

[Paper review 20]

Uncertainty in Deep Learning - Chapter 3

(Yarin Gal, 2016)

[Contents]

3. Bayesian Deep Learning
 1. Advanced techniques in VI
 1. MC estimator
 2. Variance analysis of MC estimator
 2. Practical Inference in BNN
 1. SRT
 2. SRT as approximate inference
 3. KL condition
 3. Model uncertainty in BNN
 1. Uncertainty in classification
 2. Difficulties with the approach
 4. Approximate inference in complex models
 1. Bayesian CNN
 2. Bayesian RNN

3. Bayesian Deep Learning

Based on two works

- 1) MC estimation (Graves, 2011)
- 2) VI (Hinton and Van Camp, 1993)

in a Bayesian perspective!

"BNN inference + SRTs (offer a practical inference technique)"

Steps

- step 1) analyze the variance of several stochastic estimators (used in VI)
- step 2) tie these derivations to SRTs
- step 3) propose practical techniques to obtain model uncertainty

3.1 Advanced techniques in VI

[review of VI]

$$\mathcal{L}_{VI}(\theta) := - \sum_{i=1}^N \int q_{\theta}(\omega) \log p(\mathbf{y}_i | \mathbf{f}^{\omega}(\mathbf{x}_i)) d\omega + \text{KL}(q_{\theta}(\omega) || p(\omega))$$

$$\text{expected log likelihood} : - \sum_{i=1}^N \int q_{\theta}(\omega) \log p(\mathbf{y}_i | \mathbf{f}^{\omega}(\mathbf{x}_i)) d\omega$$

problems in expected log likelihood :

- problem 1) $\sum_{i=1}^N$: perform computations over the entire dataset

- problem 2) $\int q_{\theta}(\boldsymbol{\omega}) \log p(\mathbf{y}_i | \mathbf{f}^{\boldsymbol{\omega}}(\mathbf{x}_i)) d\boldsymbol{\omega}$: not tractable

Solutions

- solution 1) data sub-sampling (mini-batch optimization)
(unbiased + stochastic estimator)

$$\widehat{\mathcal{L}}_{VI}(\theta) := -\frac{N}{M} \sum_{i \in S} \int q_{\theta}(\boldsymbol{\omega}) \log p(\mathbf{y}_i | \mathbf{f}^{\boldsymbol{\omega}}(\mathbf{x}_i)) d\boldsymbol{\omega} + \text{KL}(q_{\theta}(\boldsymbol{\omega}) || p(\boldsymbol{\omega}))$$

- solution 2) MC integration

3.1.1 MC estimators

use MC estimation to estimate "EXPECTED LOG LIKELIHOOD"

(more importantly, the "derivatives" of expected log likelihood)

Estimate "integral derivatives"

$$I(\theta) = \frac{\partial}{\partial \theta} \int f(x) p_{\theta}(x) dx$$

THREE main techniques for MC estimation for $\theta = \{\mu, \sigma\}$

- find "mean" derivative estimator
- find "standard deviation" derivative estimator

(1) Score function estimator

(= Likelihood ratio estimator, Reinforce)

(2) Path-wise derivative estimator

(= reparametrization trick)

(3) Characteristic function estimator

(1) Score function estimator (= $\hat{I}_1(\theta)$)

$$\hat{I}_1(\theta) = f(x) \frac{\partial \log p_{\theta}(x)}{\partial \theta} \text{ with } x \sim p_{\theta}(x)$$

- simple
- but high variance

$$\begin{aligned} \frac{\partial}{\partial \theta} \int f(x) p_{\theta}(x) dx &= \int f(x) \frac{\partial}{\partial \theta} p_{\theta}(x) dx \\ &= \int f(x) \frac{\partial \log p_{\theta}(x)}{\partial \theta} p_{\theta}(x) dx \end{aligned}$$

(2) Path-wise derivative estimator (= $\hat{I}_2(\theta)$)

$$\hat{I}_2(\theta) = f'(g(\theta, \epsilon)) \frac{\partial}{\partial \theta} g(\theta, \epsilon)$$

