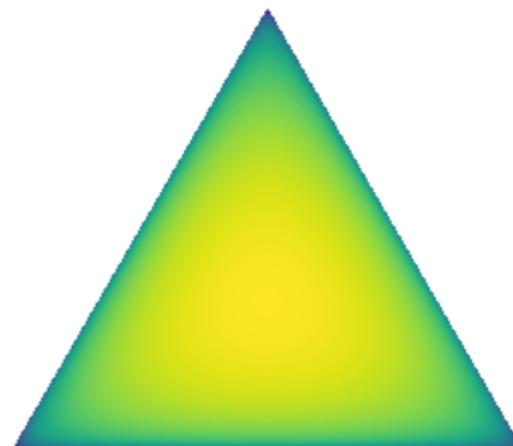


(b) $p(\mu|x^*, \mathcal{D})$

Distribution over Distributions



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



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

- **Explicitly** model $p(\mu|x^*, \mathcal{D})$ using a **Prior Network** $p(\mu|x^*; \hat{\theta})$

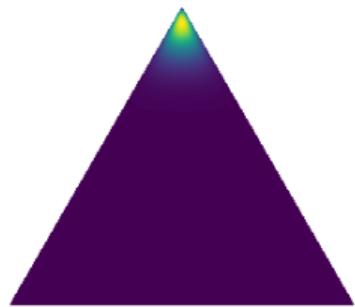
$$p(\mu|x^*; \hat{\theta}) \approx p(\mu|x^*, \mathcal{D})$$

- Predictive posterior distribution is given by expected categorical

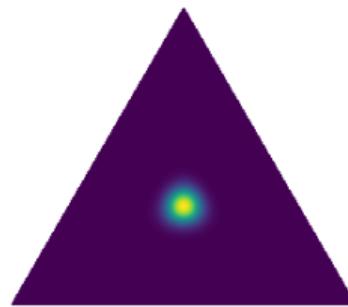
$$P(y|x^*; \hat{\theta}) = \mathbb{E}_{p(\mu|x^*; \hat{\theta})} [p(y|\mu)] = \hat{\mu}$$

Prior Networks [Malinin and Gales, 2018]

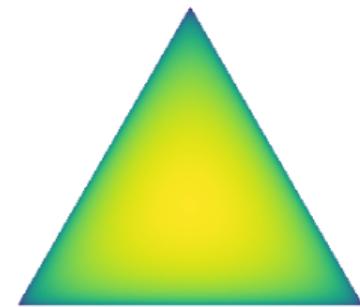
- Construct $p(\mu|x^*, \hat{\theta})$ to emulate ensemble



