

Bayesian neural networks

So far, we have used ML techniques to find weights and biases. Let's now consider a Bayesian framework:

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

Annotations:  
 -  $t$ : single output (for simplicity)  
 -  $\vec{x}$ : input  
 -  $\vec{w}$ : weights/biases  
 -  $\beta$ : precision  
 -  $y(\vec{x}, \vec{w})$ : neural network model  
 -  $\beta^{-1}$ :  $\sigma^2$

Choose a prior:

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

$N$  observations:  $\vec{x}_1, \dots, \vec{x}_N$  (continuous vars)  
 Target values:  $\mathcal{D} = \{t_1, \dots, t_N\}$

Then  $\mathcal{L} = p(\mathcal{D} | \vec{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y(\vec{x}_n, \vec{w}), \beta^{-1})$

and the posterior is given by

$$p(\vec{w} | \mathcal{D}, \alpha, \beta) \sim \underbrace{p(\vec{w} | \alpha) p(\mathcal{D} | \vec{w}, \beta)}_{\text{Gaussian}}$$

non-gaussian since  $y(\vec{x}, \vec{w})$  is a non-linear f'n of  $\vec{w}$   
 need to use Laplace approx'n

Consider

$$\log P(\vec{w} | \mathcal{D}, \alpha, \beta) = -\frac{\alpha}{2} \vec{w}^T \vec{w} - \underbrace{\frac{\beta}{2} \sum_{n=1}^N (y(\vec{x}_n, \vec{w}) - t_n)^2}_{\text{BE}} + \text{const}(\vec{w})$$

↳ can find  $\vec{w}_{\text{MAP}}$  which maximizes this (assume  $\alpha$  &  $\beta$  are fixed for now) by error backpropagation + conjugate grad. Note that  $\vec{w}_{\text{MAP}}$  is a local max in general.

Next, compute

$$A_{ij} = -\frac{\partial^2}{\partial w_i \partial w_j} \log P(\vec{w} | \mathcal{D}, \alpha, \beta) = \alpha \delta_{ij} + \beta H_{ij} = \frac{\partial^2 E}{\partial w_i \partial w_j}, \text{ Hessian}$$

needed for Laplace approx'n

already discussed how to compute it efficiently

Recall that under Laplace,

$$f(\vec{z}) \Rightarrow q_f(\vec{z}) = \mathcal{N}(\vec{z} | \vec{z}_0, A^{-1})$$

So, the posterior becomes gaussian:

$$P(\vec{w} | \mathcal{D}, \alpha, \beta) \Rightarrow q_f(\vec{w} | \mathcal{D}) = \mathcal{N}(\vec{w} | \vec{w}_{\text{MAP}}, A^{-1}).$$

The predictive distribution is

$$p(t|\vec{x}, \mathcal{D}) = \int d\vec{w} \underbrace{p(t|\vec{x}, \vec{w})}_{\text{non-gaussian}} \underbrace{q(\vec{w}|\mathcal{D})}_{\text{gaussian}}$$

↑  
explicit  $\alpha, \beta$  dependence  
suppressed for brevity

If  $q(\vec{w}|\mathcal{D})$  is sufficiently narrow,

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

where

$$\vec{g} = \nabla_{\vec{w}} y(\vec{x}, \vec{w}) \Big|_{\vec{w} = \vec{w}_{\text{MAP}}}$$

Then  $p(t|\vec{x}, \vec{w}, \beta) \approx \mathcal{N}(t|y(\vec{x}, \vec{w}_{\text{MAP}}) + \vec{g}^T \cdot (\vec{w} - \vec{w}_{\text{MAP}}), \beta^{-1})$ .

Recall:

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

$$\begin{aligned} p(\vec{y}) &= \int d\vec{x} p(\vec{y}|\vec{x}) p(\vec{x}) = \\ &= \mathcal{N}(\vec{y}|A\vec{\mu} + \vec{b}, L^{-1} + A\Lambda^{-1}A^T) \end{aligned}$$

Here,

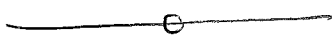
$$\begin{cases} \vec{x} \rightarrow \vec{w} \\ p(\vec{x}) \rightarrow q(\vec{w}|\mathcal{D}) \rightarrow \begin{cases} \vec{\mu} = \vec{w}_{\text{MAP}}, \\ \Lambda^{-1} = A^{-1} \quad \text{scalar} \end{cases} \\ \vec{y} \rightarrow t \\ p(\vec{y}|\vec{x}) \rightarrow \begin{cases} \vec{b} \rightarrow y(\vec{x}, \vec{w}_{\text{MAP}}) - \vec{g}^T \cdot \vec{w}_{\text{MAP}} \\ A \rightarrow \vec{g}^T, \quad L^{-1} \rightarrow \beta^{-1} \end{cases} \end{cases}$$

