

Uncertainty in Deep Learning PhD Thesis

Variational Inference: A Review for Statisticians

Bayesian Convolutional Neural Networks with Bernoulli Approximate Variational Inference *ICLR*

Dropout as a Bayesian Approximation: Insights and Applications *ICML*

Deep Bayesian Active Learning with Image data *ICML 2017*

What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? *NIPS 2017*

2018.4.4

Given training inputs $\{X_1, \dots, X_N\}$ and outputs $\{y_1, \dots, y_N\}$, estimate that a function $y = f(x)$ that is likely to have generated outputs

prior: $p(f)$ likelihood: $p(Y|f, X)$

posterior: $p(f|X, Y)$

We can predict an output for a new input x^* by integrate all possible functions f

$$p(y^*|x^*, X, Y) = \int p(y^*|f^*)p(f^*|x^*, X, Y)df^*.$$

To approximate the integral, assume that the model depends on a finite set of random variables ω alone

$$p(y^*|x^*, X, Y) = \int p(y^*|f^*)p(f^*|x^*, \omega)p(\omega|X, Y) df^* d\omega.$$

Variational inference

Let $\mathbf{X} = x_{1:n}$ be a set of observed variables and $\mathbf{Z} = z_{1:m}$ be a set of latent variables.

We specify a family of densities over the latent variables. Now our goal is to find an approximation $q(\mathbf{z})$ to the exact conditional $p(\mathbf{z}|\mathbf{x})$.

We thus minimize the Kullback–Leibler (KL) divergence, intuitively a measure of similarity between two distributions

$$q^*(\mathbf{z}) = \arg \min_{q(\mathbf{z}) \in \mathcal{Q}} \text{KL}(q(\mathbf{z}) \| p(\mathbf{z} | \mathbf{x})).$$

$$\text{KL}(q(\mathbf{z}) \| p(\mathbf{z} | \mathbf{x})) = \mathbb{E}[\log q(\mathbf{z})] - \mathbb{E}[\log p(\mathbf{z} | \mathbf{x})],$$

$$\text{KL}(q(\mathbf{z}) \| p(\mathbf{z} | \mathbf{x})) = \mathbb{E}[\log q(\mathbf{z})] - \mathbb{E}[\log p(\mathbf{z}, \mathbf{x})] + \log p(\mathbf{x}).$$

Because we cannot compute the KL, we optimize an alternative objective that is equivalent to the KL up to an added constant

$$\text{ELBO}(q) = \mathbb{E}[\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}[\log q(\mathbf{z})].$$

This function is called the evidence lower bound(ELBO)

We rewrite the ELBO as a sum of the expected log likelihood of the data and the KL divergence between the prior $p(\mathbf{z})$ and $q(\mathbf{z})$

$$\begin{aligned}\text{ELBO}(q) &= \mathbb{E}[\log p(\mathbf{z})] + \mathbb{E}[\log p(\mathbf{x} | \mathbf{z})] - \mathbb{E}[\log q(\mathbf{z})] \\ &= \mathbb{E}[\log p(\mathbf{x} | \mathbf{z})] - \text{KL}(q(\mathbf{z}) \| p(\mathbf{z})).\end{aligned}$$

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

We define an approximating variational distribution $q(\omega)$ to be as close as possible to the posterior distribution by minimizing the KL divergence.

$$\text{KL}(q(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y})).$$

Minimising the Kullback–Leibler divergence is equivalent to maximising the log evidence lower bound

$$\mathcal{L}_{\text{VI}} := \int q(\omega)p(\mathbf{F}|\mathbf{X}, \omega) \log p(\mathbf{Y}|\mathbf{F})d\mathbf{F}d\omega - \text{KL}(q(\omega)||p(\omega))$$

resulting in the approximate predictive distribution

$$q(\mathbf{y}^*|\mathbf{x}^*) = \int p(\mathbf{y}^*|\mathbf{f}^*)p(\mathbf{f}^*|\mathbf{x}^*, \omega)q(\omega)d\mathbf{f}^* d\omega.$$

Dropout as Approximate Variational Inference in Bayesian Neural Network

Let $\hat{\mathbf{y}}$ be the output of a NN with L layers and a loss function $E(\cdot, \cdot)$ such as the softmax loss. We denote by \mathbf{W}_i the NN's weight matrices of dimensions $K_i \times K_{i-1}$, and by \mathbf{b}_i the bias vectors of dimensions K_i for each layer $i = 1, \dots, L$. In addition to regularization, we have the optimization problem

$$\mathcal{L}_{\text{dropout}} := \frac{1}{N} \sum_{i=1}^N E(\mathbf{y}_i, \hat{\mathbf{y}}_i) + \lambda \sum_{i=1}^L (\|\mathbf{W}_i\|_2^2 + \|\mathbf{b}_i\|_2^2).$$

Given weight matrices \mathbf{W}_i and bias vectors \mathbf{b}_i for layer i , we often place standard matrix Gaussian prior distributions over the weight matrices, $p(\mathbf{W}_i)$:

$$\mathbf{W}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

We are thus interested in the posterior over the weights given our observables \mathbf{X}, \mathbf{Y} : $p(\omega|\mathbf{X}, \mathbf{Y})$. This posterior is not tractable for a Bayesian NN, and we use variational inference to approximate it.

