

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

Andrey Malinin and Mark Gales

18 October 2019



(a) Prof. Mark Gales



- 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 \rightarrow
 - 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 \rightarrow
 - 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



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

Sources of Uncertainty



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

Data (Aleatoric) Uncertainty



WINIVERSITY OF Yandex Research

Data Uncertainty



CAMBRIDGE Yandex Research

Data Uncertainty

• Distinct Classes



• Overlapping Classes





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

- Data Uncertainty is the entropy of the true data distribution ightarrow

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

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

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

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

Sources of Uncertainty



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

Knowledge Uncertainty



VINIVERSITY OF Vandex Research

Knowledge Uncertainty - Out-of-Distribution

• Unseen classes

• Unseen variations of seen classes





- Data Uncertainty \rightarrow Known-Unknown
 - Class overlap (complexity of decision boundaries)
 - Human labelling error
- Knowledge Uncertainty \rightarrow 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 heta captured by model posterior $\mathrm{p}(heta|\mathcal{D}) o$

$$\mathtt{p}(oldsymbol{ heta} | \mathcal{D}) = rac{\mathtt{p}(\mathcal{D} | oldsymbol{ heta}) \mathtt{p}(oldsymbol{ heta})}{\mathtt{p}(\mathcal{D})}$$

• Bayesian inference of P $(y|m{x}^*,m{ heta})
ightarrow$

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

- Can consider an ensemble of models \rightarrow

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

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

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

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

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

Expected Data Uncertainty

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

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

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

Model Uncertainty

• If the predictions from the models are consistent

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

• If the predictions from the models are diverse

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

• Difference of the two is a measure of model uncertainty

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

Model Uncertainty \rightarrow Knowledge Uncertainty

• If the predictions from the models are consistent

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

• If the predictions from the models are diverse

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

• Difference of the two is a measure of knowledge uncertainty

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

Ensemble for certain in-domain input





Ensemble for uncertain in-domain input





Ensemble for Out-of-Domain input





• Ensemble $\{P(y|\boldsymbol{x}^*, \boldsymbol{\theta}^{(m)})\}_{m=1}^M$ can be visualized on a simplex for an input \boldsymbol{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 \rightarrow use approximate inference
 - Compute approximate posterior $\mathrm{q}(m{ heta}) pprox \mathrm{p}(m{ heta} | \mathcal{D})$
 - 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]

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

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



- (a) Confident (b) Data Uncertainty (c) Knowledge Uncertainty
- Same as sampling from implicit Distribution over output Distributions

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

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




Distribution over Distributions





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

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



Distribution over Distributions





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

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



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

$$\operatorname{p}(oldsymbol{\mu}|oldsymbol{x}^*; oldsymbol{\hat{ heta}}) pprox \operatorname{p}(oldsymbol{\mu}|oldsymbol{x}^*, \mathcal{D})$$

• Predictive posterior distribution is given by expected categorical

$$\mathbb{P}(y|oldsymbol{x}^*; oldsymbol{\hat{ heta}}) = \mathbb{E}_{\mathbb{P}(oldsymbol{\mu}|oldsymbol{x}^*; oldsymbol{\hat{ heta}})} ig[\mathbb{P}(y|oldsymbol{\mu}) ig] = oldsymbol{\hat{eta}}$$

Prior Networks [Malinin and Gales, 2018]

• Construct $\mathrm{p}(\mu|\mathbf{x}^*, \hat{\mathbf{ heta}})$ to emulate ensemble





Distributions over Distributions via Prior Networks





Distributions over Distributions via Prior Networks





Distributions over Distributions via Prior Networks





Uncertainty Measures for Prior Networks [Malinin and Gales, 2018]

• Ensemble uncertainty decomposition:

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

• Prior Network uncertainty decomposition

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

Λ

- 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 $\ensuremath{\mathcal{D}}$

• A Prior Network parametrizes the Dirichlet Distribution

$$\mathrm{p}(\mu|m{x}^*; m{\hat{ heta}}) = \mathtt{Dir}(\mu|m{lpha}), \quad m{lpha} = m{f}(m{x}^*; m{\hat{ heta}})$$

- Dirichlet Distribution \rightarrow Distribution over simplex
 - · Conjugate prior to categorical distribution
 - Convenient properties \rightarrow analytically tractable

• Dirichlet is a distribution over categorical distributions

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

- Parameterised by concentration parameters: α , $\alpha_c > 0$
- Expected label posteriors given by