- $\hat{I}_2(\mu) = f'(x)$
- $\hat{I}_2(\sigma) = f'(x) \frac{(x-\mu)}{\sigma}$

reparameterization trick

- before) $p_\theta(x) = \mathcal{N}(x; \mu, \sigma^2)$
- after) $g(\theta, \epsilon) = \mu + \sigma\epsilon$ & $p(\epsilon) = \mathcal{N}(\epsilon; 0, I)$

"mean" derivative estimator : $\frac{\partial}{\partial \mu} \int f(x)p_\theta(x)dx = \int f'(x)p_\theta(x)dx$

"standard deviation" derivative estimator : $\frac{\partial}{\partial \sigma} \int f(x)p_\theta(x)dx = \int f'(x)\frac{(x-\mu)}{\sigma}p_\theta(x)dx$

(3) Characteristic function estimator (= $\hat{I}_3(\theta)$)

$$\hat{I}_2(\mu) = f'(x)$$

$$\hat{I}_3(\sigma) = \sigma f''(x) \quad \left(\frac{\partial}{\partial \sigma} \int f(x)p_\theta(x)dx = 2\sigma \cdot \frac{1}{2} \int f''(x)p_\theta(x)dx \right)$$

- rely on the characteristic function of "Gaussian distribution"
(restricts the estimator to Gaussian $p_\theta(x)$ alone)

[tip] Reparameterization Trick

- use $p_\theta(x) = \int p_\theta(x, \epsilon)d\epsilon = \int p_\theta(x | \epsilon)p(\epsilon)d\epsilon$
- where $p_\theta(x | \epsilon) = \delta(x - g(\theta, \epsilon))$
(+ $\delta(x - g(\theta, \epsilon))$ is zero for all x apart from $x = g(\theta, \epsilon)$)

$$\begin{aligned} \frac{\partial}{\partial \theta} \int f(x)p_\theta(x)dx &= \frac{\partial}{\partial \theta} \int f(x) \left(\int p_\theta(x, \epsilon)d\epsilon \right) dx \\ &= \frac{\partial}{\partial \theta} \int f(x)p_\theta(x | \epsilon)p(\epsilon)d\epsilon dx \\ &= \frac{\partial}{\partial \theta} \int \left(\int f(x)\delta(x - g(\theta, \epsilon))dx \right) p(\epsilon)d\epsilon \\ &= \frac{\partial}{\partial \theta} \int f(g(\theta, \epsilon))p(\epsilon)d\epsilon \\ &= \int \frac{\partial}{\partial \theta} f(g(\theta, \epsilon))p(\epsilon)d\epsilon \\ &= \int f'(g(\theta, \epsilon)) \frac{\partial}{\partial \theta} g(\theta, \epsilon)p(\epsilon)d\epsilon \end{aligned}$$

3.1.2 Variance Analysis of MC estimator

None of 3 estimators has the lowest variance for all functions of $f(x)$

- (1) score function
- (2) path-wise derivative function
- (3) characteristic function

Properties that (2), (3) have lower variance than (1) : in the paper

From empirical observations, (2) seems to be good!

Will continue our work using the path-wise derivative estimator

3.2 Practical Inference in BNN

in terms of "PRACTICALITY"

Graves (2011)

- (a) delta approximating distribution (use "characteristic function")
- (b) fully factorized approximating distribution
(factorized the approximating distribution for EACH WEIGHT scalar, thus "losing weight correlation" → hurt performance)

Advancement

- (a) use "path-wise derivative estimator" instead (used 're-parameterization trick ')
- (b) factorize the approximating distribution for EACH ROW WEIGHT