We define our approximating variational distribution $q(\omega_i)$ for every layer i as

$$\mathbf{W}_i = \mathbf{M}_i \cdot \text{diag}([\mathbf{z}_{i,j}]_{j=1}^{K_i})$$

$$\mathbf{z}_{i,j} \sim \text{Bernoulli}(p_i) \text{ for } i = 1, \dots, L, j = 1, \dots, K_{i-1}.$$

Here $z_{i,j}$ are Bernoulli distributed random variables with some probabilities p_i , and M_i are variational parameters to be optimized. The binary variable $z_{i,j} = 0$ corresponds to unit j in layer $i - 1$ being dropped out as an input to the i 'th layer.

The integral in eq is intractable, and we approximate the integral with Monte Carlo integration over ω .

$$\mathcal{L}_{\text{VI}} := \int q(\omega) p(\mathbf{F}|\mathbf{X}, \omega) \log p(\mathbf{Y}|\mathbf{F}) d\mathbf{F} d\omega - \text{KL}(q(\omega)||p(\omega))$$



$$\hat{\mathcal{L}}_{\text{VI}} := \sum_{i=1}^N E(y_i, \hat{\mathbf{f}}(\mathbf{x}_i, \hat{\omega}_i)) - \text{KL}(q(\omega)||p(\omega)) \quad \hat{\omega}_i \sim q(\omega)$$

$$\hat{\mathcal{L}}_{\text{VI}} := \sum_{i=1}^N E(y_i, \hat{\mathbf{f}}(\mathbf{x}_i, \hat{\omega}_i)) - \text{KL}(q(\omega) || p(\omega)) \quad \hat{\omega}_i \sim q(\omega)$$

$E(\cdot, \cdot)$ being the softmax loss (for a softmax likelihood).

Note that sampling from $q(W_i)$ is identical to performing dropout on layer i in a network whose weights are $(M_i)_{i=1}^L$. The second term can be approximated, resulting in the same objective as dropout.

$$\mathcal{L}_{\text{dropout}} := \frac{1}{N} \sum_{i=1}^N E(y_i, \hat{y}_i) + \lambda \sum_{i=1}^L (\|\mathbf{W}_i\|_2^2 + \|\mathbf{b}_i\|_2^2).$$

Dropout and Bayesian NNs, in effect, result in the same model parameters that best explain the data.

Predictions in this model follow eq replacing the posterior $p(\omega|D_{train})$ with the approximate Posterior $q(\omega)$. We can approximate the integral with Monte Carlo integration:

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



$$\begin{aligned} p(y = c|\mathbf{x}, \mathcal{D}_{train}) &= \int p(y = c|\mathbf{x}, \omega)p(\omega|\mathcal{D}_{train})d\omega \\ &\approx \int p(y = c|\mathbf{x}, \omega)q_{\theta}^*(\omega)d\omega \\ &\approx \frac{1}{T} \sum_{t=1}^T p(y = c|\mathbf{x}, \hat{\omega}_t) \end{aligned}$$

With $\hat{\omega}_t \sim q_{\theta}^*(\omega)$, where $q_{\theta}(\omega)$ is the Dropout distribution.

Deep Bayesian Active Learning

Deep learning poses several difficulties when used in an active learning setting.

- Active learning methods generally rely on being able to learn and update models from small amounts of data.
- Many AL acquisition functions rely on model uncertainty, yet deep learning methods rarely represent such model uncertainty.

Acquisition Functions and their Approximations

Given a model M , pool data D_{pool} , and inputs $x \in D_{pool}$, an acquisition function $a(x, M)$ is a function of x that the AL system uses to decide where to query next:

$$x^* = \operatorname{argmax}_{x \in \mathcal{D}_{pool}} a(x, \mathcal{M}).$$

For classification, several acquisition functions are available:

1. Choose pool points that maximize the predictive entropy (*Max Entropy*)

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

2. Choose pool points that are expected to maximize the mutual information between predictions and model posterior (*BALD*)

$$\mathbb{I}[y, \omega|\mathbf{x}, \mathcal{D}_{\text{train}}] = \mathbb{H}[y|\mathbf{x}, \mathcal{D}_{\text{train}}] - \mathbb{E}_{p(\omega|\mathcal{D}_{\text{train}})} [\mathbb{H}[y|\mathbf{x}, \omega]]$$

$\mathbb{H}[y|x, w]$ is the entropy of y given model weight ω .

