

Bayesian deep learning /
NASH & IASD

Julian ARBEL
Feb 5, 2021

Bayesian NN under a Laplace approximation

Chap 5.7 of Bishop

Univariate regression $t \in \mathbb{R}$

Gaussian $p(t|x)$ to be Gaussian

the mean = $y(\underset{\text{input}}{x}, \underset{\text{weight}}{w})$ output of NN

variance β^{-1} β precision

Data $\mathcal{D} = \{(x_n, t_n), n=1, \dots, N\}$

Model $p(t|x) = N(t | y(x, w), \beta^{-1})$

Prior $p(w|\alpha) = N(w | 0, \alpha^{-1} I)$

Posterior $p(w|\mathcal{D}, \alpha, \beta) \propto \prod_{n=1}^N N(\underset{\text{red underline}}{t_n} | y(x_n, w), \beta^{-1}) p(w|\alpha)$

w_{MAP} : numerically find local optimum.

$$\log p(w | \mathcal{D}, \alpha, \beta) = - \underbrace{\frac{\alpha}{2} w^T w}_{\text{penalty to the MLE}} - \underbrace{\frac{\beta}{2} \sum_{n=1}^N (y(x_n, w) - t_n)^2}_{\text{SSE}} + c^t$$

Laplace approx. of the posterior:

$$q(w | \mathcal{D}, \alpha, \beta) = N(w | w_{\text{MAP}}, A^{-1})$$

$$A = - \nabla \nabla \log p = \alpha I + \beta H, \text{ with } H = \text{Hess}(\text{SSE})$$

approx.
predictive

$$p(t | x, \mathcal{D}) = \int p(t | x, w) q(w | \mathcal{D}, \alpha, \beta) dw \quad (*)$$

$\underbrace{w}_{\sim (t | y(x, w), \beta^{-1})}$

Taylor approximation for NN:

$$y(x, w) \approx y(x, w_{\text{MAP}}) + g^T (w - w_{\text{MAP}})$$

$$g = \nabla_w y(x, w) |_{w=w_{\text{MAP}}}$$

$$p(t | x, w, \beta) \approx N(t | y(x, w_{\text{MAP}}) + g^T (w - w_{\text{MAP}}), \beta^{-1})$$

plug this in (*):

Plus use

$$\begin{cases} p(x) = N(x | \mu, \Lambda^{-1}) \\ p(y | x) = N(y | Ax + b, L^{-1}) \end{cases}$$

$$\Rightarrow p(y) = N(y | A\mu + b, \underbrace{L^{-1} + A\Lambda^{-1}A^T}_{\text{epistemic}})$$

$$p(t | x, \alpha, \beta) \approx N(t | y(x, w_{\text{MAP}}), \sigma^2(x))$$

$$\sigma^2(x) = \underbrace{\beta^{-1}}_{\text{aleatoric}} + \underbrace{g^T A^{-1} g}_{\text{x-dependent : epistemic}}$$

aleatoric x-dependent : epistemic.

α, β hyperparameters

Marginal likelihood $p(\mathcal{D} | \alpha, \beta)$, aka evidence

$\leadsto w$ is integrated out.

$$(\hat{\alpha}, \hat{\beta}) = \underset{(\alpha, \beta)}{\operatorname{argmax}} p(\mathcal{D} | \alpha, \beta).$$

$$p(\mathcal{D} | \alpha, \beta) = \int p(\mathcal{D} | w, \beta) p(w | \alpha) dw$$

Laplace approximation: $z_0 \operatorname{argmax}$, $A = -\nabla \nabla \log f(z) |_{z=z_0}$

$$\begin{aligned} \int f(z) dz &= f(z_0) \int \exp\left(-\frac{1}{2} (z - z_0)^T A (z - z_0)\right) dz \\ &= f(z_0) \frac{(2\pi)^{n/2}}{|A|^{1/2}} \quad n = \#(z) \end{aligned}$$

$$\Rightarrow \log p(\mathcal{D} | \alpha, \beta) = - \underbrace{E(w_{\text{MAP}})}_{\text{regularized SSE}} - \frac{1}{2} \log |A| + \frac{W}{2} \log \alpha + \frac{N}{2} \log \beta + c$$

per.

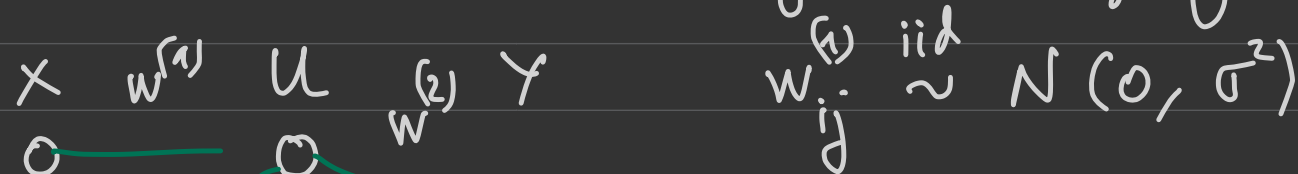
optimization in (α, β) is done by analogy with linear regression. \rightarrow Bishop 5.7.

Other techniques for practical BNN:

variational inference: with various assumptions on approx family.