(a) Certain

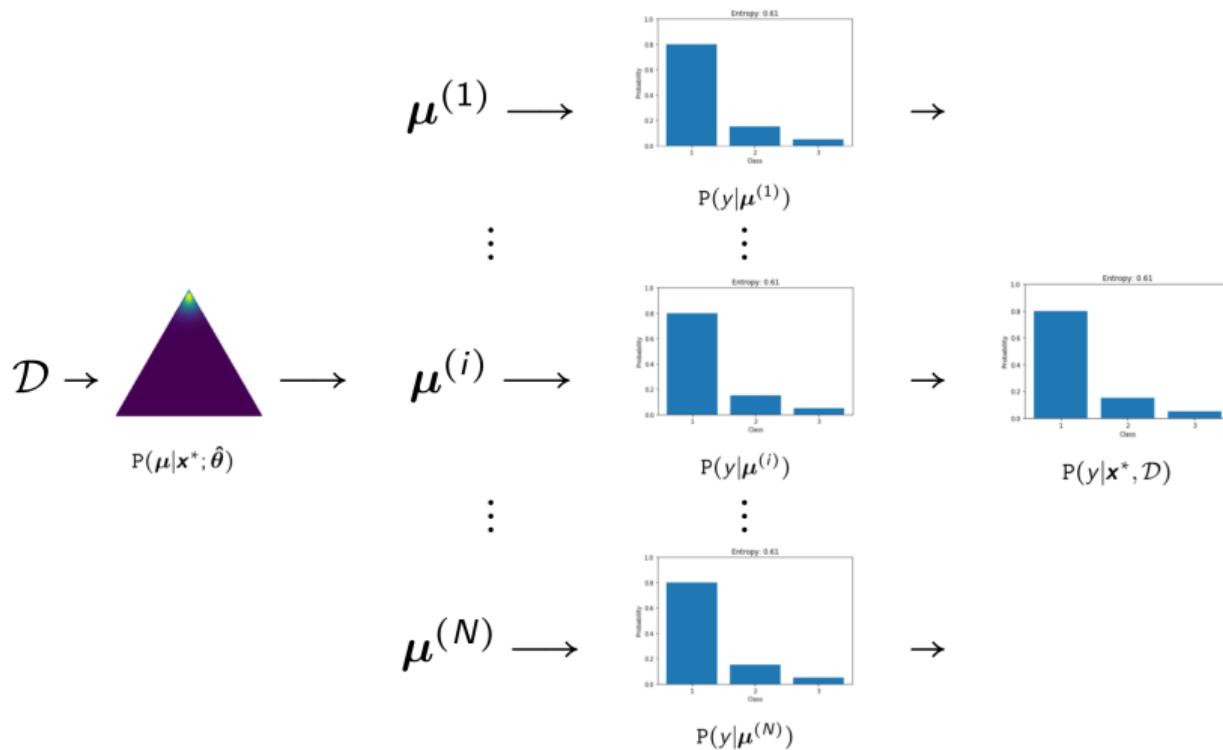


(b) Data Uncertainty

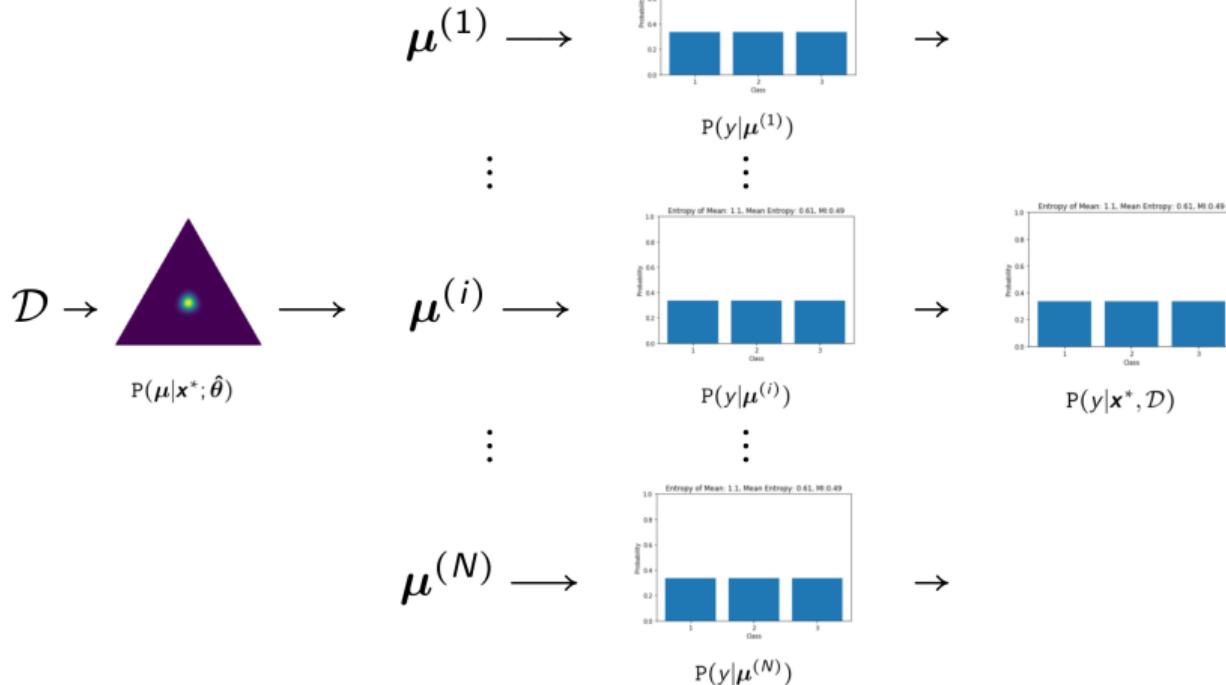


(c) Knowledge Uncertainty

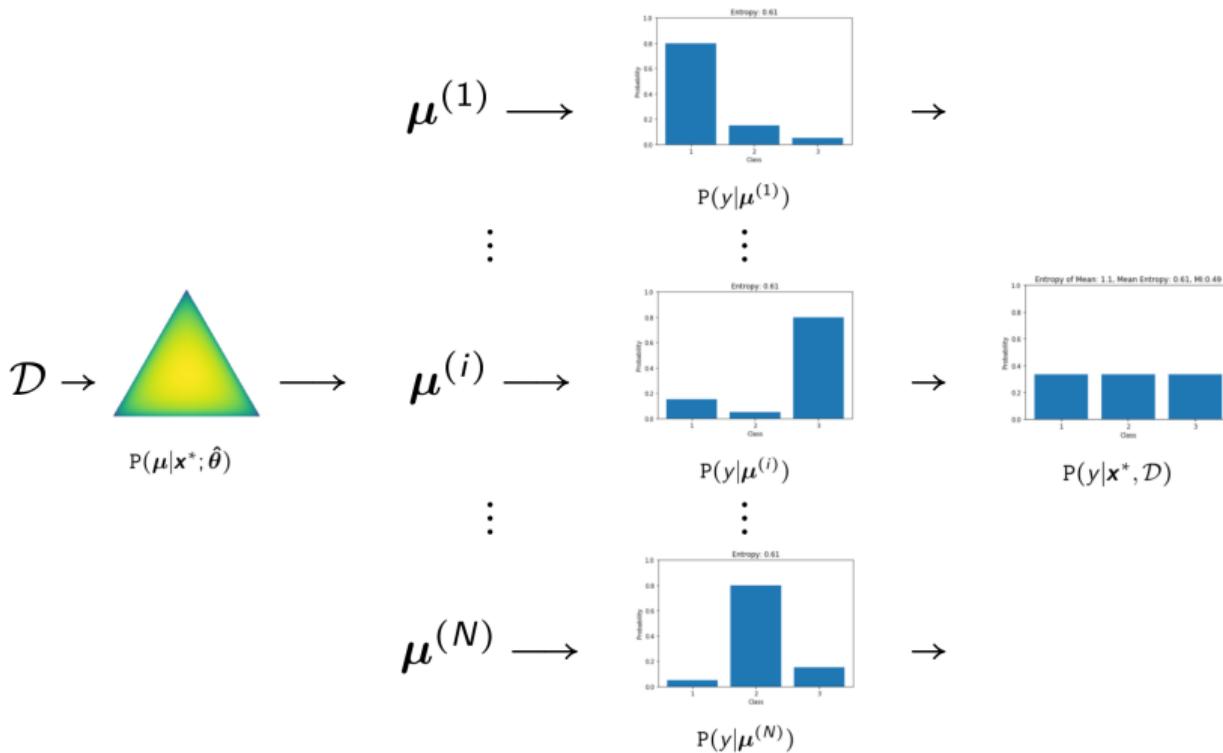
Distributions over Distributions via Prior Networks