3. Maximize the *Variation Ratios*

$$\text{variation-ratio}[\mathbf{x}] := 1 - \max_y p(y|x, \mathcal{D}_{\text{train}})$$

4. Maximize mean standard deviation (*Mean STD*)

$$\sigma_c = \sqrt{\mathbb{E}_{q(\boldsymbol{\omega})}[p(y = c|\mathbf{x}, \boldsymbol{\omega})^2] - \mathbb{E}_{q(\boldsymbol{\omega})}[p(y = c|\mathbf{x}, \boldsymbol{\omega})]^2}$$
$$\sigma(\mathbf{x}) = \frac{1}{C} \sum_c \sigma_c$$

5. Random acquisition (baseline)

We can approximate each of these acquisition functions using our approximate distribution $q_{\theta}^*(\omega)$. For BALD, for example, we can write the acquisition function as follows :

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

$\mathbb{I}[y, w | x, D_{\text{train}}]$ can be approximated in our setting using

$$p(y = c | \mathbf{x}, \mathcal{D}_{\text{train}}) = \int p(y = c | \mathbf{x}, \omega) p(\omega | \mathcal{D}_{\text{train}}) d\omega$$

then,

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

Swapping the posterior $p(\mathbf{w}|D_{train})$ with our approximate posterior $q_{\theta}^*(\omega)$ and through MC sampling, we then have

$$\begin{aligned}
&\approx - \sum_c \int p(y = c|\mathbf{x}, \omega) q_{\theta}^*(\omega) d\omega \\
&\quad \cdot \log \int p(y = c|\mathbf{x}, \omega) q_{\theta}^*(\omega) d\omega \\
&\quad + \mathbb{E}_{q_{\theta}^*(\omega)} \left[\sum_c p(y = c|\mathbf{x}, \omega) \log p(y = c|\mathbf{x}, \omega) \right] \\
&\approx - \sum_c \left(\frac{1}{T} \sum_t \hat{p}_c^t \right) \log \left(\frac{1}{T} \sum_t \hat{p}_c^t \right) \\
&\quad + \frac{1}{T} \sum_{c,t} \hat{p}_c^t \log \hat{p}_c^t =: \hat{\mathbb{I}}[y, \omega|\mathbf{x}, \mathcal{D}_{train}]
\end{aligned}$$

with \hat{p}_c^t the probability of input x with model parameters $\hat{\mathbf{w}}_t \sim q_{\theta}^*(\omega)$ to take class c :

$$\hat{\mathbf{p}}^t = [\hat{p}_1^t, \dots, \hat{p}_C^t] = \text{softmax}(\mathbf{f}^{\hat{\mathbf{w}}_t}(\mathbf{x})).$$

We then have

$$\begin{aligned}\widehat{\mathbb{I}}[y, \boldsymbol{\omega} | \mathbf{x}, \mathcal{D}_{\text{train}}] &\xrightarrow{T \rightarrow \infty} \mathbb{H}[y | \mathbf{x}, q_{\theta}^*] - \mathbb{E}_{q_{\theta}^*}(\boldsymbol{\omega}) [\mathbb{H}[y | \mathbf{x}, \boldsymbol{\omega}]] \\ &\approx \mathbb{I}[y, \boldsymbol{\omega} | \mathbf{x}, \mathcal{D}_{\text{train}}],\end{aligned}$$