Distributional properties of BNN

Neal's result: 1-hidden layer NN, fully-connected



$$w_{ij}^{(2)} \sim N(0, \sigma_H^2)$$

ϕ

$$\|X\|_2^2 = \sum x_h^2$$

$m = 1 \dots H$

$$g_m(x), \text{ pre non linearity} = \sum_{h=1}^{H_0} w_{mh} \cdot x_h \stackrel{\text{iid}}{\sim} N(0, \|X\|_2^2 \sigma^2)$$

post non linearity: $h_m(x) = \phi(g_m(x))$

$$y = \sum_{m=1}^H \underbrace{w_m^{(2)} h_m(x)}_{Y_m} \quad Y_m \text{ are iid}$$

$$\forall m: E[Y_m] = E[w_m^{(2)}] E[h_m(x)] = 0 \quad E[h_m(x)] = 0$$

$$V[Y_m] = \underbrace{E[W_m^{(2)2}]}_{\sigma_H^2} \underbrace{E[h_m^2(x)]}_{c^2}$$

$$y \rightarrow N(0, H \underbrace{\sigma_H^2}) \quad \text{by CLT}$$

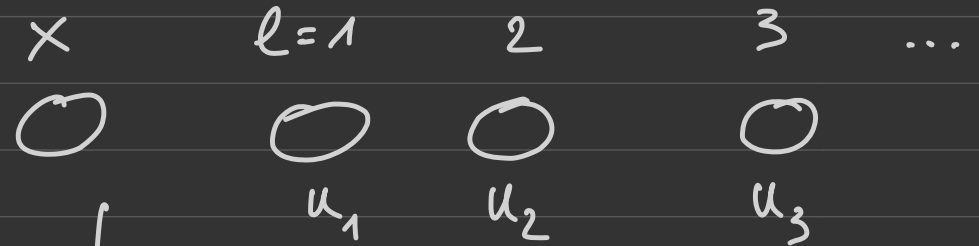
$$\sigma_H^2 = \frac{1}{4}$$

Gaussian limit of output when width $\rightarrow \infty$.

Extended to any deeper NN in 2018.

Another distributional property of BNN

what is the distrib. of units under Gaussian prior
on weights?



Assumption on prior:

indep priors, centered

eg: $w_{ij}^l \stackrel{iid}{\sim} N(0, \sigma^2)$

$w \sim N$

Assumption on nonlinearity ϕ

envelope condition for ϕ



Theorem Conditional on x , the prior for unit at
layer l is Sub Weibull $(l/2)$.

More on Sub-Weibull: generalizes sub-Gaussian &
sub Exp properties

Let X be a.v., let $\|X\|_k = (E[X^k])^{1/k}$

X is sub-Gaussian if $\exists C > 0$ s.t.

$$\forall k, \|X\|_k \leq C \sqrt{k} = C k^{1/2}$$

X is sub-Exponential if $\exists C > 0$ s.t.

$$\forall k, \|X\|_k \leq C \underline{k}$$

X is sub-Weibull(θ) $\forall k, \|X\|_k \leq C k^\theta$

X	1	2	3	$\in \text{sub W}(\frac{\ell}{2})$.
\bigcirc	\bigcirc	\bigcirc	\bigcirc	
w				

\bigcirc	\bigcirc	\bigcirc	\bigcirc
------------	------------	------------	------------

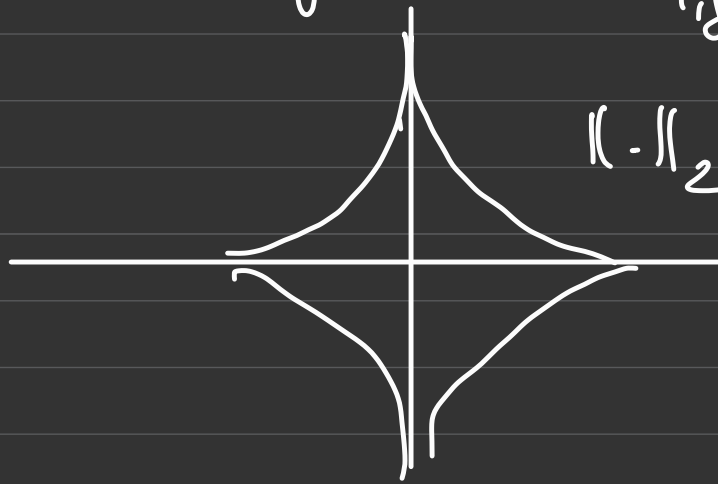
\bigcirc	\bigcirc	\bigcirc	\bigcirc
------------	------------	------------	------------

$$N = \text{sub G} = \text{sub W}(\frac{1}{2})$$

$$\text{sub W}(\frac{\ell}{2}) = \text{sub W}(1) = \text{sub E}$$

density for Weibull^(θ) $g(w) \propto e^{-|w|^{1/\theta}}$
induces a penalty $\rightarrow |w|^{1/\theta}$ $\theta = \frac{\ell}{2}$

Layer ℓ $\sum_{i,j} |w_{ij}^\ell|^{2/\ell} = \|w^\ell\|_{2/\ell}$



$\|\cdot\|_{2/\ell}$ unit ball in \mathbb{R}^2 .

Vladimirava, 2018

PyTorch, pyro.