

Uncertainty in Deep Learning: From Bayesian Neural Networks to Dropout

Samar Hadou

Department of Electrical and Systems Engineering
University of Pennsylvania

April 14, 2022

This talk is based on Yarin Gal's PhD work:

- ▶ Y. Gal. **Uncertainty in Deep Learning**. PhD Thesis, University of Cambridge, 2016.
- ▶ Y. Gal and Z. Ghahramani. **Dropout as a Bayesian approximate: Representing Model Uncertainty in Deep Learning**. ICML, 2016.
- ▶ A. Kendall and Y. Gal. **What uncertainties do we need in Bayesian deep learning for computer vision?** NeurIPS, 2017.

Standard deep learning (CNNs, RNNs, ..., etc.)

- ▶ Imagine that we trained a CNN to classify **dog breeds**
- ▶ What happens when executed on this one?



Yarin's question:
Do we need to replace our standard models?

Bayesian deep learning

- ▶ Gives the ability to models to say "I am not sure!"
- ▶ Learns distributions over the weights to tell how likely a model fits the data, and
- ▶ Provides uncertainty estimates

Standard deep learning (CNNs, RNNs, ..., etc.)

- ▶ Imagine that we trained a CNN to classify **dog breeds**
- ▶ What happens when executed on this one?



Yarin's question:
Do we need to replace our standard models?

Bayesian deep learning

- ▶ Gives the ability to models to say "I am not sure!"
- ▶ Learns distributions over the weights to tell how likely a model fits the data, and
- ▶ Provides uncertainty estimates

Standard deep learning (CNNs, RNNs, ..., etc.)

- ▶ Imagine that we trained a CNN to classify **dog breeds**
- ▶ What happens when executed on this one?



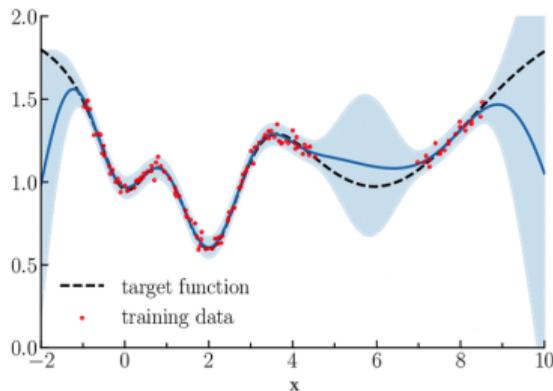
Yarin's question:
Do we need to replace our standard models?

Bayesian deep learning

- ▶ Gives the ability to models to say "I am not sure!"
- ▶ Learns distributions over the weights to tell how likely a model fits the data, and
- ▶ Provides uncertainty estimates

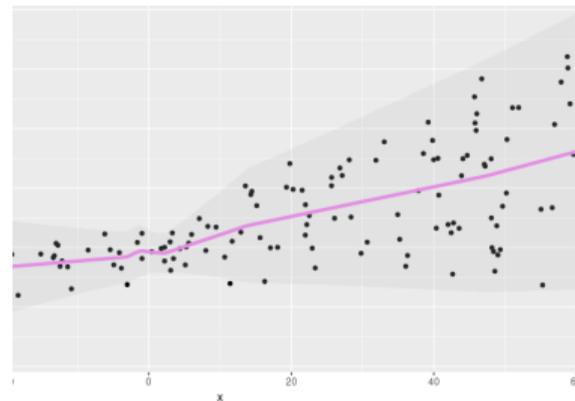
I. Epistemic Uncertainty

- Uncertainty that we can mitigate by adding more data



II. Aleatoric Uncertainty

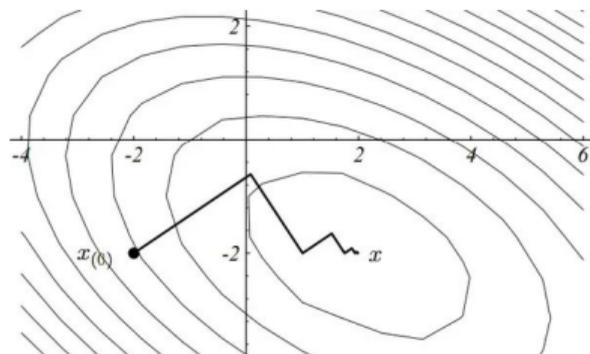
- Adding more data would not decrease uncertainty



Bayesian Deep Learning

Standard deep learning

- Finds a point estimate (model) that **minimizes** some loss function



Bayesian deep learning

- Evaluates predictive distributions via **marginalization**

$$p(\mathbf{y}|\mathbf{x}, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}|\mathbf{x}, \omega)p(\omega|\mathbf{X}, \mathbf{Y})d\omega$$

- Assumption: $p(\mathbf{y}|\mathbf{x}, \omega) = \mathcal{N}(\mathbf{y}; f^\omega(\mathbf{x}), \tau^{-1}\mathbf{I})$, with τ is said to be the model precision. (We'll revisit this assumption.)

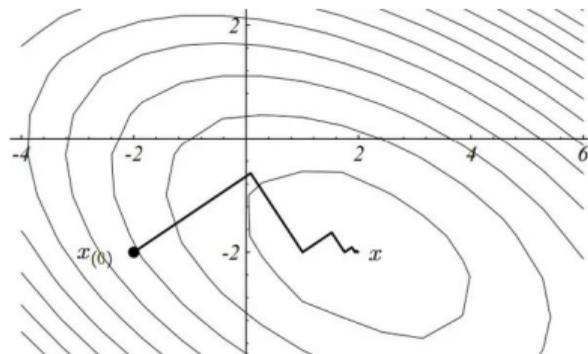
- However, the posterior distribution

$$p(\omega|\mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{Y}|\mathbf{X}, \omega)p(\omega)}{p(\mathbf{Y}|\mathbf{X})}$$

is intractable.

Standard deep learning

- Finds a point estimate (model) that **minimizes** some loss function



Bayesian deep learning

- Evaluates predictive distributions via **marginalization**

$$p(\mathbf{y}|\mathbf{x}, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\omega})p(\boldsymbol{\omega}|\mathbf{X}, \mathbf{Y})d\boldsymbol{\omega}$$

- Assumption: $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\omega}) = \mathcal{N}(\mathbf{y}; f^{\boldsymbol{\omega}}(\mathbf{x}), \tau^{-1}\mathbf{I})$, with τ is said to be the model precision. (We'll revisit this assumption.)

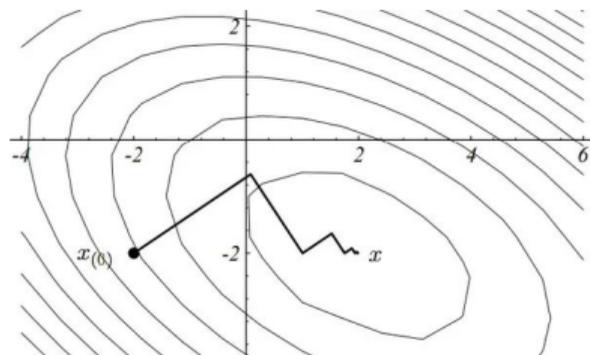
- However, the posterior distribution

$$p(\boldsymbol{\omega}|\mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega})p(\boldsymbol{\omega})}{p(\mathbf{Y}|\mathbf{X})}$$

is intractable.

Standard deep learning

- Finds a point estimate (model) that **minimizes** some loss function



Bayesian deep learning

- Evaluates predictive distributions via **marginalization**

$$p(\mathbf{y}|\mathbf{x}, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\omega})p(\boldsymbol{\omega}|\mathbf{X}, \mathbf{Y})d\boldsymbol{\omega}$$

- Assumption: $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\omega}) = \mathcal{N}(\mathbf{y}; f^{\boldsymbol{\omega}}(\mathbf{x}), \tau^{-1}\mathbf{I})$, with τ is said to be the model precision. (We'll revisit this assumption.)