Comparison of various acquisition functions

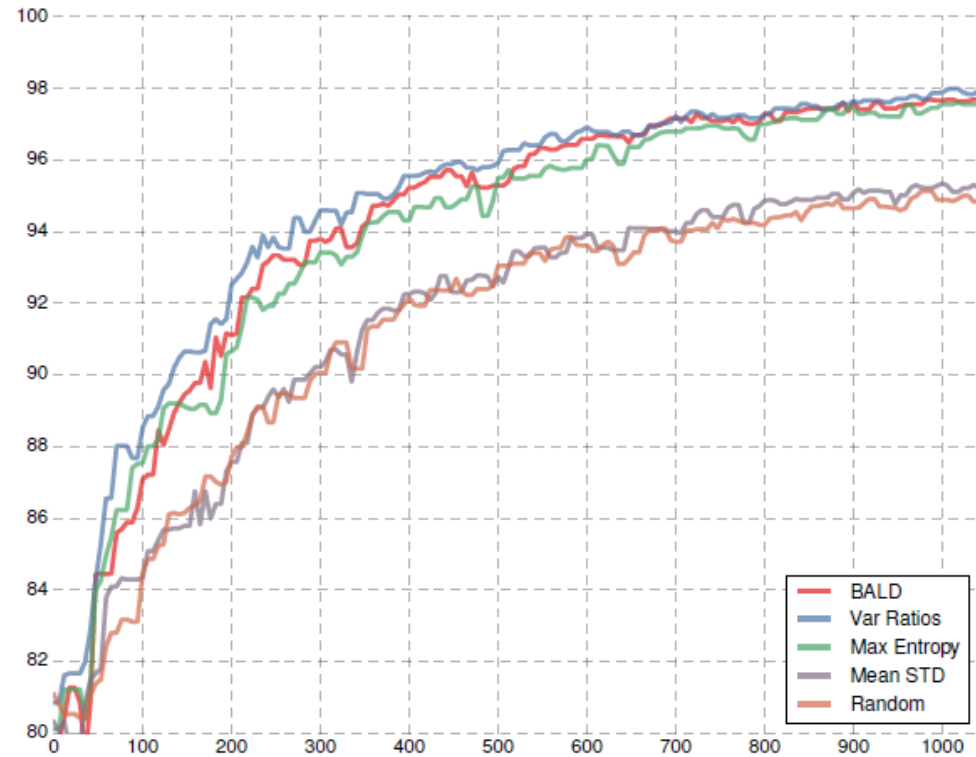


Figure 1. MNIST test accuracy as a function of number of acquired images from the pool set (up to 1000 images, using validation set size 100, and averaged over 3 repetitions). Four acquisition functions (*BALD*, *Variation Ratios*, *Max Entropy*, and *Mean STD*) are evaluated and compared to a *Random* acquisition function.

Comparison to current active learning techniques with image data

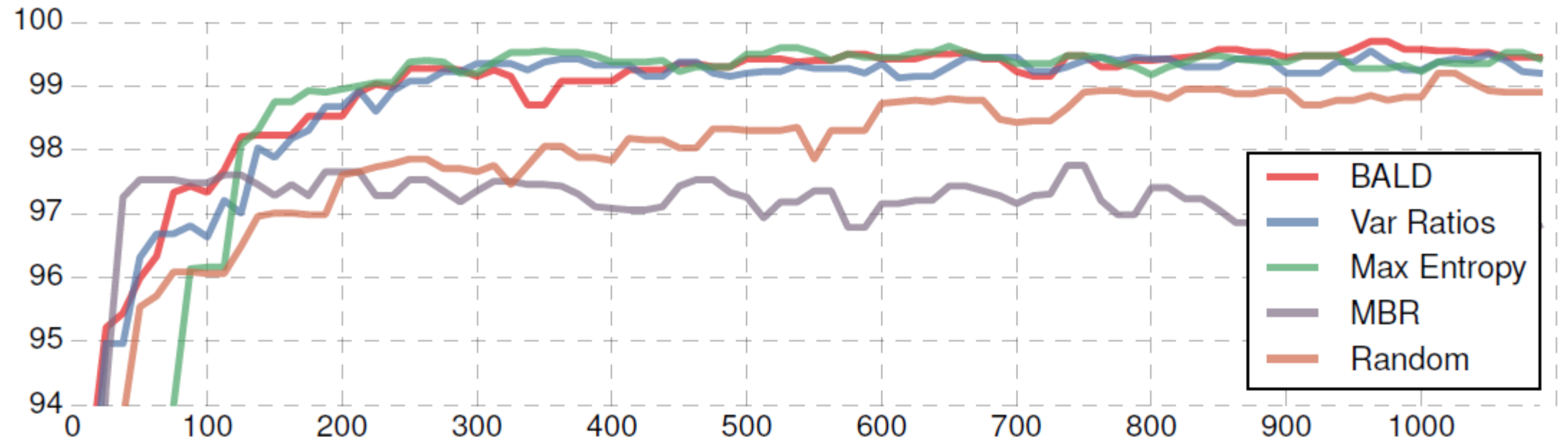


Figure 3. MNIST test accuracy (two digit classification) as a function of number acquired images, compared to a current technique for active learning of image data: MBR (Zhu et al., 2003).

Comparison to semi-supervised learning

Technique	Test error
Semi-supervised:	
Semi-sup. Embedding (Weston et al., 2012)	5.73%
Transductive SVM (Weston et al., 2012)	5.38%
MTC (Rifai et al., 2011)	3.64%
Pseudo-label (Lee, 2013)	3.46%
AtlasRBF (Pitelis et al., 2014)	3.68%
DGN (Kingma et al., 2014)	2.40%
Virtual Adversarial (Miyato et al., 2015)	1.32%
Ladder Network (Γ -model) (Rasmus et al., 2015)	1.53%
Ladder Network (full) (Rasmus et al., 2015)	0.84%
Active learning with various acquisitions:	
Random	4.66%
BALD	1.80%
Max Entropy	1.74%
Var Ratios	1.64%