ELBO using (1) reparam trick & (2) MC estimation

$$\begin{aligned}
 \widehat{\mathcal{L}}_{\text{VI}}(\theta) &= -\frac{N}{M} \sum_{i \in S} \int q_{\theta}(\omega) \log p(\mathbf{y}_i | \mathbf{f}^{\omega}(\mathbf{x}_i)) d\omega + \text{KL}(q_{\theta}(\omega) || p(\omega)) \\
 &= -\frac{N}{M} \sum_{i \in S} \int p(\epsilon) \log p(\mathbf{y}_i | \mathbf{f}^{g(\theta, \epsilon)}(\mathbf{x}_i)) d\epsilon + \text{KL}(q_{\theta}(\omega) || p(\omega)) \\
 &\approx -\frac{N}{M} \sum_{i \in S} \log p(\mathbf{y}_i | \mathbf{f}^{g(\theta, \epsilon)}(\mathbf{x}_i)) + \text{KL}(q_{\theta}(\omega) || p(\omega)) \\
 &\equiv \widehat{\mathcal{L}}_{\text{MC}}(\theta) \quad \text{where } \mathbb{E}_{S, \epsilon}(\widehat{\mathcal{L}}_{\text{MC}}(\theta)) = \mathcal{L}_{\text{VI}}(\theta)
 \end{aligned}$$

Predictive distribution

$$\begin{aligned}
 \tilde{q}_{\theta}(\mathbf{y}^* | \mathbf{x}^*) &:= \frac{1}{T} \sum_{t=1}^T p(\mathbf{y}^* | \mathbf{x}^*, \widehat{\omega}_t) \xrightarrow{T \rightarrow \infty} \int p(\mathbf{y}^* | \mathbf{x}^*, \omega) q_{\theta}(\omega) d\omega \\
 &\approx \int p(\mathbf{y}^* | \mathbf{x}^*, \omega) p(\omega | \mathbf{X}, \mathbf{Y}) d\omega \\
 &= p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{X}, \mathbf{Y})
 \end{aligned}$$

where $\widehat{\omega}_t \sim q_{\theta}(\omega)$

[Summary]

optimizing $\widehat{\mathcal{L}}_{\text{MC}}(\theta)$ w.r.t θ = optimizing $\widehat{\mathcal{L}}_{\text{VI}}(\theta)$ w.r.t θ

Algorithm 1 Minimise divergence between $q_{\theta}(\omega)$ and $p(\omega | X, Y)$

- 1: Given dataset \mathbf{X}, \mathbf{Y} ,
- 2: Define learning rate schedule η ,
- 3: Initialise parameters θ randomly.
- 4: **repeat**
- 5: Sample M random variables $\widehat{\epsilon}_i \sim p(\epsilon)$, S a random subset of $\{1, \dots, N\}$ of size M .
- 6: Calculate stochastic derivative estimator w.r.t. θ :

$$\widehat{\Delta\theta} \leftarrow -\frac{N}{M} \sum_{i \in S} \frac{\partial}{\partial \theta} \log p(\mathbf{y}_i | \mathbf{f}^{g(\theta, \widehat{\epsilon}_i)}(\mathbf{x}_i)) + \frac{\partial}{\partial \theta} \text{KL}(q_{\theta}(\omega) || p(\omega)).$$

- 7: Update θ :
 $\theta \leftarrow \theta + \eta \widehat{\Delta\theta}$.
 - 8: **until** θ has converged.
-

3.2.1 SRT (Stochastic Regularization Techniques)

"REGULARIZE models through injection of STOCHASTIC NOISE"

ex) dropout, multiplicative Gaussian Noise, drop connect ...

notation

- M : deterministic matrix
- W : random variable defined over the set of real matrices
- \hat{W} : realization of W

(1) Dropout

- two binary vectors $\hat{\epsilon}_1, \hat{\epsilon}_2$
(each with dimension Q (=input dim) and K (=intermediate dim))
- parameters : $\theta = \{M_1, M_2, b\}$
- $\hat{y} = \hat{h}M_2$
 - $\hat{h} = h \odot \hat{\epsilon}_2$
 - $h = \sigma(\hat{x}M_1 + b)$
 - $\hat{x} = x \odot \hat{\epsilon}_1^3$
- sample $\hat{\epsilon}_i$ for every input & every forward pass
use the same value for backward pass
- test time : do not sample any variables & just use the original units x, h scaled by $\frac{1}{1-p_i}$