Distributions over Distributions via Prior Networks



Distributions over Distributions via Prior Networks



- Ensemble uncertainty decomposition:

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

- Prior Network uncertainty decomposition

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

Prior Networks vs. Ensembles [Malinin and Gales, 2018]

- 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}

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

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

- Dirichlet Distribution → Distribution over simplex
 - Conjugate prior to categorical distribution
 - Convenient properties → analytically tractable

Reminder - Dirichlet Distribution

- Dirichlet is a distribution over categorical distributions

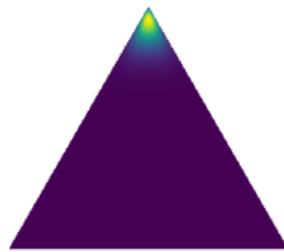
$$\text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\prod_{c=1}^K \Gamma(\alpha_c)} \prod_{c=1}^K \pi_c^{\alpha_c-1}; \quad \alpha_0 = \sum_{c=1}^K \alpha_c$$

- Parameterised by **concentration** parameters: $\boldsymbol{\alpha}, \alpha_c > 0$
- Expected label posteriors given by

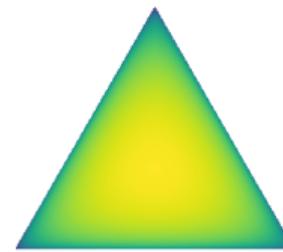
$$\hat{P}(y = \omega_c) = \hat{\mu}_c = \frac{\alpha_c}{\sum_{k=1}^K \alpha_k}$$

Prior Network Construction [Malinin and Gales, 2018]

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



(a) In-Domain Target



(b) OOD Target

Target Concentration Parameters [Malinin and Gales, 2018]

- To train the prior network we need a target distribution $p(\boldsymbol{\mu}|\boldsymbol{\beta})$ for $\mathbf{x}^{(i)}$
 - We **want** training data $\{\boldsymbol{\beta}^{(i)}, \mathbf{x}^{(i)}\}_{i=1}^N$
 - ... but **have** training data $\{y^{(i)}, \mathbf{x}^{(i)}\}_{i=1}^N$, where $y^{(i)} \in \{\omega_1, \dots, \omega_K\}$
- Solution → specify concentration parameters $\boldsymbol{\beta}^{(c)}$ as a function of target class y

$$p(\boldsymbol{\mu}|\boldsymbol{\beta}^{(c)}) = \text{Dir}(\boldsymbol{\mu}|\boldsymbol{\beta}^{(c)})$$

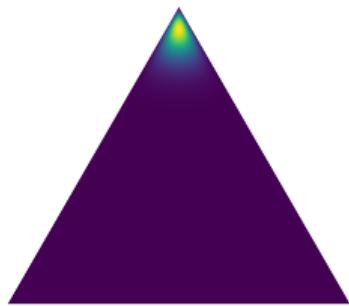
- Need $\boldsymbol{\beta}^{(c)}$ to yield **correct class**
- Need $\boldsymbol{\beta}^{(c)}$ to reflect “confidence” in sample
- $\beta_k > 0 \forall k$

- Consider setting $\beta^{(c)}$ as follows →

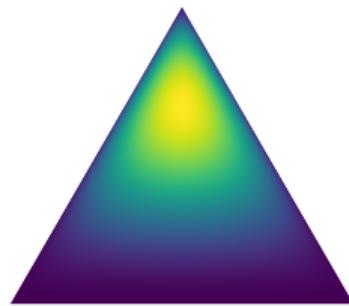
$$\beta_k^{(c)} = \begin{cases} \beta + 1 & \text{if } c = k \\ 1 & \text{if } c \neq k \end{cases}$$

- If β is large →
 - Sharp Dirichlet at corner of simplex corresponding to target class.
- If β is low →
 - Wide Dirichlet with the mode near the corner corresponding to target class.
- If β is zero →
 - Flat (uniform) Dirichlet distribution.

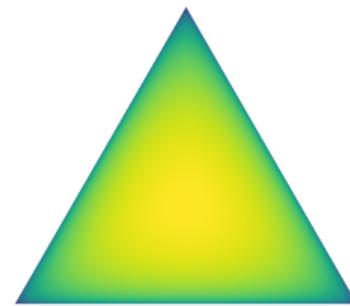
Target Concentration Parameters [Malinin and Gales, 2018]



(c) $\beta = 30$



(d) $\beta = 2$



(e) $\beta = 0$

- We can consider two loss functions - *Forward KL-Divergence* →

$$\mathcal{L}^{KL}(y, \mathbf{x}, \boldsymbol{\theta}; \beta) = \sum_{c=1}^K \mathcal{I}(y = \omega_c) \cdot \text{KL}[\text{Dir}(\boldsymbol{\mu}|\boldsymbol{\beta}^{(c)}) || \mathbf{p}(\boldsymbol{\mu}|\mathbf{x}; \boldsymbol{\theta})]$$

- ... or *reverse KL-Divergence* →

$$\mathcal{L}^{RKL}(y, \mathbf{x}, \boldsymbol{\theta}; \beta) = \sum_{c=1}^K \mathcal{I}(y = \omega_c) \cdot \text{KL}[\mathbf{p}(\boldsymbol{\mu}|\mathbf{x}; \boldsymbol{\theta}) || \text{Dir}(\boldsymbol{\mu}|\boldsymbol{\beta}^{(c)})]$$