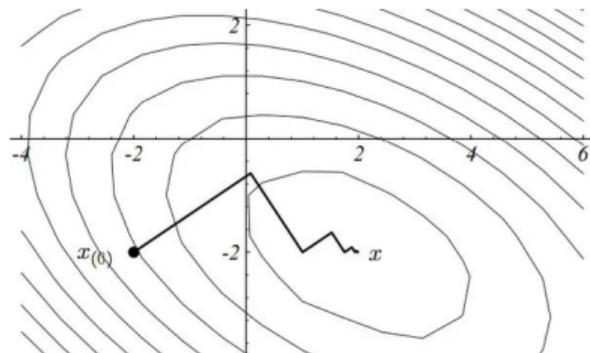
- However, the posterior distribution

$$p(\boldsymbol{\omega}|\mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega})p(\boldsymbol{\omega})}{p(\mathbf{Y}|\mathbf{X})}$$

is intractable.

Standard deep learning

- Finds a point estimate (model) that **minimizes** some loss function



Bayesian deep learning

- Evaluates predictive distributions via **marginalization**

$$p(\mathbf{y}|\mathbf{x}, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\omega})p(\boldsymbol{\omega}|\mathbf{X}, \mathbf{Y})d\boldsymbol{\omega}$$

- Assumption: $p(\mathbf{y}|\mathbf{x}, \boldsymbol{\omega}) = \mathcal{N}(\mathbf{y}; f^{\boldsymbol{\omega}}(\mathbf{x}), \tau^{-1}\mathbf{I})$, with τ is said to be the model precision. (We'll revisit this assumption.)

- However, the posterior distribution

$$p(\boldsymbol{\omega}|\mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega})p(\boldsymbol{\omega})}{p(\mathbf{Y}|\mathbf{X})}$$

is intractable.

- ▶ Approximate the posterior $p(\omega|\mathbf{X}, \mathbf{Y})$ with a **parameterized** distribution $q_\theta(\omega)$

- ▶ The goal is to solve

$$\theta^* \in \underset{\theta}{\operatorname{argmin}} \quad \operatorname{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y}))$$

- ▶ Finding an equivalent problem

$$\begin{aligned} \operatorname{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y})) &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)}{p(\omega|\mathbf{X}, \mathbf{Y})} d\omega \\ &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)p(\mathbf{Y}|\mathbf{X})}{p(\mathbf{Y}|\mathbf{X}, \omega)p(\omega)} d\omega \quad (\text{Bayes' rule}) \\ &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)}{p(\omega)} d\omega - \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega + \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}) d\omega \\ &= \underbrace{\operatorname{KL}(q_\theta(\omega) \parallel p(\omega)) - \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega}_{-\mathcal{L}_{\text{VI}}(\theta)} + \underbrace{\log p(\mathbf{Y}|\mathbf{X})}_{\text{log evidence}} \\ &= -\text{ELBO: Evidence Lower BOund} \end{aligned}$$

- ▶ Approximate the posterior $p(\omega|\mathbf{X}, \mathbf{Y})$ with a **parameterized** distribution $q_\theta(\omega)$

- ▶ The goal is to solve

$$\theta^* \in \operatorname{argmin}_{\theta} \operatorname{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y}))$$

- ▶ Finding an equivalent problem

$$\begin{aligned} \operatorname{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y})) &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)}{p(\omega|\mathbf{X}, \mathbf{Y})} d\omega \\ &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)p(\mathbf{Y}|\mathbf{X})}{p(\mathbf{Y}|\mathbf{X}, \omega)p(\omega)} d\omega \quad (\text{Bayes' rule}) \\ &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)}{p(\omega)} d\omega - \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega + \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}) d\omega \\ &= \underbrace{\operatorname{KL}(q_\theta(\omega) \parallel p(\omega)) - \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega}_{-\mathcal{L}_{\text{VI}}(\theta)} + \underbrace{\log p(\mathbf{Y}|\mathbf{X})}_{\text{log evidence}} \\ &= -\text{ELBO: Evidence Lower BOund} \end{aligned}$$

- ▶ Approximate the posterior $p(\omega|\mathbf{X}, \mathbf{Y})$ with a **parameterized** distribution $q_\theta(\omega)$

- ▶ The goal is to solve

$$\theta^* \in \operatorname{argmin}_{\theta} \operatorname{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y}))$$

- ▶ Finding an equivalent problem

$$\begin{aligned} \operatorname{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y})) &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)}{p(\omega|\mathbf{X}, \mathbf{Y})} d\omega \\ &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)p(\mathbf{Y}|\mathbf{X})}{p(\mathbf{Y}|\mathbf{X}, \omega)p(\omega)} d\omega \quad (\text{Bayes' rule}) \\ &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)}{p(\omega)} d\omega - \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega + \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}) d\omega \\ &= \underbrace{\operatorname{KL}(q_\theta(\omega) \parallel p(\omega))}_{-\mathcal{L}_{\text{VI}}(\theta)} - \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega + \underbrace{\log p(\mathbf{Y}|\mathbf{X})}_{\text{log evidence}} \\ &= -\text{ELBO: Evidence Lower BOund} \end{aligned}$$

- ▶ Approximate the posterior $p(\omega|\mathbf{X}, \mathbf{Y})$ with a **parameterized** distribution $q_\theta(\omega)$

- ▶ The goal is to solve

$$\theta^* \in \operatorname{argmin}_{\theta} \operatorname{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y}))$$

- ▶ Finding an equivalent problem

$$\begin{aligned} \operatorname{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y})) &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)}{p(\omega|\mathbf{X}, \mathbf{Y})} d\omega \\ &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)p(\mathbf{Y}|\mathbf{X})}{p(\mathbf{Y}|\mathbf{X}, \omega)p(\omega)} d\omega \quad (\text{Bayes' rule}) \\ &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)}{p(\omega)} d\omega - \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega + \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}) d\omega \\ &= \underbrace{\operatorname{KL}(q_\theta(\omega) \parallel p(\omega))}_{-\mathcal{L}_{\text{VI}}(\theta)} - \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega + \underbrace{\log p(\mathbf{Y}|\mathbf{X})}_{\text{log evidence}} \\ &= -\text{ELBO: Evidence Lower BOund} \end{aligned}$$

- ▶ Approximate the posterior $p(\omega|\mathbf{X}, \mathbf{Y})$ with a **parameterized** distribution $q_\theta(\omega)$

- ▶ The goal is to solve

$$\theta^* \in \operatorname{argmin}_{\theta} \operatorname{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y}))$$

- ▶ Finding an equivalent problem

$$\begin{aligned} \operatorname{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y})) &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)}{p(\omega|\mathbf{X}, \mathbf{Y})} d\omega \\ &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)p(\mathbf{Y}|\mathbf{X})}{p(\mathbf{Y}|\mathbf{X}, \omega)p(\omega)} d\omega \quad (\text{Bayes' rule}) \\ &= \int q_\theta(\omega) \log \frac{q_\theta(\omega)}{p(\omega)} d\omega - \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega + \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}) d\omega \\ &= \underbrace{\operatorname{KL}(q_\theta(\omega) \parallel p(\omega)) - \int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega}_{-\mathcal{L}_{\text{VI}}(\theta)} + \underbrace{\log p(\mathbf{Y}|\mathbf{X})}_{\text{log evidence}} \\ &= -\text{ELBO: Evidence Lower BOund} \end{aligned}$$