(2) Multiplicative Gaussian Noise

- same as (1), except $\hat{\epsilon}_i \sim N(1, \alpha)$

3.2.2 SRT as approximate inference

inject noise to feature space (=the input features to each layer, which are x and h)

$$\begin{aligned} \hat{y} &= \hat{h}M_2 \\ &= (h \odot \hat{\epsilon}_2) M_2 \\ &= (h \cdot \text{diag}(\hat{\epsilon}_2)) M_2 \\ &= h (\text{diag}(\hat{\epsilon}_2) M_2) \\ &= \sigma(\hat{x}M_1 + b) (\text{diag}(\hat{\epsilon}_2) M_2) \\ &= \sigma((x \odot \hat{\epsilon}_1) M_1 + b) (\text{diag}(\hat{\epsilon}_2) M_2) \\ &= \sigma(x (\text{diag}(\hat{\epsilon}_1) M_1) + b) (\text{diag}(\hat{\epsilon}_2) M_2) \\ &= \sigma(x \hat{W}_1 + b) \hat{W}_2 =: f^{\hat{W}_1}, \hat{W}_2, b(x) \end{aligned}$$

(let $\hat{W}_1 := \text{diag}(\hat{\epsilon}_1)M_1$ and $\hat{W}_2 := \text{diag}(\hat{\epsilon}_2)M_2$)

(random variable realization as weights : $\hat{\omega} = \{ \hat{W}_1, \hat{W}_2, b \}$)

Loss function

$$\hat{\mathcal{L}}_{\text{dropout}}(M_1, M_2, b) := \frac{1}{M} \sum_{i \in S} E^{\hat{W}_1^i, \hat{W}_2^i, b}(\mathbf{x}_i, \mathbf{y}_i) + \lambda_1 \|M_1\|^2 + \lambda_2 \|M_2\|^2 + \lambda_3 \|b\|^2$$

where \hat{W}_1^i, \hat{W}_2^i corresponding to new masks $\hat{\epsilon}_1^i, \hat{\epsilon}_2^i$

Example) Negative Log Likelihood

$$E^{\mathbf{M}_1, \mathbf{M}_2, \mathbf{b}}(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{y} - \mathbf{f}^{\mathbf{M}_1, \mathbf{M}_2, \mathbf{b}}(\mathbf{x})\|^2 = -\frac{1}{\tau} \log p(\mathbf{y} | \mathbf{f}^{\mathbf{M}_1, \mathbf{M}_2, \mathbf{b}}(\mathbf{x})) + \text{const}$$

- where $p(\mathbf{y} | \mathbf{f}^{\mathbf{M}_1, \mathbf{M}_2, \mathbf{b}}(\mathbf{x})) = \mathcal{N}(\mathbf{y}; \mathbf{f}^{\mathbf{M}_1, \mathbf{M}_2, \mathbf{b}}(\mathbf{x}), \tau^{-1} I)$ with τ^{-1} observation noise
- (for classification) $\tau = 1$

Reparametrization trick

$$\widehat{\omega}_i = \{\widehat{\mathbf{W}}_1^i, \widehat{\mathbf{W}}_2^i, \mathbf{b}\} = \{\text{diag}(\widehat{\epsilon}_1^i) \mathbf{M}_1, \text{diag}(\widehat{\epsilon}_2^i) \mathbf{M}_2, \mathbf{b}\} =: g(\theta, \widehat{\epsilon}_i)$$

where $\theta = \{\mathbf{M}_1, \mathbf{M}_2, \mathbf{b}\}$, $\widehat{\epsilon}_1^i \sim p(\epsilon_1)$, and $\widehat{\epsilon}_2^i \sim p(\epsilon_2)$ for $1 \leq i \leq N$

Loss function :

$$\widehat{\mathcal{L}}_{\text{dropout}}(\mathbf{M}_1, \mathbf{M}_2, \mathbf{b}) = -\frac{1}{M\tau} \sum_{i \in S} \log p(\mathbf{y}_i | \mathbf{f}^{g(\theta, \widehat{\epsilon}_i)}(\mathbf{x})) + \lambda_1 \|\mathbf{M}_1\|^2 + \lambda_2 \|\mathbf{M}_2\|^2 + \lambda_3 \|\mathbf{b}\|^2$$