Reminder - Kullback-Leibler Divergence

- Standard “measure” between distributions

$$\text{KL}[P_{\text{tr}}(y|x) || P(y|x; \theta)] = \mathbb{E}_{P_{\text{tr}}(y|x)} [\ln P_{\text{tr}}(y|x) - \ln P(y|x; \theta)]$$

- Variational optimization often yields reverse KL for training

$$\text{KL}[P(y|x; \theta) || P_{\text{tr}}(y|x)] = \mathbb{E}_{P(y|x; \theta)} [\ln P(y|x; \theta) - \ln P_{\text{tr}}(y|x)]$$

- Measures have different properties →
 - Forward KL is **zero-avoiding**
 - Reverse KL is **zero-forcing**

Forward KL-divergence Loss [Malinin and Gales, 2018]

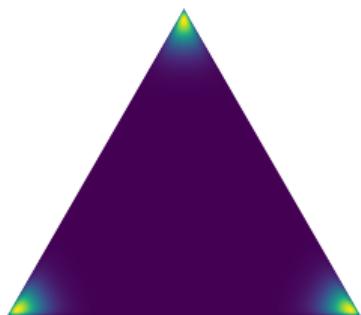
- Consider expectation of *forward* KL-div loss wrt. empirical distribution $\hat{p}(\mathbf{x}, y) \rightarrow$

$$\begin{aligned}\mathcal{L}^{KL}(\boldsymbol{\theta}; \beta) &= \mathbb{E}_{\hat{p}_{\text{tr}}(\mathbf{x}, y)} \left[\sum_{c=1}^K \mathcal{I}(y = \omega_c) \cdot \text{KL}[\text{Dir}(\boldsymbol{\mu} | \boldsymbol{\beta}^{(c)}) || p(\boldsymbol{\mu} | \mathbf{x}; \boldsymbol{\theta})] \right] \\ &= \mathbb{E}_{\hat{p}_{\text{tr}}(\mathbf{x})} \left[\sum_{c=1}^K \mathbb{E}_{\hat{p}_{\text{tr}}(y | \mathbf{x})} [\mathcal{I}(y = \omega_c)] \cdot \text{KL}[\text{Dir}(\boldsymbol{\mu} | \boldsymbol{\beta}^{(c)}) || p(\boldsymbol{\mu} | \mathbf{x}; \boldsymbol{\theta})] \right] \\ &= \mathbb{E}_{\hat{p}_{\text{tr}}(\mathbf{x})} \left[\text{KL} \left[\sum_{c=1}^K P_{\text{tr}}(y = \omega_c) \cdot \text{Dir}(\boldsymbol{\mu} | \boldsymbol{\beta}^{(c)}) || p(\boldsymbol{\mu} | \mathbf{x}; \boldsymbol{\theta}) \right] \right] + C\end{aligned}$$

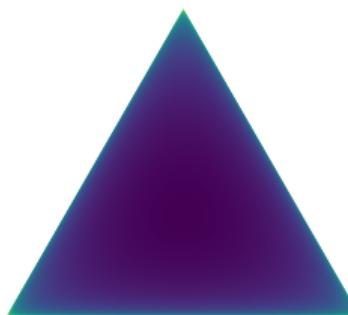
- Target distribution becomes a **mixture of Dirichlets!**

Forward KL-divergence Loss [Malinin and Gales, 2018]

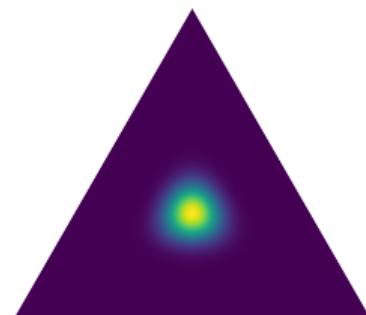
- Target distribution becomes a **mixture of Dirichlets**!



(a) Induced Target



(b) Model



(c) Want

- Forward KL-divergence is zero avoiding → Model will try to cover **each mode**!
 - Leads to **undesired behaviour** → bad performance!
 - **Doesn't scale** to datasets with more than **10 classes**!

Reverse KL-divergence Loss [Malinin and Gales, 2019]

- Consider expectation of reverse KL-div loss wrt. empirical distribution $\hat{p}(\mathbf{x}, y) \rightarrow$

$$\begin{aligned}\mathcal{L}^{RKL}(\boldsymbol{\theta}; \beta) &= \mathbb{E}_{\hat{p}_{\text{tr}}(\mathbf{x})} \left[\sum_{c=1}^K \hat{P}_{\text{tr}}(y = \omega_c | \mathbf{x}) \text{KL}[\mathbf{p}(\boldsymbol{\mu} | \mathbf{x}; \boldsymbol{\theta}) || \text{Dir}(\boldsymbol{\mu} | \boldsymbol{\beta}^{(c)})] \right] \\ &= \mathbb{E}_{\hat{p}_{\text{tr}}(\mathbf{x})} \left[\mathbb{E}_{\mathbf{p}(\boldsymbol{\mu} | \mathbf{x}; \boldsymbol{\theta})} \left[\ln \mathbf{p}(\boldsymbol{\mu} | \mathbf{x}; \boldsymbol{\theta}) - \ln \prod_{c=1}^K \text{Dir}(\boldsymbol{\mu} | \boldsymbol{\beta}^{(c)})^{\hat{P}_{\text{tr}}(y = \omega_c | \mathbf{x})} \right] \right] \\ &= \mathbb{E}_{\hat{p}_{\text{tr}}(\mathbf{x})} \left[\text{KL}[\mathbf{p}(\boldsymbol{\mu} | \mathbf{x}; \boldsymbol{\theta}) || \text{Dir}(\boldsymbol{\mu} | \sum_{c=1}^K \hat{P}_{\text{tr}}(y = \omega_c | \mathbf{x}) \cdot \boldsymbol{\beta}^{(c)})] \right] + C\end{aligned}$$

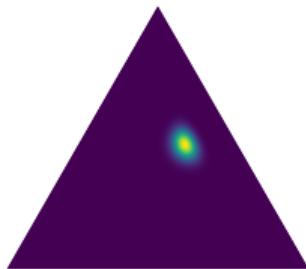
- Expectation induces product of target Dirichlet distributions.