- ▶ The variational inference problem is translated into ELBO maximization

$$\begin{aligned}
 \theta^* &\in \operatorname{argmax}_{\theta} \mathcal{L}_{\text{VI}}(\theta) \\
 &= \operatorname{argmax}_{\theta} \underbrace{\int q_{\theta}(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega}_{\text{Expected log likelihood}} - \text{KL}(q_{\theta}(\omega) \parallel p(\omega)) \\
 &= \operatorname{argmax}_{\theta} \sum_{i=1}^N \int q_{\theta}(\omega) \log p(\mathbf{y}_i | f^{\omega}(\mathbf{x}_i)) d\omega - \text{KL}(q_{\theta}(\omega) \parallel p(\omega)),
 \end{aligned}$$

assuming that $(\mathbf{x}_i, \mathbf{y}_i)$ are drawn independently from the data distribution

- ▶ $f^{\omega}(\mathbf{x}_i)$ is the model's stochastic output
- ▶ However, the integral is still **intractable** for models with more than one hidden layer

- ▶ The variational inference problem is translated into ELBO maximization

$$\begin{aligned}\boldsymbol{\theta}^* &\in \operatorname{argmax}_{\boldsymbol{\theta}} \mathcal{L}_{\text{VI}}(\boldsymbol{\theta}) \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} \underbrace{\int q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \log p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega}) d\boldsymbol{\omega}}_{\text{Expected log likelihood}} - \text{KL}(q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \parallel p(\boldsymbol{\omega})) \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^N \int q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \log p(\mathbf{y}_i | f^{\boldsymbol{\omega}}(\mathbf{x}_i)) d\boldsymbol{\omega} - \text{KL}(q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \parallel p(\boldsymbol{\omega})),\end{aligned}$$

assuming that $(\mathbf{x}_i, \mathbf{y}_i)$ are drawn independently from the data distribution

- ▶ $f^{\boldsymbol{\omega}}(\mathbf{x}_i)$ is the model's stochastic output

- ▶ However, the integral is still **intractable** for models with more than one hidden layer

- ▶ The variational inference problem is translated into ELBO maximization

$$\begin{aligned}
 \boldsymbol{\theta}^* &\in \operatorname{argmax}_{\boldsymbol{\theta}} \mathcal{L}_{\text{VI}}(\boldsymbol{\theta}) \\
 &= \operatorname{argmax}_{\boldsymbol{\theta}} \underbrace{\int q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \log p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega}) d\boldsymbol{\omega}}_{\text{Expected log likelihood}} - \text{KL}(q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \parallel p(\boldsymbol{\omega})) \\
 &= \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^N \int q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \log p(\mathbf{y}_i | f^{\boldsymbol{\omega}}(\mathbf{x}_i)) d\boldsymbol{\omega} - \text{KL}(q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \parallel p(\boldsymbol{\omega})),
 \end{aligned}$$

assuming that $(\mathbf{x}_i, \mathbf{y}_i)$ are drawn independently from the data distribution

- ▶ $f^{\boldsymbol{\omega}}(\mathbf{x}_i)$ is the model's stochastic output
- ▶ However, the integral is still **intractable** for models with more than one hidden layer

- ▶ The variational inference problem is translated into ELBO maximization

$$\begin{aligned}
 \boldsymbol{\theta}^* &\in \operatorname{argmax}_{\boldsymbol{\theta}} \mathcal{L}_{\text{VI}}(\boldsymbol{\theta}) \\
 &= \operatorname{argmax}_{\boldsymbol{\theta}} \underbrace{\int q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \log p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\omega}) d\boldsymbol{\omega}}_{\text{Expected log likelihood}} - \text{KL}(q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \parallel p(\boldsymbol{\omega})) \\
 &= \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^N \int q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \log p(\mathbf{y}_i | f^{\boldsymbol{\omega}}(\mathbf{x}_i)) d\boldsymbol{\omega} - \text{KL}(q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \parallel p(\boldsymbol{\omega})),
 \end{aligned}$$

assuming that $(\mathbf{x}_i, \mathbf{y}_i)$ are drawn independently from the data distribution

- ▶ $f^{\boldsymbol{\omega}}(\mathbf{x}_i)$ is the model's stochastic output
- ▶ However, the integral is still **intractable** for models with more than one hidden layer

- ▶ **Monte-Carlo integration:** Sample T realizations of the weights from the distribution $q_\theta(\omega)$

$$\frac{1}{T} \sum_{t=1}^T \log p(\mathbf{y}_i | f^{\hat{\omega}_t}(\mathbf{x}_i)) \xrightarrow{T \rightarrow \infty} \int q_\theta(\omega) \log p(\mathbf{y}_i | f^\omega(\mathbf{x}_i)) d\omega$$

with $\hat{\omega}_t \sim q_\theta(\omega)$.¹ However, the distribution $q_\theta(\omega)$ is not known yet

- ▶ **Re-parametrize $q_\theta(\omega)$:** Let $\omega = \{\mathbf{W}_\ell\}_{\ell=1}^L$, and re-write each column in \mathbf{W}_ℓ as

$$\mathbf{W}_{\ell,k} = g(\theta_{\ell,k}, \epsilon_{\ell,k})$$

with some distribution $p(\epsilon_{\ell,k})$ that is parameter free

Ex: If $q_\theta(\omega) = \mathcal{N}(\omega; \mu, \sigma^2)$, then we can define $\omega = \mu + \sigma\epsilon$ with $\epsilon \sim \mathcal{N}(0, 1)$

¹The hat denotes a realization of a random variable

- ▶ **Monte-Carlo integration:** Sample T realizations of the weights from the distribution $q_\theta(\boldsymbol{\omega})$

$$\frac{1}{T} \sum_{t=1}^T \log p(\mathbf{y}_i | f^{\hat{\boldsymbol{\omega}}_t}(\mathbf{x}_i)) \xrightarrow{T \rightarrow \infty} \int q_\theta(\boldsymbol{\omega}) \log p(\mathbf{y}_i | f^\boldsymbol{\omega}(\mathbf{x}_i)) d\boldsymbol{\omega}$$

with $\hat{\boldsymbol{\omega}}_t \sim q_\theta(\boldsymbol{\omega})$.¹ However, the distribution $q_\theta(\boldsymbol{\omega})$ is not known yet

- ▶ **Re-parametrize $q_\theta(\boldsymbol{\omega})$:** Let $\boldsymbol{\omega} = \{\mathbf{W}_\ell\}_{\ell=1}^L$, and re-write each column in \mathbf{W}_ℓ as

$$\mathbf{W}_{\ell,k} = g(\boldsymbol{\theta}_{\ell,k}, \epsilon_{\ell,k})$$

with some distribution $p(\epsilon_{\ell,k})$ that is parameter free

Ex: If $q_\theta(\boldsymbol{\omega}) = \mathcal{N}(\boldsymbol{\omega}; \boldsymbol{\mu}, \sigma^2)$, then we can define $\boldsymbol{\omega} = \boldsymbol{\mu} + \sigma\boldsymbol{\epsilon}$ with $\boldsymbol{\epsilon} \sim \mathcal{N}(0, 1)$

¹The hat denotes a realization of a random variable

- ▶ **Monte-Carlo integration:** Sample T realizations of the weights from the distribution $q_\theta(\boldsymbol{\omega})$

$$\frac{1}{T} \sum_{t=1}^T \log p(\mathbf{y}_i | f^{\hat{\boldsymbol{\omega}}_t}(\mathbf{x}_i)) \xrightarrow{T \rightarrow \infty} \int q_\theta(\boldsymbol{\omega}) \log p(\mathbf{y}_i | f^\boldsymbol{\omega}(\mathbf{x}_i)) d\boldsymbol{\omega}$$

with $\hat{\boldsymbol{\omega}}_t \sim q_\theta(\boldsymbol{\omega})$.¹ However, the distribution $q_\theta(\boldsymbol{\omega})$ is not known yet