Derivative of the loss function :

$$\frac{\partial}{\partial \theta} \widehat{\mathcal{L}}_{\text{dropout}}(\theta) = -\frac{1}{M\tau} \sum_{i \in S} \frac{\partial}{\partial \theta} \log p(\mathbf{y}_i | \mathbf{f}^{g(\theta, \widehat{\epsilon}_i)}(\mathbf{x})) + \frac{\partial}{\partial \theta} (\lambda_1 \|\mathbf{M}_1\|^2 + \lambda_2 \|\mathbf{M}_2\|^2 + \lambda_3 \|\mathbf{b}\|^2)$$

[Summary]

optimizing $\widehat{\mathcal{L}}_{\text{dropout}}(\mathbf{M}_1, \mathbf{M}_2, \mathbf{b})$ with dropout :

Algorithm 2 Optimisation of a neural network with dropout

- 1: Given dataset \mathbf{X}, \mathbf{Y} ,
- 2: Define learning rate schedule η ,
- 3: Initialise parameters θ randomly.
- 4: **repeat**
- 5: Sample M random variables $\widehat{\epsilon}_i \sim p(\epsilon)$, S a random subset of $\{1, \dots, N\}$ of size M .
- 6: Calculate derivative w.r.t. θ :

$$\widehat{\Delta \theta} \leftarrow -\frac{1}{M\tau} \sum_{i \in S} \frac{\partial}{\partial \theta} \log p(\mathbf{y}_i | \mathbf{f}^{g(\theta, \widehat{\epsilon}_i)}(\mathbf{x})) + \frac{\partial}{\partial \theta} (\lambda_1 \|\mathbf{M}_1\|^2 + \lambda_2 \|\mathbf{M}_2\|^2 + \lambda_3 \|\mathbf{b}\|^2).$$

- 7: Update θ :
 $\theta \leftarrow \theta + \eta \widehat{\Delta \theta}.$
 - 8: **until** θ has converged.
-

instead of Dropout...

1) Multiplicative Gaussian Noise

- $g(\theta, \epsilon) = \{\text{diag}(\epsilon_1) \mathbf{M}_1, \text{diag}(\epsilon_2) \mathbf{M}_2, \mathbf{b}\}$
with $p(\epsilon_l)$ (for $l = 1, 2$) a product of $\mathcal{N}(1, \alpha)$ with positive α

2) Drop connect

- $g(\theta, \epsilon) = \{\mathbf{M}_1 \odot \epsilon_1, \mathbf{M}_2 \odot \epsilon_2, \mathbf{b}\}$
with $p(\epsilon_l)$ a product of Bernoulli random variables

3) Additive Gaussian Noise

- $g(\theta, \epsilon) = \{\mathbf{M}_1 + \epsilon_1, \mathbf{M}_2 + \epsilon_2, \mathbf{b}\}$
with $p(\epsilon_i)$ a product of $\mathcal{N}(0, \alpha)$ for each weight scalar

[Algorithm summary]

[3.2.1] Algorithm 1) Minimize divergence between $q_\theta(w) = p(w | X, Y)$

[3.2.2] Algorithm 2) Optimization of a NN with Dropout

For algorithm 1 = 2....

- 1) "regularization term derivatives" should be same (= KL condition)

$$\frac{\partial}{\partial \theta} \text{KL}(q_\theta(\omega) \| p(\omega)) = \frac{\partial}{\partial \theta} N\tau \left(\lambda_1 \|\mathbf{M}_1\|^2 + \lambda_2 \|\mathbf{M}_2\|^2 + \lambda_3 \|\mathbf{b}\|^2 \right)$$

- 2) scale of derivatives

$$\frac{\partial}{\partial \theta} \widehat{\mathcal{L}}_{\text{dropout}}(\theta) = \frac{1}{N\tau} \frac{\partial}{\partial \theta} \widehat{\mathcal{L}}_{\text{MC}}(\theta)$$

Summary :

- "Optimizing any NN with DROPOUT" = "APPROXIMATE INFERENCE in a probabilistic interpretation of the model"
- NN trained with dropout is "Bayesian NN"

3.2.3 KL condition

condition for "VI" = "DO" → depends on the model specification (choice of $p(w)$ and $q_\theta(w)$)

Example 1)