Reverse KL-divergence Loss [Malinin and Gales, 2019]

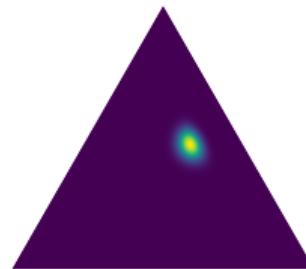
- Expectation induces product of target Dirichlet distributions

$$\mathcal{L}^{RKL}(\theta; \beta) = \mathbb{E}_{\hat{p}_{\text{tr}}(\mathbf{x})} \left[\text{KL}[\mathbf{p}(\mu|\mathbf{x}; \theta) || \text{Dir}(\mu | \sum_{c=1}^K \hat{P}_{\text{tr}}(y = \omega_c | \mathbf{x}) \cdot \beta^{(c)})] \right] + C$$

- Target becomes a **uni-modal Dirichlet distribution** at appropriate location!



(a) Induced Target

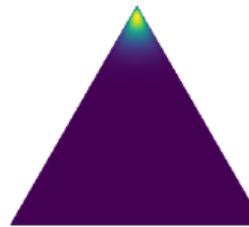


(b) Model

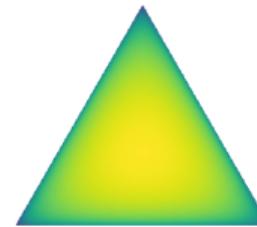
Prior Network Construction

- Reverse KL loss $\mathcal{L}^{RKL}(\theta, \mathcal{D}; \beta) \rightarrow$ full control over behaviour of model!

$$\mathcal{L}(\theta, \mathcal{D}; \beta_{in}, \beta_{out}, \gamma) = \mathcal{L}_{in}^{RKL}(\theta, \mathcal{D}_{trn}; \beta_{in}) + \gamma \cdot \mathcal{L}_{out}^{RKL}(\theta, \mathcal{D}_{out}; \beta_{out})$$



(a) In-Domain Target

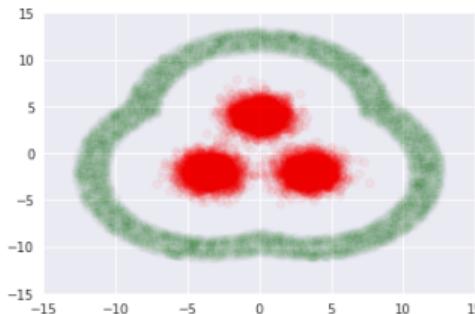


(b) OOD Target

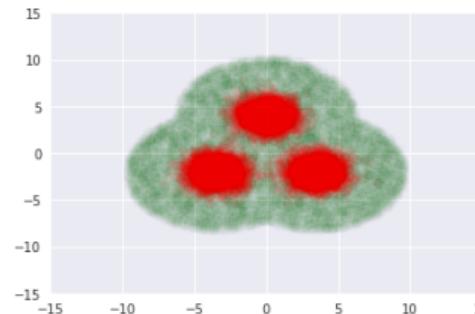
- But how to obtain out-of-domain training data \mathcal{D}_{out} ?
 - Use a different dataset or **adversarial attacks**

Prior Network Construction

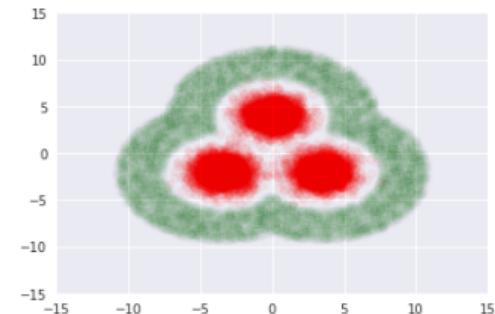
- Out-of-domain (OOD) training data must be on *boundary* on in-domain region →
 - Too loose → Some OOD might be considered in-domain
 - Too tight → Some in-domain might be considered OOD



(a) Too Loose

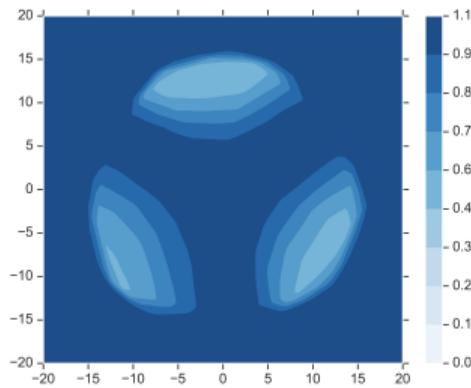


(b) Too Tight

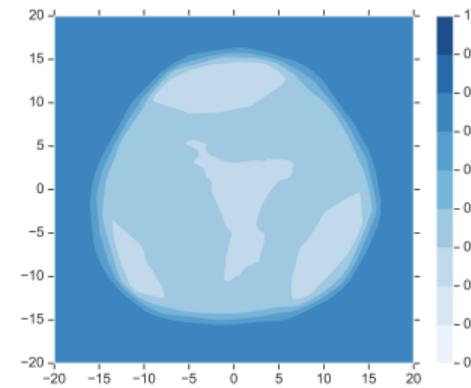


(c) Good

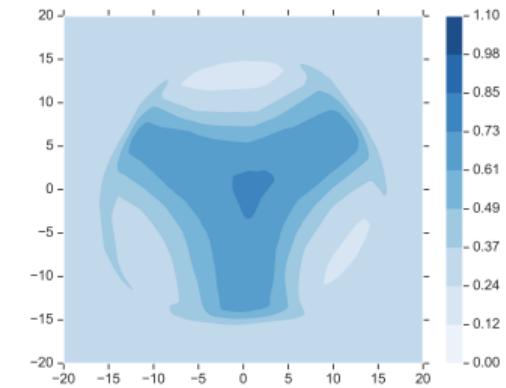
Prior Networks trained with *forward* KL-divergence loss on Artificial Data