- ▶ **Re-parametrize $q_\theta(\boldsymbol{\omega})$:** Let $\boldsymbol{\omega} = \{\mathbf{W}_\ell\}_{\ell=1}^L$, and re-write each column in \mathbf{W}_ℓ as

$$\mathbf{W}_{\ell,k} = g(\boldsymbol{\theta}_{\ell,k}, \boldsymbol{\epsilon}_{\ell,k})$$

with some distribution $p(\boldsymbol{\epsilon}_{\ell,k})$ that is parameter free

Ex: If $q_\theta(\boldsymbol{\omega}) = \mathcal{N}(\boldsymbol{\omega}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2)$, then we can define $\boldsymbol{\omega} = \boldsymbol{\mu} + \boldsymbol{\sigma}\boldsymbol{\epsilon}$ with $\boldsymbol{\epsilon} \sim \mathcal{N}(0, 1)$

¹The hat denotes a realization of a random variable

► **Change of variables:** $\omega = g(\theta, \epsilon) \iff q_\theta(\omega|\epsilon) = \delta(\omega - g(\theta, \epsilon))$

► Re-formulating our integral

$$\begin{aligned} \int q_\theta(\omega) \log p(y_i | f^\omega(x_i)) d\omega &= \int \int q_\theta(\omega|\epsilon) p(\epsilon) \log p(y_i | f^\omega(x_i)) d\omega d\epsilon \quad (\text{total probability}) \\ &= \int p(\epsilon) \left(\int \delta(\omega - g(\theta, \epsilon)) \log p(y_i | f^\omega(x_i)) d\omega \right) d\epsilon \\ &= \int p(\epsilon) \log p(y_i | f^{g(\theta, \epsilon)}(x_i)) d\epsilon \end{aligned}$$

with $p(\epsilon) = \prod_{\ell, k} p(\epsilon_{\ell, k})$

► Next step is to use Monte-Carlo integration with a **single sample**, i.e., $T = 1$

► **Change of variables:** $\omega = g(\theta, \epsilon) \iff q_\theta(\omega|\epsilon) = \delta(\omega - g(\theta, \epsilon))$

► Re-formulating our integral

$$\begin{aligned}\int q_\theta(\omega) \log p(\mathbf{y}_i | f^\omega(\mathbf{x}_i)) d\omega &= \int \int q_\theta(\omega|\epsilon) p(\epsilon) \log p(\mathbf{y}_i | f^\omega(\mathbf{x}_i)) d\omega d\epsilon \quad (\text{total probability}) \\ &= \int p(\epsilon) \left(\int \delta(\omega - g(\theta, \epsilon)) \log p(\mathbf{y}_i | f^\omega(\mathbf{x}_i)) d\omega \right) d\epsilon \\ &= \int p(\epsilon) \log p(\mathbf{y}_i | f^{g(\theta, \epsilon)}(\mathbf{x}_i)) d\epsilon\end{aligned}$$

with $p(\epsilon) = \prod_{\ell, k} p(\epsilon_{\ell, k})$

► Next step is to use Monte-Carlo integration with a **single sample**, i.e., $T = 1$

► **Change of variables:** $\omega = g(\theta, \epsilon) \iff q_\theta(\omega|\epsilon) = \delta(\omega - g(\theta, \epsilon))$

► Re-formulating our integral

$$\begin{aligned}
 \int q_\theta(\omega) \log p(\mathbf{y}_i | f^\omega(\mathbf{x}_i)) d\omega &= \int \int q_\theta(\omega|\epsilon) p(\epsilon) \log p(\mathbf{y}_i | f^\omega(\mathbf{x}_i)) d\omega d\epsilon \quad (\text{total probability}) \\
 &= \int p(\epsilon) \left(\int \delta(\omega - g(\theta, \epsilon)) \log p(\mathbf{y}_i | f^\omega(\mathbf{x}_i)) d\omega \right) d\epsilon \\
 &= \int p(\epsilon) \log p(\mathbf{y}_i | f^{g(\theta, \epsilon)}(\mathbf{x}_i)) d\epsilon
 \end{aligned}$$

with $p(\epsilon) = \prod_{\ell, k} p(\epsilon_{\ell, k})$

► Next step is to use Monte-Carlo integration with a **single sample**, i.e., $T = 1$

► **Change of variables:** $\omega = g(\theta, \epsilon) \iff q_\theta(\omega|\epsilon) = \delta(\omega - g(\theta, \epsilon))$

► Re-formulating our integral

$$\begin{aligned}\int q_\theta(\omega) \log p(\mathbf{y}_i | f^\omega(\mathbf{x}_i)) d\omega &= \int \int q_\theta(\omega|\epsilon) p(\epsilon) \log p(\mathbf{y}_i | f^\omega(\mathbf{x}_i)) d\omega d\epsilon \quad (\text{total probability}) \\ &= \int p(\epsilon) \left(\int \delta(\omega - g(\theta, \epsilon)) \log p(\mathbf{y}_i | f^\omega(\mathbf{x}_i)) d\omega \right) d\epsilon \\ &= \int p(\epsilon) \log p(\mathbf{y}_i | f^{g(\theta, \epsilon)}(\mathbf{x}_i)) d\epsilon\end{aligned}$$

with $p(\epsilon) = \prod_{\ell, k} p(\epsilon_{\ell, k})$

► Next step is to use Monte-Carlo integration with a **single sample**, i.e., $T = 1$

► **Change of variables:** $\omega = g(\theta, \epsilon) \iff q_\theta(\omega|\epsilon) = \delta(\omega - g(\theta, \epsilon))$

► Re-formulating our integral

$$\begin{aligned}\int q_\theta(\omega) \log p(\mathbf{y}_i | f^\omega(\mathbf{x}_i)) d\omega &= \int \int q_\theta(\omega|\epsilon) p(\epsilon) \log p(\mathbf{y}_i | f^\omega(\mathbf{x}_i)) d\omega d\epsilon \quad (\text{total probability}) \\ &= \int p(\epsilon) \left(\int \delta(\omega - g(\theta, \epsilon)) \log p(\mathbf{y}_i | f^\omega(\mathbf{x}_i)) d\omega \right) d\epsilon \\ &= \int p(\epsilon) \log p(\mathbf{y}_i | f^{g(\theta, \epsilon)}(\mathbf{x}_i)) d\epsilon\end{aligned}$$

with $p(\epsilon) = \prod_{\ell, k} p(\epsilon_{\ell, k})$

► Next step is to use Monte-Carlo integration with **a single sample**, i.e., $T = 1$

- ▶ The ELBO optimization problem becomes

$$\boldsymbol{\theta}^* \in \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^N \log p(\mathbf{y}_i | f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i)) - \text{KL}(q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \parallel p(\boldsymbol{\omega}))$$

- ▶ For regression tasks: If we assume that $p(\mathbf{y}_i | f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i)) = \mathcal{N}(\mathbf{y}_i; f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i), \tau^{-1}\mathbf{I})$, then

$$\log p(\mathbf{y}_i | f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i)) = -\frac{\tau}{2} \left\| \mathbf{y}_i - f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i) \right\|^2 + \text{const}$$

- ▶ Looks familiar, Huh?

- ▶ The ELBO optimization problem becomes

$$\boldsymbol{\theta}^* \in \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{i=1}^N \log p(\mathbf{y}_i | f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i)) - \operatorname{KL}(q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \parallel p(\boldsymbol{\omega}))$$

- ▶ For regression tasks: If we assume that $p(\mathbf{y}_i | f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i)) = \mathcal{N}(\mathbf{y}_i; f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i), \tau^{-1} \mathbf{I})$, then

$$\log p(\mathbf{y}_i | f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i)) = -\frac{\tau}{2} \left\| \mathbf{y}_i - f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i) \right\|^2 + \text{const}$$

- ▶ Looks familiar, Huh?