- prior : $p(\omega) = \prod_{i=1}^L p(\mathbf{W}_i) = \prod_{i=1}^L \mathcal{N}(0, \mathbf{I}/l_i^2)$, where $l_i^2 = \frac{2N\tau\lambda_i}{1-p_i}$ (prior length scale)
- then

$$\frac{\partial}{\partial \theta} \text{KL}(q_\theta(\omega) \| p(\omega)) \approx \frac{\partial}{\partial \theta} N\tau \left(\lambda_1 \|\mathbf{M}_1\|^2 + \lambda_2 \|\mathbf{M}_2\|^2 + \lambda_3 \|\mathbf{b}\|^2 \right)$$

Example 2) discrete prior

- $p(\mathbf{w}) \propto e^{-\frac{l^2}{2} \mathbf{w}^T \mathbf{w}}$

Example 3) improper log-uniform prior

- for Multiplicative Gaussian Noise

3.3 Model Uncertainty in BNN

approximate predictive distribution

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

- $\omega = \{\mathbf{W}_i\}_{i=1}^L$ is our set of random variables for a model with L layers
- \mathbf{f}^ω : model's stochastic output
- $q_\theta^*(\omega)$: optimum

check if FIRST & SECOND MOMENT matches!

First Moment

$$\tilde{\mathbb{E}}[\mathbf{y}^*] := \frac{1}{T} \sum_{t=1}^T \mathbf{f}^{\hat{\omega}_t}(\mathbf{x}^*) \xrightarrow{T \rightarrow \infty} \mathbb{E}_{q_{\theta^*}(\mathbf{y}^*|\mathbf{x}^*)}[\mathbf{y}^*]$$

with $\hat{\omega}_t \sim q_{\theta^*}(\omega)$

- unbiased estimator, following MC integration with T samples
- when use Dropout \rightarrow "MC Dropout" (= model averaging)

(proof)

$$\begin{aligned} \mathbb{E}_{q_{\theta^*}(\mathbf{y}^*|\mathbf{x}^*)}[\mathbf{y}^*] &= \int \mathbf{y}^* q_{\theta^*}(\mathbf{y}^* | \mathbf{x}^*) d\mathbf{y}^* \\ &= \iint \mathbf{y}^* \mathcal{N}(\mathbf{y}^*; \mathbf{f}^{\omega}(\mathbf{x}^*), \tau^{-1}\mathbf{I}) q_{\theta^*}(\omega) d\omega d\mathbf{y}^* \\ &= \int \left(\int \mathbf{y}^* \mathcal{N}(\mathbf{y}^*; \mathbf{f}^{\omega}(\mathbf{x}^*), \tau^{-1}\mathbf{I}) d\mathbf{y}^* \right) q_{\theta^*}(\omega) d\omega \\ &= \int \mathbf{f}^{\omega}(\mathbf{x}^*) q_{\theta^*}(\omega) d\omega \end{aligned}$$

Second Moment

$$\tilde{\mathbb{E}}[(\mathbf{y}^*)^T (\mathbf{y}^*)] := \tau^{-1}\mathbf{I} + \frac{1}{T} \sum_{t=1}^T \mathbf{f}^{\hat{\omega}_t}(\mathbf{x}^*)^T \mathbf{f}^{\hat{\omega}_t}(\mathbf{x}^*) \xrightarrow{T \rightarrow \infty} \mathbb{E}_{p(\mathbf{y}^*|\mathbf{x}^*)}[(\mathbf{y}^*)^T (\mathbf{y}^*)]$$

with $\hat{\omega}_t \sim q_{\theta^*}(\omega)$ and $\mathbf{y}^*, \mathbf{f}^{\hat{\omega}_t}(\mathbf{x}^*)$ row vectors