(a) Total Uncertainty

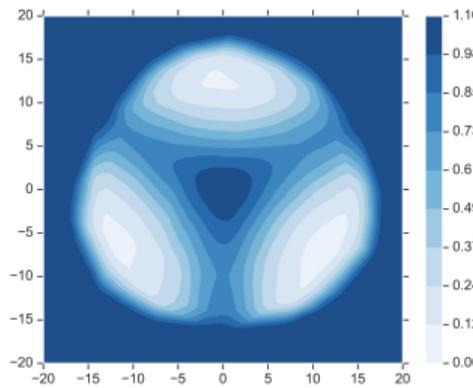


(b) Data Uncertainty

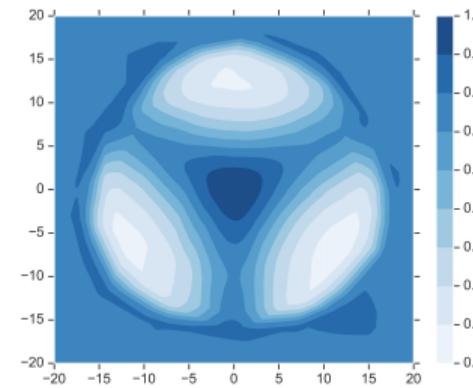


(c) Knowledge Uncertainty

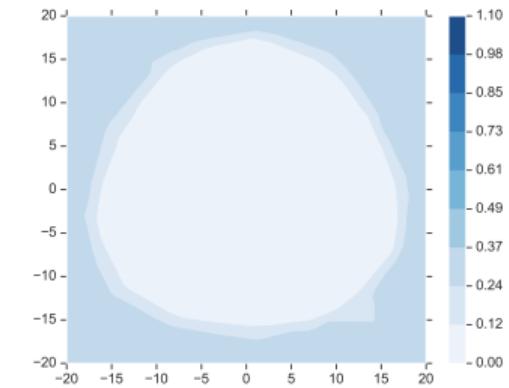
Prior Networks trained with *reverse KL-divergence loss* on Artificial Data



(a) Total Uncertainty



(b) Data Uncertainty



(c) Knowledge Uncertainty

Classification Error Rate

Dataset	DNN	PN-KL	PN-RKL	Ensemble
CIFAR-10	8.0	14.7	7.5	6.6
CIFAR-100	30.4	-	28.1	26.9
TinyImageNet	41.7	-	40.3	36.9

Out-of-Distribution Detection

Model	CIFAR-10/CIFAR-100			CIFAR-100/TinyImageNet		
	SVHN	LSUN	TinyImageNet	SVHN	LSUN	CIFAR-10
Ensemble	89.5	93.2	90.3	78.9	85.6	76.5
PN-KL	97.8	91.6	92.4	-	-	-
PN-RKL	98.2	95.7	95.7	84.8	100.0	57.8

Table: Out-of-domain detection results (mean % AUROC across 10 rand. inits).

Overview of the Talk

1. Context: Why do we need Uncertainty Estimation?
2. Sources of Uncertainty in Predictions
3. Ensemble Approaches
4. Prior Networks
5. **Adversarial Attack Detection**

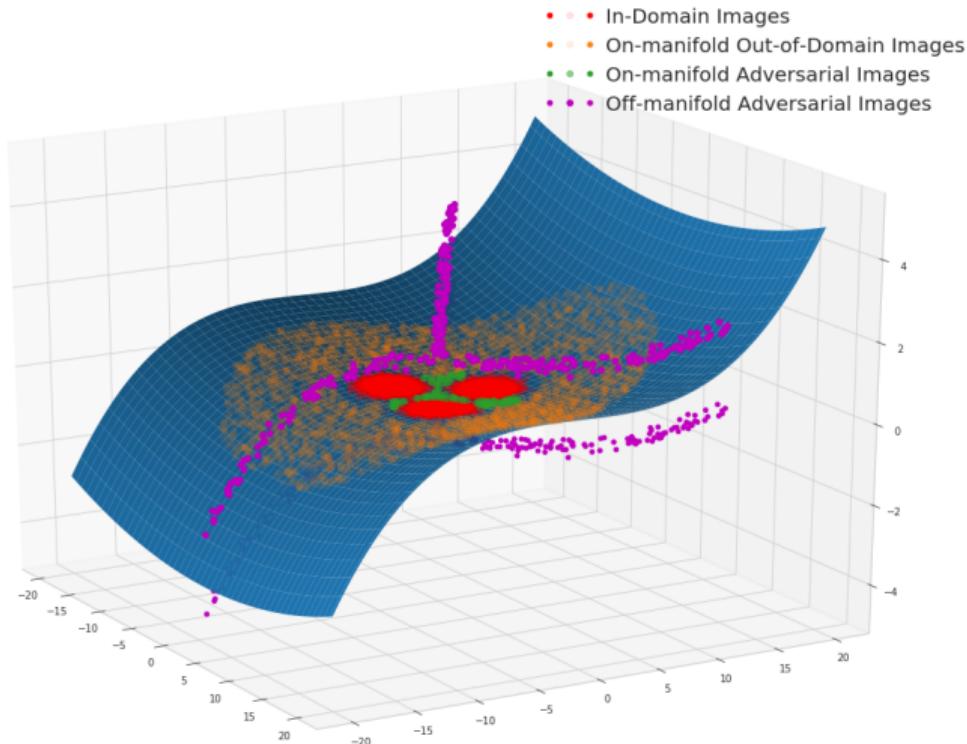
- Adversarial Attacks → Small perturbation of the input x^* which affects prediction
 - Exist for **many modalities** → images, text and audio
 - **Transferable** between models
 - Can be deployed in **real world**.
- Adversarial Attacks are a **security concern!**
 - Detect using measures of **uncertainty?**