- ▶ The ELBO optimization problem becomes

$$\boldsymbol{\theta}^* \in \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^N \log p(\mathbf{y}_i | f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i)) - \text{KL}(q_{\boldsymbol{\theta}}(\boldsymbol{\omega}) \parallel p(\boldsymbol{\omega}))$$

- ▶ For regression tasks: If we assume that $p(\mathbf{y}_i | f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i)) = \mathcal{N}(\mathbf{y}_i; f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i), \tau^{-1} \mathbf{I})$, then

$$\log p(\mathbf{y}_i | f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i)) = -\frac{\tau}{2} \left\| \mathbf{y}_i - f^{g(\boldsymbol{\theta}, \hat{\boldsymbol{\epsilon}}_i)}(\mathbf{x}_i) \right\|^2 + \text{const}$$

- ▶ Looks familiar, Huh?

- ▶ The output of a feed-forward NN (without dropout) can be written as

$$f^{\mathbf{M}_1, \dots, \mathbf{M}_L}(\mathbf{x}) = \sigma\left(\dots \sigma(\mathbf{M}_2 \underbrace{\sigma(\mathbf{M}_1 \mathbf{x})}_{\mathbf{h}_2}) \dots\right)$$

where \mathbf{M}_ℓ is a deterministic weight matrix and $\sigma(\cdot)$ is an activation function

- ▶ Dropout injects stochastic noise in the feature space $\{\mathbf{x}, \mathbf{h}_2, \dots\}$, i.e.,

$$\mathbf{h}_{\ell+1} = \sigma(\mathbf{M}_\ell(\mathbf{h}_\ell \odot \hat{\epsilon}_\ell)) = \sigma(\mathbf{M}_\ell \text{diag}(\hat{\epsilon}_\ell) \mathbf{h}_\ell) = \sigma(\widehat{\mathbf{W}}_\ell \mathbf{h}_\ell)$$

with $\epsilon_\ell \sim \text{Ber}(p_\ell)$, $\ell = 1, \dots, L$ and $\mathbf{W}_{\ell,k} = \epsilon_\ell \mathbf{M}_{\ell,k}$

$$\underset{\mathbf{M}_1, \dots, \mathbf{M}_L}{\text{argmin}} \quad \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{y}_i - f^{\widehat{\mathbf{W}}_1, \dots, \widehat{\mathbf{W}}_L}(\mathbf{x}_i) \right\|^2 + \sum_{i=1}^L \lambda_i \|\mathbf{M}_i\|^2$$

- ▶ The output of a feed-forward NN (without dropout) can be written as

$$f^{\mathbf{M}_1, \dots, \mathbf{M}_L}(\mathbf{x}) = \sigma\left(\dots \sigma(\mathbf{M}_2 \underbrace{\sigma(\mathbf{M}_1 \mathbf{x})}_{\mathbf{h}_2}) \dots\right)$$

where \mathbf{M}_ℓ is a deterministic weight matrix and $\sigma(\cdot)$ is an activation function

- ▶ Dropout injects stochastic noise in the feature space $\{\mathbf{x}, \mathbf{h}_2, \dots\}$, i.e.,

$$\mathbf{h}_{\ell+1} = \sigma(\mathbf{M}_\ell(\mathbf{h}_\ell \odot \hat{\boldsymbol{\epsilon}}_\ell)) = \sigma(\mathbf{M}_\ell \text{diag}(\hat{\boldsymbol{\epsilon}}_\ell) \mathbf{h}_\ell) = \sigma(\widehat{\mathbf{W}}_\ell \mathbf{h}_\ell)$$

with $\epsilon_\ell \sim \text{Ber}(p_\ell)$, $\ell = 1, \dots, L$ and $\mathbf{W}_{\ell,k} = \epsilon_\ell \mathbf{M}_{\ell,k}$

$$\underset{\mathbf{M}_1, \dots, \mathbf{M}_L}{\text{argmin}} \quad \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{y}_i - f^{\widehat{\mathbf{W}}_1, \dots, \widehat{\mathbf{W}}_L}(\mathbf{x}_i) \right\|^2 + \sum_{i=1}^L \lambda_i \|\mathbf{M}_i\|^2$$

- ▶ The output of a feed-forward NN (without dropout) can be written as

$$f^{\mathbf{M}_1, \dots, \mathbf{M}_L}(\mathbf{x}) = \sigma\left(\dots \sigma(\mathbf{M}_2 \underbrace{\sigma(\mathbf{M}_1 \mathbf{x})}_{\mathbf{h}_2}) \dots\right)$$

where \mathbf{M}_ℓ is a deterministic weight matrix and $\sigma(\cdot)$ is an activation function

- ▶ Dropout injects stochastic noise in the feature space $\{\mathbf{x}, \mathbf{h}_2, \dots\}$, i.e.,

$$\mathbf{h}_{\ell+1} = \sigma(\mathbf{M}_\ell(\mathbf{h}_\ell \odot \hat{\boldsymbol{\epsilon}}_\ell)) = \sigma(\mathbf{M}_\ell \text{diag}(\hat{\boldsymbol{\epsilon}}_\ell) \mathbf{h}_\ell) = \sigma(\widehat{\mathbf{W}}_\ell \mathbf{h}_\ell)$$

with $\epsilon_\ell \sim \text{Ber}(p_\ell)$, $\ell = 1, \dots, L$ and $\mathbf{W}_{\ell,k} = \epsilon_\ell \mathbf{M}_{\ell,k}$

$$\underset{\mathbf{M}_1, \dots, \mathbf{M}_L}{\text{argmin}} \quad \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{y}_i - f^{\widehat{\mathbf{W}}_1, \dots, \widehat{\mathbf{W}}_L}(\mathbf{x}_i) \right\|^2 + \sum_{i=1}^L \lambda_i \|\mathbf{M}_i\|^2$$

- ▶ The two problems are equivalent when we pick a prior distribution $p(\boldsymbol{\omega})$ and a family of distributions $q_\theta(\boldsymbol{\omega})$ that satisfy

$$\text{KL}(q_\theta(\boldsymbol{\omega}) \parallel p(\boldsymbol{\omega})) = \frac{N\tau}{2} \sum_{i=1}^L \lambda_i \|\mathbf{M}_i\|^2$$

- ▶ **When is it achievable?**

- ▶ Under a prior distribution

$$p(\boldsymbol{\omega}) = \prod_{\ell=1}^L p(\mathbf{W}_\ell) = \prod_{\ell=1}^L \mathcal{N}(\mathbf{0}, \mathbf{I}/l_\ell^2),$$

with $l_\ell^2 = \frac{2N\tau\lambda_\ell}{p_\ell}$, and

- ▶ High dimensional random vectors $\mathbf{W}_{\ell,k}, \forall \ell, k$, which means the number of neurons at each layer should be sufficiently large

Sketch of the proof:

1. Under the above assumptions, $q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k})$ is a **mixture of two Gaussian distributions**:

► In Dropout, $\mathbf{W}_{\ell,k} = \epsilon_{\ell} \mathbf{M}_{\ell,k} \iff p(\mathbf{W}_{\ell,k} | \epsilon_{\ell}) = \delta(\mathbf{W}_{\ell,k} - \epsilon_{\ell} \mathbf{M}_{\ell,k})$

$$q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k}) = \sum_{\epsilon_{\ell}=0}^1 p(\mathbf{W}_{\ell,k} | \epsilon_{\ell}) p(\epsilon_{\ell}) = p_{\ell} \mathcal{N}(\mathbf{W}_{\ell,k}; \mathbf{M}_{\ell,k}, \sigma^2 \mathbf{I}) + (1 - p_{\ell}) \mathcal{N}(\mathbf{W}_{\ell,k}; \mathbf{0}, \sigma^2 \mathbf{I})$$

