# Bayesian deep learning,
## MASH & IASD

Julyan ARBEL

Feb 5, 2021.

# Bayesian neural network, predictive distribution:

$$p(t \mid x, \mathcal{D}) = \int p(t \mid x, w) \, p(w \mid \mathcal{D}) \, dw$$

Approximations to the posterior $p(w \mid \mathcal{D})$:

- Laplace (yesterday)
- Variational inference
- Monte Carlo dropout

VI (Blundell, 2015)

parameters are weights $w = (w_1, \ldots, w_{\overline{W}})$

Approximate family : Gaussian : $N(\mu, \sigma^2)$

$w_i$ approximated by $\theta_i = (\mu_i, \sigma_i^2)$, $i = 1, \ldots, \overline{W}$

$$q_\theta(w) = \prod_{i=1}^{\overline{W}} N(w_i \mid \mu_i, \sigma_i^2)$$

$$Q = \{ q_\theta(w), \ \mu_i \in \mathbb{R}, \ \sigma_i > 0 \}$$

$$KL\left( q_\theta(w) \parallel p(w \mid X, Y) \right) \qquad X, Y \begin{cases} X = (X_i) \\ Y = (Y_i) \end{cases} i = 1 \ldots N$$

$$= \int q_\theta(w) \log\left( \frac{q_\theta(w)}{p(w \mid X, Y)} \right) dw$$

$$= \int q_\theta(w) \log\left( \frac{q_\theta(w) \, p(Y \mid X)}{p(w) \, p(Y \mid w, X)} \right) dw$$

$$= \boxed{-} \int q_\theta(w) \log \underbrace{p(Y|w,X)}_{\text{likelih.}} + \underbrace{KL\left(q_\theta(w) \| p(w)\right)}_{\text{KL for prior}} + \underbrace{\log p(Y|X)}_{\text{evidence}}$$

$$\text{expected log-like.}$$

$$= \boxed{-} \underbrace{\mathcal{L}_{VI}(\theta)}_{} + \underbrace{\log p(Y|X)}_{c^+}$$

Since $KL \geq 0$, $\mathcal{L}_{VI}(\theta) \leq \log p(Y|X)$

so $\mathcal{L}_{VI}$ is called evidence lower bound ELBO

Minimize VI $\iff$ Maximize ELBO

ELBO: $\mathcal{L}_{VI}(\theta) = \mathbb{E}_{q_\theta(w)}\left[\log p(Y|X,w)\right] - \underbrace{KL\left(q_\theta(w) \| p(w)\right)}_{\text{closed-form}}$

$\mathcal{L}_{VI}(\theta) = \sum_{i=1}^{N} \mathbb{E}_{q_\theta(w)}\left[\log p(Y_i|w,X_i)\right] - KL$

approximate by sampling : $\hat{w} \sim q_\theta(\cdot)$

In order to estimate the gradient of $\mathcal{L}_{VI}(\theta)$, use the reparametrization trick.

idea : $w = g(\theta, \varepsilon)$ , $\begin{cases} g \text{ is deterministic} \\ \varepsilon \perp\!\!\!\perp \theta \end{cases}$

$$w_j \sim N(w_j \mid \mu_j, \sigma_j^2) = q_{\theta_j}(w)$$

$$w_j = g(\theta_j, \varepsilon_j) = \mu_j + \sigma_j \varepsilon_j \text{ , with } \varepsilon_j \sim N(0, 1).$$

$$ELBO(\theta) = \mathcal{L}_{VI}(\theta) \approx \underbrace{\sum_{i=1}^{N} \lg p(y_i \mid y(x_i, g(\theta, \hat{\varepsilon}_i))}_{\text{MC approx of } \mathbb{E}} - KL$$

Algorithm : Stochastic VI

☑ Given $X, Y$ data, $\eta$: learning rate, initialise $\theta$

☑ Repeat: Sample $\begin{cases} \varepsilon_j \sim p(\varepsilon) & , \quad j \in S \\ S \text{ subset from } \{1, \ldots, N\} \text{ of size } M \end{cases}$

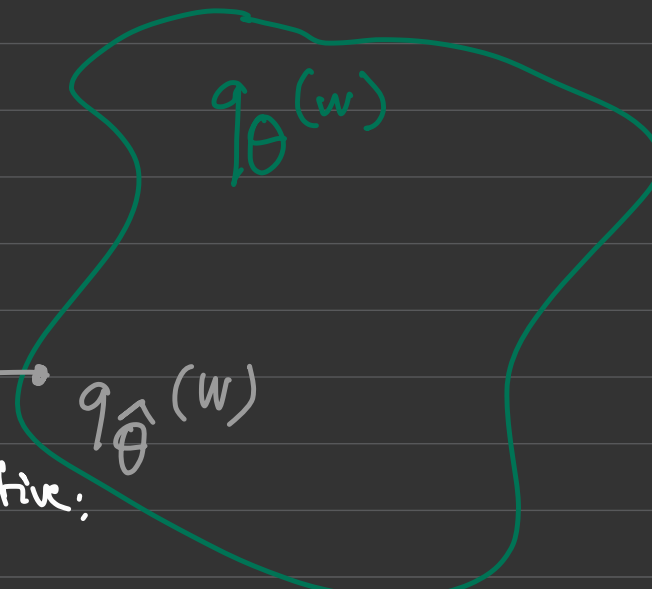Stochastic derivative estimator wrt $\theta$:

$$\widehat{\Delta\theta} = -\frac{N}{M} \sum_{i \in S} \frac{\partial}{\partial \theta} \log p\left(Y_i \mid \gamma(g(\theta, \hat{\varepsilon}_i), X_i)\right) + \frac{\partial KL}{\partial \theta}$$

avail. by reparam. trick.    closed-form.

$$\theta \leftarrow \theta + \eta \widehat{\Delta\theta}.$$

Until (some) convergence.

$$\mathcal{L}_{VI}(\theta)$$

$q_\theta(w)$

$p(w \mid x, y)$

$q_{\hat\theta}(w)$

$q_{\hat\theta}(w)$ can be plugged-in predictive:

$$p(\epsilon \mid x, \mathcal{D}) \approx \int \underbrace{p(\epsilon \mid x, w)}_{\downarrow} \underbrace{q_{\hat\theta}(w)}\, dw$$

can be made Gaussian
by some Taylor expansion
as yesterday

Gaussian approx by VI

Monte Carlo dropout ( _Gal_, Ghahramani, ICML, 2016)

based on: Dropout (Hinton, 2012)

☒ training a NN with dropout

⟺ training a BNN with variational posterior $q_\theta(w)$

☒ MC dropout:

sampling several passes of NN with dropout

⟺ MC approx. inference with $q_\theta(w)$.