- unbiased estimator, following MC integration with T samples

(proof)

$$\begin{aligned} \mathbb{E}_{q_{\theta^*}(\mathbf{y}^*|\mathbf{x}^*)}[(\mathbf{y}^*)^T (\mathbf{y}^*)] &= \int \left(\int (\mathbf{y}^*)^T (\mathbf{y}^*) p(\mathbf{y}^* | \mathbf{x}^*, \omega) d\mathbf{y}^* \right) q_{\theta^*}(\omega) d\omega \\ &= \int \left(\text{Cov}_{p(\mathbf{y}^*|\mathbf{x}^*, \omega)}[\mathbf{y}^*] + \mathbb{E}_{p(\mathbf{y}^*|\mathbf{x}^*, \omega)}[\mathbf{y}^*]^T \mathbb{E}_{p(\mathbf{y}^*|\mathbf{x}^*, \omega)}[\mathbf{y}^*] \right) q_{\theta^*}(\omega) d\omega \\ &= \int \left(\tau^{-1}\mathbf{I} + \mathbf{f}^{\omega}(\mathbf{x}^*)^T \mathbf{f}^{\omega}(\mathbf{x}^*) \right) q_{\theta^*}(\omega) d\omega \end{aligned}$$

Variance

$$\widehat{\text{Var}}[\mathbf{y}^*] := \tau^{-1}\mathbf{I} + \frac{1}{T} \sum_{t=1}^T \mathbf{f}^{\hat{\omega}_t}(\mathbf{x}^*)^T \mathbf{f}^{\hat{\omega}_t}(\mathbf{x}^*) - \tilde{\mathbb{E}}[\mathbf{y}^*]^T \tilde{\mathbb{E}}[\mathbf{y}^*] \xrightarrow{T \rightarrow \infty} \text{Var}_{q_{\theta^*}(\mathbf{y}^*|\mathbf{x}^*)}[\mathbf{y}^*]$$

How to find model precision τ ?

- with grid search, find weight-decay λ
- $\tau = \frac{(1-p)l_i^2}{2N\lambda_i}$

Predictive Log-likelihood

(approximated by MC integration)

$$\log p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{X}, \mathbf{Y}) := \widehat{\log} \left(\frac{1}{T} \sum_{t=1}^T p(\mathbf{y}^* | \mathbf{x}^*, \omega_t) \right) \xrightarrow{T \rightarrow \infty}$$

$$\begin{aligned} & \log \int p(\mathbf{y}^* | \mathbf{x}^*, \boldsymbol{\omega}) q_{\theta}^*(\boldsymbol{\omega}) d\boldsymbol{\omega} \\ & \approx \log \int p(\mathbf{y}^* | \mathbf{x}^*, \boldsymbol{\omega}) p(\boldsymbol{\omega} | \mathbf{X}, \mathbf{Y}) d\boldsymbol{\omega} \\ & = \log p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{X}, \mathbf{Y}) \end{aligned}$$

- for regression : $\widetilde{\log p}(\mathbf{y}^* | \mathbf{x}^*, \mathbf{X}, \mathbf{Y}) = \text{logsumexp}\left(-\frac{1}{2}\tau \|\mathbf{y} - \mathbf{f}^{\hat{\omega}_t}(\mathbf{x}^*)\|^2\right) - \log T - \frac{1}{2}\log 2\pi + \frac{1}{2}\log \tau$

3.3.1 Uncertainty in Classification

(regression) find predictive uncertainty by "looking at the sample variance of multiple stochastic forward pass"

(classification)

- 1) variation ratios
- 2) predictive entropy
- 3) mutual information

1) Variation Ratio

$$\text{variation-ratio}[\mathbf{x}] := 1 - \frac{f_{\mathbf{x}}}{T}$$

- sample a label from softmax probabilities
- collecting a set of T labels y_t from multiple stochastic forward passes (of the same input)
- $f_{\mathbf{x}} = \sum_t 1[y^t = c^*]$, where $c^* = \arg \max_{c=1, \dots, C} \sum_t 1[y^t = c]$

variation-ratio can be seen as approximating the quantity : $1 - p(y = c^* | \mathbf{x}, \mathcal{D}_{\text{train}})$

2) Predictive entropy

$$\mathbb{H}[y | \mathbf{x}, \mathcal{D}_{\text{train}}] := - \sum_c p(y = c | \mathbf{x}, \mathcal{D}_{\text{train}}) \log p(y = c | \mathbf{x}, \mathcal{D}_{\text{train}})$$

- foundation in "information theory"
(captures the information contained in the predictive distribution)
- summing over all possible classes c that y can take