2. KL divergence between $q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k})$ and $p(\mathbf{W}_{\ell,k})$ is

$$\text{KL}(q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k}) \parallel p(\mathbf{W}_{\ell,k})) = l_{\ell}^2 \frac{p_{\ell}}{2} \|\mathbf{M}_{\ell,k}\|^2 + \text{const}$$

since the mixture components do not overlap in high dimension spaces

3. Total KL divergence is

$$\text{KL}(q_{\theta}(\omega) \parallel p(\omega)) = \sum_{\ell,k} \text{KL}(q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k}) \parallel p(\mathbf{W}_{\ell,k}))$$

Sketch of the proof:

1. Under the above assumptions, $q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k})$ is a **mixture of two Gaussian distributions**:

► In Dropout, $\mathbf{W}_{\ell,k} = \epsilon_{\ell} \mathbf{M}_{\ell,k} \iff p(\mathbf{W}_{\ell,k} | \epsilon_{\ell}) = \delta(\mathbf{W}_{\ell,k} - \epsilon_{\ell} \mathbf{M}_{\ell,k})$

$$q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k}) = \sum_{\epsilon_{\ell}=0}^1 p(\mathbf{W}_{\ell,k} | \epsilon_{\ell}) p(\epsilon_{\ell}) = p_{\ell} \mathcal{N}(\mathbf{W}_{\ell,k}; \mathbf{M}_{\ell,k}, \sigma^2 \mathbf{I}) + (1 - p_{\ell}) \mathcal{N}(\mathbf{W}_{\ell,k}; \mathbf{0}, \sigma^2 \mathbf{I})$$

2. KL divergence between $q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k})$ and $p(\mathbf{W}_{\ell,k})$ is

$$\text{KL}(q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k}) \parallel p(\mathbf{W}_{\ell,k})) = l_{\ell}^2 \frac{p_{\ell}}{2} \|\mathbf{M}_{\ell,k}\|^2 + \text{const}$$

since the mixture components do not overlap in high dimension spaces

3. Total KL divergence is

$$\text{KL}(q_{\theta}(\omega) \parallel p(\omega)) = \sum_{\ell,k} \text{KL}(q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k}) \parallel p(\mathbf{W}_{\ell,k}))$$

Sketch of the proof:

1. Under the above assumptions, $q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k})$ is a **mixture of two Gaussian distributions**:

► In Dropout, $\mathbf{W}_{\ell,k} = \epsilon_{\ell} \mathbf{M}_{\ell,k} \iff p(\mathbf{W}_{\ell,k} | \epsilon_{\ell}) = \delta(\mathbf{W}_{\ell,k} - \epsilon_{\ell} \mathbf{M}_{\ell,k})$

$$q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k}) = \sum_{\epsilon_{\ell}=0}^1 p(\mathbf{W}_{\ell,k} | \epsilon_{\ell}) p(\epsilon_{\ell}) = p_{\ell} \mathcal{N}(\mathbf{W}_{\ell,k}; \mathbf{M}_{\ell,k}, \sigma^2 \mathbf{I}) + (1 - p_{\ell}) \mathcal{N}(\mathbf{W}_{\ell,k}; \mathbf{0}, \sigma^2 \mathbf{I})$$

2. KL divergence between $q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k})$ and $p(\mathbf{W}_{\ell,k})$ is

$$\text{KL}(q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k}) \parallel p(\mathbf{W}_{\ell,k})) = l_{\ell}^2 \frac{p_{\ell}}{2} \|\mathbf{M}_{\ell,k}\|^2 + \text{const}$$

since the mixture components do not overlap in high dimension spaces

3. Total KL divergence is

$$\text{KL}(q_{\theta}(\boldsymbol{\omega}) \parallel p(\boldsymbol{\omega})) = \sum_{\ell,k} \text{KL}(q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k}) \parallel p(\mathbf{W}_{\ell,k}))$$

- ▶ Gal and his co-authors showed that a standard NN trained with dropout is equivalent to a Bayesian NN.
- ▶ A standard NN learns a model that minimizes a loss function, while BNNs learn a distribution over the models that maximizes an expected loglikelihood (plus regularization terms)
- ▶ Under the above conditions, the optimal weights \mathbf{M}_ℓ in a standard NN = the optimal parameter $\boldsymbol{\theta}$ in a Bayesian NN

Uncertainty Estimates

- ▶ Recall that the predictive distribution is

$$p(\mathbf{y}_* | \mathbf{x}_*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}_* | f^\omega(\mathbf{x}_*)) p(\omega | \mathbf{X}, \mathbf{Y}) d\omega$$

replaced with $q_\theta(\mathbf{y}_* | \mathbf{x}_*) = \int p(\mathbf{y}_* | f^\omega(\mathbf{x}_*)) q_\theta(\omega) d\omega,$

with $p(\mathbf{y}_* | f^\omega(\mathbf{x}_*)) = \mathcal{N}(\mathbf{y}_*; f^\omega(\mathbf{x}_*), \tau^{-1} \mathbf{I})$

- ▶ Predictive Mean:

$$\begin{aligned} \mathbb{E}_{q_\theta(\mathbf{y}_* | \mathbf{x}_*)}[\mathbf{y}_*] &= \int \mathbf{y}_* q_\theta(\mathbf{y}_* | \mathbf{x}_*) d\mathbf{y}_* = \int \left(\int \mathbf{y}_* p(\mathbf{y}_* | f^\omega(\mathbf{x}_*)) d\mathbf{y}_* \right) q_\theta(\omega) d\omega \\ &= \int f^\omega(\mathbf{x}_*) q_\theta(\omega) d\omega \quad (\text{approximated by Monte-Carlo integration}) \end{aligned}$$

- ▶ In testing, sample T realization of the model (using dropout) to have $\tilde{\mathbf{E}}[\mathbf{y}_*] = \frac{1}{T} \sum_{t=1}^T f^{\hat{\omega}_t}(\mathbf{x}_*)$

- ▶ Recall that the predictive distribution is

$$p(\mathbf{y}_* | \mathbf{x}_*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}_* | f^\omega(\mathbf{x}_*)) p(\omega | \mathbf{X}, \mathbf{Y}) d\omega$$

replaced with $q_\theta(\mathbf{y}_* | \mathbf{x}_*) = \int p(\mathbf{y}_* | f^\omega(\mathbf{x}_*)) q_\theta(\omega) d\omega,$

with $p(\mathbf{y}_* | f^\omega(\mathbf{x}_*)) = \mathcal{N}(\mathbf{y}_*; f^\omega(\mathbf{x}_*), \tau^{-1} \mathbf{I})$

- ▶ **Predictive Mean:**

$$\begin{aligned} \mathbb{E}_{q_\theta(\mathbf{y}_* | \mathbf{x}_*)}[\mathbf{y}_*] &= \int \mathbf{y}_* q_\theta(\mathbf{y}_* | \mathbf{x}_*) d\mathbf{y}_* = \int \left(\int \mathbf{y}_* p(\mathbf{y}_* | f^\omega(\mathbf{x}_*)) d\mathbf{y}_* \right) q_\theta(\omega) d\omega \\ &= \int f^\omega(\mathbf{x}_*) q_\theta(\omega) d\omega \quad (\text{approximated by Monte-Carlo integration}) \end{aligned}$$

- ▶ In testing, sample T realization of the model (using dropout) to have $\tilde{\mathbf{E}}[\mathbf{y}_*] = \frac{1}{T} \sum_{t=1}^T f^{\hat{\omega}_t}(\mathbf{x}_*)$

- **Predictive variance** as a measure of uncertainty

$$\widetilde{\text{Var}}[\mathbf{y}_*] = \tau^{-1} \mathbf{I} + \underbrace{\frac{1}{T} \sum_{t=1}^T f^{\hat{\omega}_t}(\mathbf{x}_*) f^{\hat{\omega}_t}(\mathbf{x}_*)^\top}_{\text{second moment}} - \widetilde{\mathbf{E}}[\mathbf{y}_*] \widetilde{\mathbf{E}}[\mathbf{y}_*]^\top$$