- Adversarial attacks: generate sample $\tilde{\mathbf{x}}$:
 1. swaps to target class $\tilde{\omega}$
 2. requires minimum changes to original sample \mathbf{x}
- Requirements expressed as

$$\mathcal{A}_{\text{adv}}(\mathbf{x}, \tilde{\omega}) = \arg \min_{\tilde{\mathbf{x}} \in \mathcal{R}^D} \left\{ \mathcal{L}(y = \tilde{\omega}, \tilde{\mathbf{x}}, \hat{\theta}) \right\} : \delta(\mathbf{x}, \tilde{\mathbf{x}}) < \epsilon$$

- ϵ is number of swapped bits (for images)

Manifold Interpretation of Adversarial Attacks



- Consider an **uncertainty** based detection scheme:

$$\hat{\mathcal{I}}_T(\mathbf{x}) = \begin{cases} 1, & T > \mathcal{H}(\mathbf{x}) \\ 0, & T \leq \mathcal{H}(\mathbf{x}) \\ 0, & \mathbf{x} = \emptyset \end{cases}$$

- Successful attacks are able to :
 - Both **affect prediction** and **avoid detection**.
- Can assess using false positive and true positives:

$$t_p(T) = \frac{1}{N} \sum_{i=1}^N \mathcal{I}_T(\mathbf{x}_i), \quad f_p(T) = \frac{1}{N} \sum_{i=1}^N \mathcal{I}_T(\mathcal{A}_{\text{adv}}(\mathbf{x}_i, \omega_t))$$

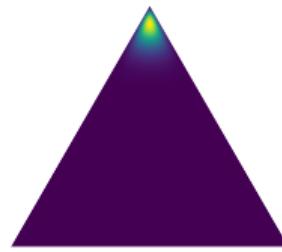
- Joint Success Rate** is where $t_p(T) = f_p(T)$

Adversarial Attack Detection via Prior Networks

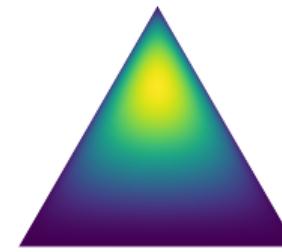
- Prior Networks yield rich measures of uncertainty
- Greatly constrain space of successful *adaptive* attacks
 - Confidence → constraint on max logit
 - Total Uncertainty → constraint on relative of magnitude of logits
 - Knowledge Uncertainty → constraint on relative and absolute magnitude of logits
 - .. and attack must also affect predicted class!
- Explicit control behaviour via training data →
 - Further constrain space of successful adversarial solutions

Prior Network Adversarial Training

$$\mathcal{L}(\theta, \mathcal{D}; \beta_{\text{nat}}, \beta_{\text{adv}}) = \mathcal{L}_{\text{nat}}^{RKL}(\theta, \mathcal{D}_{\text{trn}}; \beta_{\text{nat}} = 1e2) + \gamma \cdot \mathcal{L}_{\text{adv}}^{RKL}(\theta, \mathcal{D}_{\text{adv}}; \beta_{\text{adv}} = 1)$$



(a) Natural Target



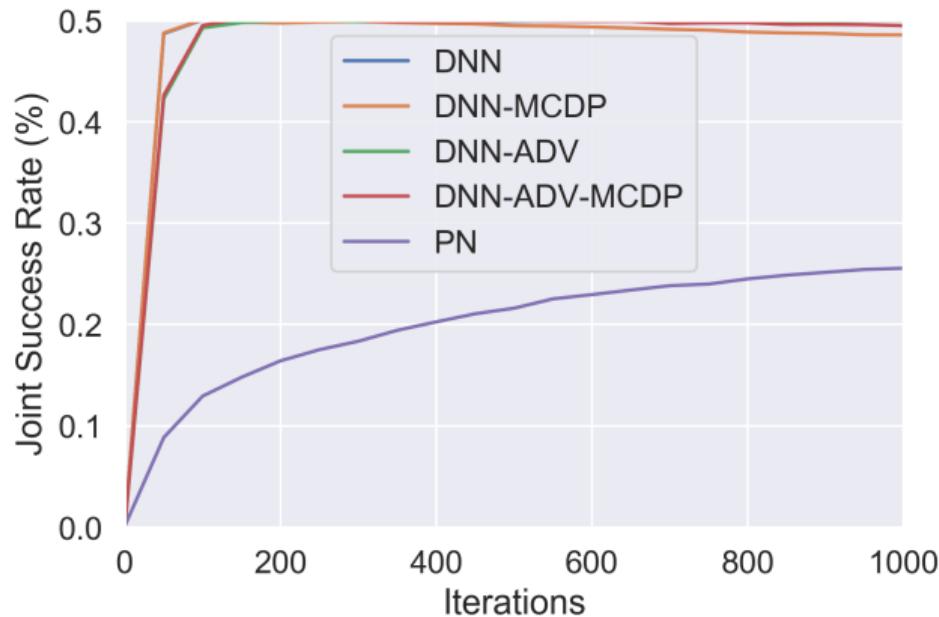
(b) Adversarial Target

- Standard adversarial training →
 - Correct Prediction for adversarial inputs
- Prior Network adversarial training →
 - Correct Prediction + High Uncertainty for adversarial inputs