3) Mutual Information

$$\begin{aligned} \mathbb{I}[y, \boldsymbol{\omega} | \mathbf{x}, \mathcal{D}_{\text{train}}] & := \mathbb{H}[y | \mathbf{x}, \mathcal{D}_{\text{train}}] - \mathbb{E}_{p(\boldsymbol{\omega} | \mathcal{D}_{\text{train}})}[\mathbb{H}[y | \mathbf{x}, \boldsymbol{\omega}]] \\ & = - \sum_c p(y = c | \mathbf{x}, \mathcal{D}_{\text{train}}) \log p(y = c | \mathbf{x}, \mathcal{D}_{\text{train}}) + \mathbb{E}_{p(\boldsymbol{\omega} | \mathcal{D}_{\text{train}})} \left[\sum_c p(y = c | \mathbf{x}, \boldsymbol{\omega}) \log p(y = c | \mathbf{x}, \boldsymbol{\omega}) \right] \end{aligned}$$

- mutual information between prediction y and posterior (over the w)

Example (with binary output)

- case 1) all equal to 1 ((1,0), (1,0), ... (1,0))
- case 2) all equal to 0.5 ((0.5,0.5), (0.5,0.5), ... (0.5,0.5))
- case 3) half 1, half 0 ((1,0), (0,1), ... (1,0))

example	Predictive Uncertainty	Model Uncertainty
case 1	LOW (=0)	LOW (=0)
case 2	HIGH (=0.5)	LOW (=0)
case 3	HIGH (=0.5)	HIGH (=0.5)

in case 2)

- variation ratio & predictive entropy = 0.5
- mutual information = 0

3.3.2 Difficulties with the approach

simple! just several stochastic forward pass & find sample mean and variance

but have 3 shortcomings...

- 1) test time is scaled by T
(but not a real concern in a real world application ... transferring an input to a GPU)
- 2) model's uncertainty is not calibrated
(calibrated model : predictive probabilities match the empirical frequency of the data)
(GP's uncertainty is known to not be calibrated)
(lack of calibration = scale is different! can not compare...)
- 3) limitation of VI : underestimation of predictive variance
(but not a real concern in a real world application)

3.4 Approximate inference in complex models

apply it to CNN & RNN

3.4.1 Bayesian CNN

- also apply dropout after all convolution layers

3.4.2 Bayesian RNN

- inference with Bernoulli variational distributions for RNNs

$$f_y(h_T) = h_T W_y + b_y$$

$$\text{where } h_t = f_h(x_t, h_{t-1}) = \sigma(x_t W_h + h_{t-1} U_h + b_h)$$

- view the under RNN model as a probabilistic model

regard $\omega = \{W_h, U_h, b_h, W_y, b_y\}$

$$\begin{aligned} \int q(\omega) \log p(y | f_y^\omega(h_T)) d\omega &= \int q(\omega) \log p(y | f_y^\omega(f_h^\omega(x_T, h_{T-1}))) d\omega \\ &= \int q(\omega) \log p(y | f_y^\omega(f_h^\omega(x_T, f_h^\omega(\dots f_h^\omega(x_1, h_0) \dots))) d\omega \\ &\approx \log p(y | f_y^{\hat{\omega}}(f_h^{\hat{\omega}}(x_T, f_h^{\hat{\omega}}(\dots f_h^{\hat{\omega}}(x_1, h_0) \dots)))) \end{aligned}$$

where $\hat{\omega} \sim q(\omega)$

Final objective function :

$$\hat{\mathcal{L}}_{\text{MC}} = -\sum_{i=1}^N \log p\left(\mathbf{y}_i \mid \mathbf{f}_{\mathbf{y}}^{\hat{\omega}_i}\left(\mathbf{f}_{\mathbf{h}}^{\hat{\omega}_i}\left(\mathbf{x}_{i,T}, \mathbf{f}_{\mathbf{h}}^{\hat{\omega}_i}\left(\dots \mathbf{f}_{\mathbf{h}}^{\hat{\omega}_i}\left(\mathbf{x}_{i,1}, \mathbf{h}_0\right)\dots\right)\right)\right)\right)\right) + \text{KL}(q(\boldsymbol{\omega})\|p(\boldsymbol{\omega}))$$