- **Predictive log likelihood** as a measure of uncertainty

$$\begin{aligned} \log q_\theta(\mathbf{y}_* | \mathbf{x}_*) &= \log \int p(\mathbf{y}_* | f^\omega(\mathbf{x}_*)) q_\theta(\omega) d\omega \\ \implies \widetilde{\log} q_\theta(\mathbf{y}_* | \mathbf{x}_*) &= \log \left(\frac{1}{T} \sum_{t=1}^T p(\mathbf{y}_* | f^{\hat{\omega}_t}(\mathbf{x}_*)) \right) \\ &= \text{logsumexp} \left(\frac{-\tau}{2} \left\| \mathbf{y}_* - f^{\hat{\omega}_t}(\mathbf{x}_*) \right\|^2 \right) + \frac{1}{2} \log \tau + \text{const} \end{aligned}$$

- High uncertainty = low τ = high penalty from the second term
- Over-confident model (high τ) with poor mean estimation = high penalty from the first term

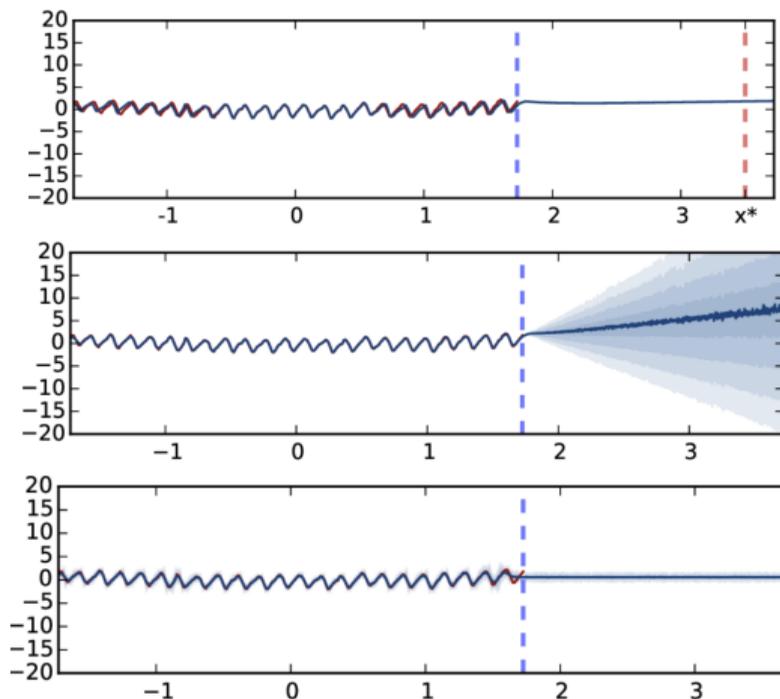
- **Predictive variance** as a measure of uncertainty

$$\widetilde{\text{Var}}[\mathbf{y}_*] = \tau^{-1} \mathbf{I} + \underbrace{\frac{1}{T} \sum_{t=1}^T f^{\hat{\omega}_t}(\mathbf{x}_*) f^{\hat{\omega}_t}(\mathbf{x}_*)^\top}_{\text{second moment}} - \widetilde{\mathbf{E}}[\mathbf{y}_*] \widetilde{\mathbf{E}}[\mathbf{y}_*]^\top$$

- **Predictive log likelihood** as a measure of uncertainty

$$\begin{aligned} \log q_\theta(\mathbf{y}_* | \mathbf{x}_*) &= \log \int p(\mathbf{y}_* | f^\omega(\mathbf{x}_*)) q_\theta(\omega) d\omega \\ \implies \widetilde{\log} q_\theta(\mathbf{y}_* | \mathbf{x}_*) &= \log \left(\frac{1}{T} \sum_{t=1}^T p(\mathbf{y}_* | f^{\hat{\omega}_t}(\mathbf{x}_*)) \right) \\ &= \text{logsumexp} \left(\frac{-\tau}{2} \left\| \mathbf{y}_* - f^{\hat{\omega}_t}(\mathbf{x}_*) \right\|^2 \right) + \frac{1}{2} \log \tau + \text{const} \end{aligned}$$

- High uncertainty = low τ = high penalty from the second term
- Over-confident model (high τ) with poor mean estimation = high penalty from the first term



CO2 concentrations dataset. (Top) Standard Dropout, (Middle) MC dropout with Relu, and (Bottom) with Tanh. Different shades of blue represent half a standard deviation.

- ▶ Add a softmax layer, $\text{Softmax}(f^\omega(\mathbf{x}))$, to predict the likelihood of each class
- ▶ Sample T realizations of the model and let the prediction be

$$c^* = \underset{c=1,\dots,C}{\operatorname{argmax}} \sum_{t=1}^T \mathbf{1}[\hat{y}_t = c]$$

- ▶ Predictive Entropy as an uncertainty estimate

$$\mathbb{H}[y_* | \mathbf{x}_*] = - \sum_c q_\theta(y_* = c | \mathbf{x}_*) \log(q_\theta(y_* = c | \mathbf{x}_*))$$

where we approximate the predictive distribution with

$$q_\theta(y_* = c | \mathbf{x}_*) = \int \underbrace{p(y_* = c | \mathbf{x}_*, \omega)}_{\text{Softmax outputs}} q_\theta(\omega) d\omega \approx \frac{1}{T} \sum_{t=1}^T p(y_* = c | \mathbf{x}_*, \hat{\omega}_t)$$

- ▶ $\exists c, q_\theta(y_* = c | \mathbf{x}_*) = 1 \implies \mathbb{H}[y_* | \mathbf{x}_*] \downarrow$, and $q_\theta(y_* = c | \mathbf{x}_*) \sim \text{Unif} \implies \mathbb{H}[y_* | \mathbf{x}_*] \uparrow$

- ▶ Add a softmax layer, $\text{Softmax}(f^\omega(\mathbf{x}))$, to predict the likelihood of each class
- ▶ Sample T realizations of the model and let the prediction be

$$c^* = \underset{c=1, \dots, C}{\operatorname{argmax}} \sum_{t=1}^T \mathbf{1}[\hat{y}_t = c]$$

- ▶ **Predictive Entropy** as an uncertainty estimate

$$\mathbb{H}[y_* | \mathbf{x}_*] = - \sum_c q_\theta(y_* = c | \mathbf{x}_*) \log(q_\theta(y_* = c | \mathbf{x}_*))$$

where we approximate the predictive distribution with

$$q_\theta(y_* = c | \mathbf{x}_*) = \int \underbrace{p(y_* = c | \mathbf{x}_*, \boldsymbol{\omega})}_{\text{Softmax outputs}} q_\theta(\boldsymbol{\omega}) d\boldsymbol{\omega} \approx \frac{1}{T} \sum_{t=1}^T p(y_* = c | \mathbf{x}_*, \hat{\boldsymbol{\omega}}_t)$$