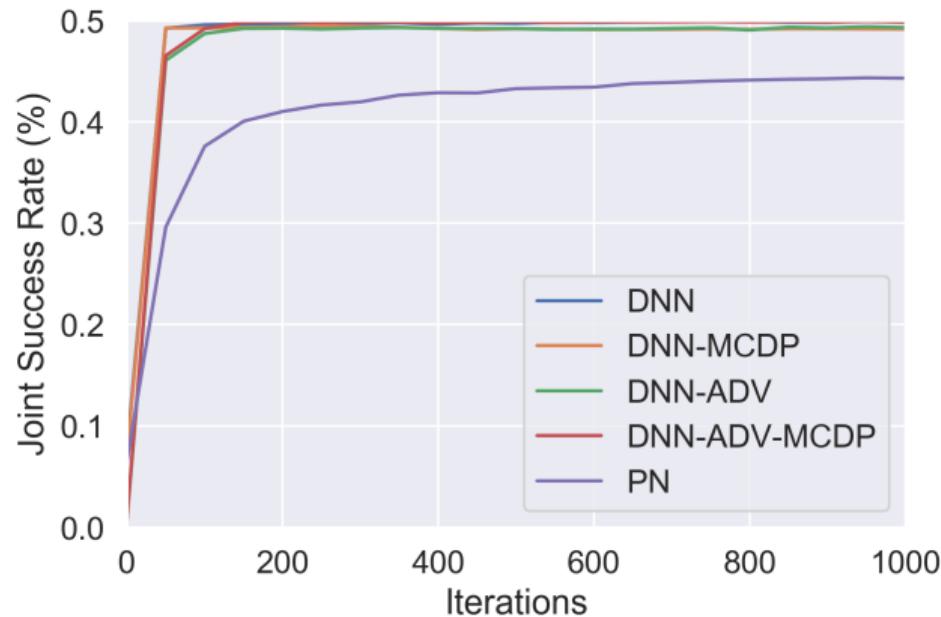
Experiments

- Baselines: MC-Dropout and Standard Adversarial Training
- Prior Network adversarial training → 1-step FGSM attacks
- Evaluation Attack → strong *adaptive* whitebox PGD-MIM attack
- Datasets: CIFAR10 and CIFAR100

Joint Success Rate - CIFAR10



Joint Success Rate - CIFAR100



Thank You!

Any questions?

References I

- [Blundell et al., 2015] Blundell, C., Cornebise, J., Kavukcuoglu, K., and Wierstra, D. (2015).
Weight uncertainty in neural networks.
arXiv preprint arXiv:1505.05424.
- [Gal, 2016] Gal, Y. (2016).
Uncertainty in Deep Learning.
PhD thesis, University of Cambridge.
- [Gal and Ghahramani, 2016] Gal, Y. and Ghahramani, Z. (2016).
Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning.
In *Proc. 33rd International Conference on Machine Learning (ICML-16)*.

References II

- [Garipov et al., 2018] Garipov, T., Izmailov, P., Podoprikhin, D., Vetrov, D. P., and Wilson, A. G. (2018).
Loss surfaces, mode connectivity, and fast ensembling of dnns.
In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems 31*, pages 8789–8798. Curran Associates, Inc.
- [Hendrycks and Gimpel, 2016] Hendrycks, D. and Gimpel, K. (2016).
A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks.
<http://arxiv.org/abs/1610.02136>.
arXiv:1610.02136.

References III

[Hernández-Lobato and Adams, 2015] Hernández-Lobato, J. M. and Adams, R. (2015).

Probabilistic backpropagation for scalable learning of bayesian neural networks.
In *International Conference on Machine Learning*, pages 1861–1869.

[Lakshminarayanan et al., 2017] Lakshminarayanan, B., Pritzel, A., and Blundell, C. (2017).

Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles.
In *Proc. Conference on Neural Information Processing Systems (NIPS)*.

[Maddox et al., 2019] Maddox, W., Garipov, T., Izmailov, P., Vetrov, D. P., and Wilson, A. G. (2019).

A simple baseline for bayesian uncertainty in deep learning.
CoRR, abs/1902.02476.

References IV

- [Malinin and Gales, 2018] Malinin, A. and Gales, M. (2018).
Predictive uncertainty estimation via prior networks.
In *Advances in Neural Information Processing Systems*, pages 7047–7058.
- [Malinin and Gales, 2019] Malinin, A. and Gales, M. (2019).
Reverse kl-divergence training of prior networks: Improved uncertainty and
adversarial robustness.
arXiv preprint arXiv:1905.13472.
- [Osband et al., 2016] Osband, I., Blundell, C., Pritzel, A., and Van Roy, B. (2016).
Deep exploration via bootstrapped dqn.
In Lee, D. D., Sugiyama, M., Luxburg, U. V., Guyon, I., and Garnett, R., editors,
Advances in Neural Information Processing Systems 29, pages 4026–4034. Curran
Associates, Inc.

[Welling and Teh, 2011] Welling, M. and Teh, Y. W. (2011).
Bayesian Learning via Stochastic Gradient Langevin Dynamics.
In *Proc. International Conference on Machine Learning (ICML)*.