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

Prior Network Construction [Malinin and Gales, 2018]





Target Concentration Parameters [Malinin and Gales, 2018]

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

$$\operatorname{p}(\mu|eta^{(c)})=\operatorname{Dir}(\mu|eta^{(c)})$$

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

Target Concentration Parameters [Malinin and Gales, 2018]

- Consider setting $eta^{(c)}$ as follows ightarrow

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

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

Target Concentration Parameters [Malinin and Gales, 2018]





- We can consider two loss functions - Forward KL-Divergence ightarrow

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

• ... or reverse KL-Divergence \rightarrow

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

• Standard "measure" between distributions

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

· Variational optimization often yields reverse KL for training

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

- Measures have different properties ightarrow
 - 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 $\widehat{p}(\pmb{x}, y) \rightarrow$

$$\begin{split} \mathcal{L}^{KL}(\boldsymbol{\theta};\beta) &= \mathbb{E}_{\tilde{p}_{tr}(\boldsymbol{X},y)} \Big[\sum_{c=1}^{K} \mathcal{I}(y = \omega_{c}) \cdot \mathrm{KL}[\mathrm{Dir}(\boldsymbol{\mu}|\boldsymbol{\beta}^{(c)})||p(\boldsymbol{\mu}|\boldsymbol{x};\boldsymbol{\theta})] \Big] \\ &= \mathbb{E}_{\tilde{p}_{tr}(\boldsymbol{X})} \Big[\sum_{c=1}^{K} \mathbb{E}_{\tilde{p}_{tr}(y|\boldsymbol{X})} [\mathcal{I}(y = \omega_{c})] \cdot \mathrm{KL}[\mathrm{Dir}(\boldsymbol{\mu}|\boldsymbol{\beta}^{(c)})||p(\boldsymbol{\mu}|\boldsymbol{x};\boldsymbol{\theta})] \Big] \\ &= \mathbb{E}_{\tilde{p}_{tr}(\boldsymbol{X})} \Big[\mathrm{KL}[\sum_{c=1}^{K} \mathbb{P}_{tr}(y = \omega_{c}) \cdot \mathrm{Dir}(\boldsymbol{\mu}|\boldsymbol{\beta}^{(c)})||p(\boldsymbol{\mu}|\boldsymbol{x};\boldsymbol{\theta})] \Big] + C \end{split}$$

Target distribution becomes a mixture of Dirichlets!

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

• Target distribution becomes a mixture of Dirichlets!



- Forward KL-divergence is zero avoiding \rightarrow 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 $\widehat{p}(\pmb{x}, y) \rightarrow$

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

• Expectation induces product of target Dirichlet distributions.

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

Expectation induces product of target Dirichlet distributions

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

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





Prior Network Construction

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

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



- 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 ightarrow
 - Too loose \rightarrow Some OOD might be considered in-domain
 - Too tight \rightarrow Some in-domain might be considered OOD



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



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



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

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).

- 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 ightarrow Small perturbation of the input \pmb{x}^* which affects prediction
 - Exist for many modalities \rightarrow 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{x} :
 - **1.** swaps to target class $\tilde{\omega}$
 - 2. requires minimum changes to original sample x
- Requirements expressed as

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

• ϵ is number of swapped bits (for images)

Manifold Interpretation of Adversarial Attacks





Adversarial Attack Detection

• Consider an uncertainty based detection scheme:

$$\hat{\mathcal{I}}_{\mathcal{T}}(oldsymbol{x}) = egin{cases} 1, & \mathcal{T} > \mathcal{H}(oldsymbol{x}) \ 0, & \mathcal{T} \leq \mathcal{H}(oldsymbol{x}) \ 0, & oldsymbol{x} = \emptyset \end{cases}$$

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

$$t_{p}(T) = rac{1}{N}\sum_{i=1}^{N}\mathcal{I}_{T}(oldsymbol{x}_{i}), \quad f_{p}(T) = rac{1}{N}\sum_{i=1}^{N}\mathcal{I}_{T}(\mathcal{A}_{ extsf{adv}}(oldsymbol{x}_{i}, \omega_{t}))$$

• Joint Success Rate is where $t_{\rho}(T) = f_{\rho}(T)$

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

Prior Network Adversarial Training

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



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

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

Joint Success Rate - CIFAR10



WINIVERSITY OF Yandex Research
Joint Success Rate - CIFAR100



Thank You!

Any questions?



[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).

[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. [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. [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).