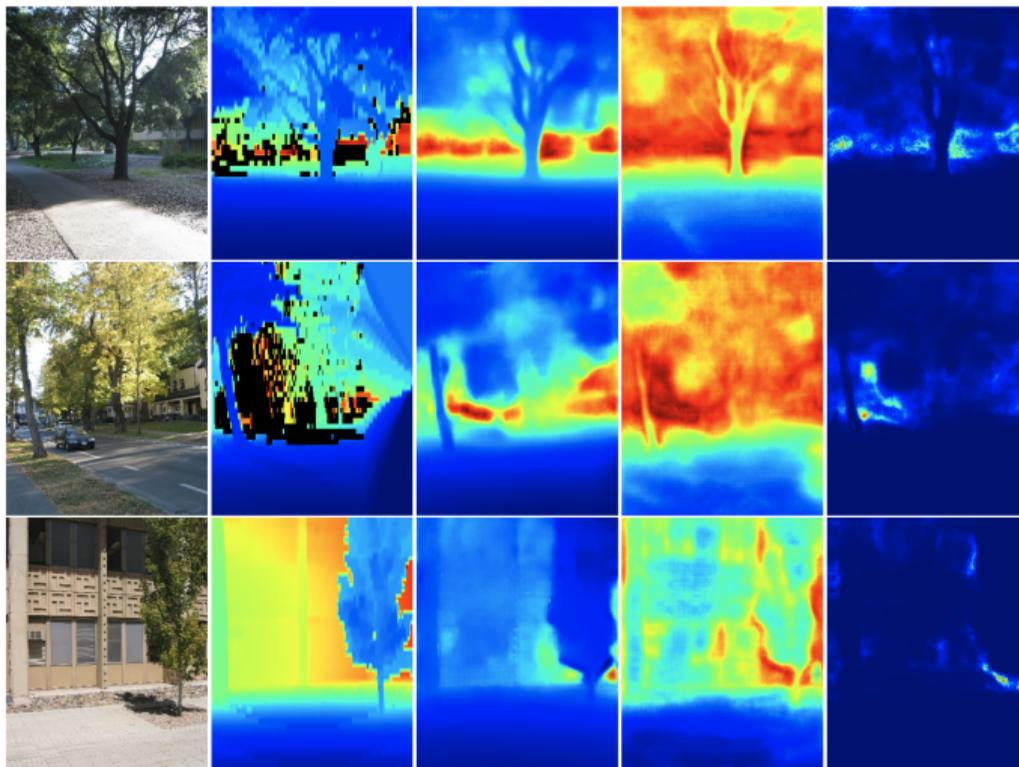
- ▶ $\exists c, q_\theta(y_* = c | \mathbf{x}_*) = 1 \implies \mathbb{H}[y_* | \mathbf{x}_*] \downarrow$, and $q_\theta(y_* = c | \mathbf{x}_*) \sim \text{Unif} \implies \mathbb{H}[y_* | \mathbf{x}_*] \uparrow$

- ▶ In order to let BNNs learn aleatoric uncertainty, we parameterize τ as $\tau^\omega(\mathbf{x})$

$$\mathcal{N}(\mathbf{y}; f^\omega(\mathbf{x}), \tau^{-1}\mathbf{I}) \implies \mathcal{N}(\mathbf{y}; f^\omega(\mathbf{x}), \tau^\omega(\mathbf{x})^{-1})$$

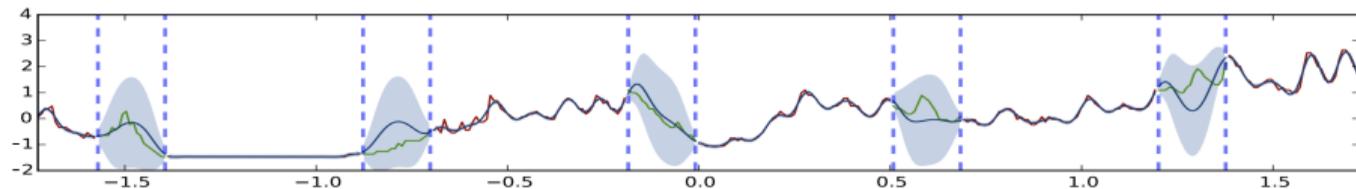
- ▶ The goal is to learn distributions over the weights used for both $f^\omega(\mathbf{x})$ and $\tau^\omega(\mathbf{x})$ following the same framework
- ▶ Predictive variance is calculated as

$$\widetilde{\text{Var}}[\mathbf{y}_*] = \frac{1}{T} \sum_{t=1}^T \tau^{\hat{\omega}_t}(\mathbf{x})\mathbf{I} + \frac{1}{T} \sum_{t=1}^T f^{\hat{\omega}_t}(\mathbf{x}_*)f^{\hat{\omega}_t}(\mathbf{x}_*)^\top - \widetilde{\mathbf{E}}[\mathbf{y}_*]\widetilde{\mathbf{E}}[\mathbf{y}_*]^\top$$

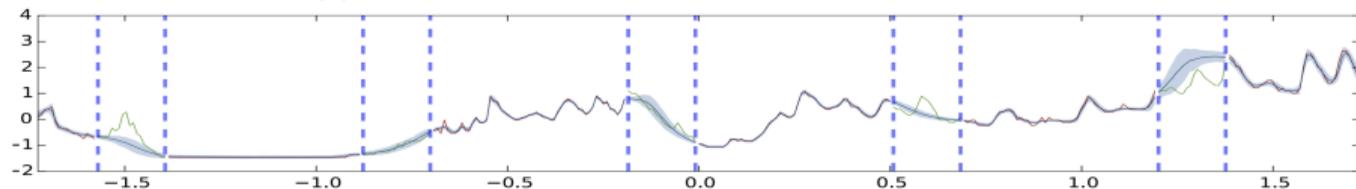


Left to Right: input image, ground truth, **depth prediction**, aleatoric uncertainty, epistemic uncertainty.
Make3D does not provide labels for depth greater than 70m

Final Remarks



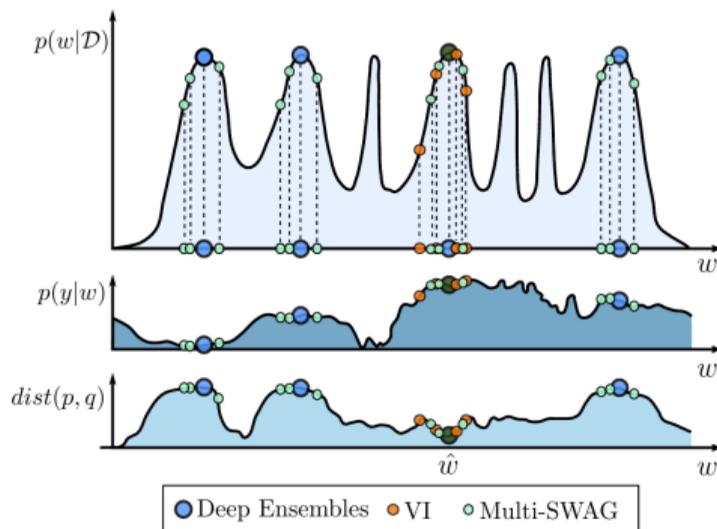
(a) Gaussian process with SE covariance function



(b) MC dropout with ReLU non-linearities

Predictive mean and uncertainties on the reconstructed solar irradiance dataset with missing segments, for the GP and MC dropout approximation. In red is the observed function and in green are the missing segments. In blue is the predictive mean plus/minus two standard deviations of the various approximations.

- ▶ Every time we train a standard NN, we reach a local minimum \implies **multimodal posterior**



- ▶ MC dropout:

$$\begin{aligned}
 q_{\theta_{\ell,k}}(\mathbf{W}_{\ell,k}) &= p_{\ell} \mathcal{N}(\mathbf{W}_{\ell,k}; \mathbf{M}_{\ell,k}, \sigma^2 \mathbf{I}) \\
 &+ (1 - p_{\ell}) \underbrace{\mathcal{N}(\mathbf{W}_{\ell,k}; \mathbf{0}, \sigma^2 \mathbf{I})}_{?}
 \end{aligned}$$

Large gap in computing the predictive distribution

$$p(\mathbf{y}_* | \mathbf{x}_*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}_* | f^{\omega}(\mathbf{x}_*)) p(\omega | \mathbf{X}, \mathbf{Y}) d\omega$$

- ▶ Deep Ensembles

Check Rahul's presentation for more details!