Consequently,

$$p(t|\vec{x}, \mathcal{D}, \alpha, \beta) = \mathcal{N}(t | y(\vec{x}, \vec{w}_{\text{MAP}}), \sigma^2(\vec{x}))$$

restored explicit dependence

$$\left\{ \begin{aligned} A\vec{\mu} + \vec{b} &\rightarrow \vec{g}^T \cdot \vec{w}_{\text{MAP}} + y(\vec{x}, \vec{w}_{\text{MAP}}) - \vec{g}^T \cdot \vec{w}_{\text{MAP}} = \\ &= y(\vec{x}, \vec{w}_{\text{MAP}}) \quad \text{uncertainty due to uncertainty in } \vec{w} \\ L^{-1} + A\Lambda^{-1}A^T &\rightarrow \underbrace{\beta^{-1}}_{\text{uncertainty due to noise in target variables}} + \underbrace{\vec{g}^T A^{-1} \vec{g}}_{\sum_{ij} g_i A^{-1}_{ij} g_j} \equiv \underbrace{\sigma^2(\vec{x})}_{\substack{\vec{x}\text{-dependence through} \\ \vec{g} = \vec{g}(\vec{x})}} \end{aligned} \right.$$



How to find  $\alpha$  and  $\beta$ ?

Consider  $p(\mathcal{D}|\alpha, \beta) = \int d\vec{w} \underbrace{p(\mathcal{D}|\vec{w}, \beta)}_{\mathcal{Z}} \underbrace{p(\vec{w}|\alpha)}_{\text{prior}}$

evidence for  $\alpha, \beta$   
 [normalization factor for posterior in weights]

Recall that under Laplace approx'n:

$$\int d\vec{z} \underbrace{f(\vec{z})}_{M \text{ dims}} \approx f(\vec{z}_0) \frac{(2\pi)^{M/2}}{|A|^{1/2}}, \text{ where}$$

$$A = -\vec{\nabla} \vec{\nabla} \log f(\vec{z}) \Big|_{\vec{z} = \vec{z}_0}$$

Here,  $\log p(\mathcal{D}|\alpha, \beta) \approx -\frac{\beta}{2} \sum_{n=1}^N (y(\vec{x}_n, \vec{w}_{\text{MAP}}) - t_n)^2 - \frac{\alpha}{2} \vec{w}_{\text{MAP}}^T \vec{w}_{\text{MAP}} - \frac{N}{2} \log(2\pi\beta^{-1}) - \frac{W}{2} \log(2\pi\alpha^{-1})$   $\oplus$

$\oplus \frac{W}{2} \log(2\pi) - \frac{1}{2} \log |A| \ominus$

$$\textcircled{=} - E(\vec{w}_{\text{MAP}}) - \frac{N}{2} \log(2\pi) + \frac{N}{2} \log \beta + \frac{W}{2} \log \alpha - \\ - \frac{1}{2} \log |A| \quad \underline{\underline{=}} \quad W = \text{total \# prms in } \vec{w}$$

— This is very similar to (3.82) & (3.86).

We need to maximize  $\log p(\mathcal{D} | \alpha, \beta)$  to get point estimates  $\hat{\alpha}$  &  $\hat{\beta}$ .  
(ML)

Define  $(\beta H) \vec{u}_i = \lambda_i \vec{u}_i$   
 $\uparrow$   
 Hessian at  $\vec{w}_{\text{MAP}}$   $i=1, \dots, W$

Then, as before,  $\frac{\partial}{\partial \alpha} \log |A| = \sum_i \frac{1}{\lambda_i + \alpha}$

$$\frac{\partial}{\partial \alpha} \log p(\mathcal{D} | \alpha, \beta) = \frac{W}{2\alpha} - \frac{1}{2} \vec{w}_{\text{MAP}}^T \vec{w}_{\text{MAP}} -$$

$$- \frac{1}{2} \sum_i \frac{1}{\lambda_i + \alpha} = 0, \text{ or}$$

$$\alpha \vec{w}_{\text{MAP}}^T \vec{w}_{\text{MAP}} = W - \alpha \sum_i \frac{1}{\lambda_i + \alpha} = \\ = \sum_i \frac{\lambda_i}{\lambda_i + \alpha} = \gamma$$

$$\text{So, } \hat{\alpha} = \frac{\gamma}{\vec{w}_{\text{MAP}}^T \vec{w}_{\text{MAP}}}$$

Ignored that  $\begin{cases} H = H(\alpha) \Rightarrow \lambda_i = \lambda_i(\alpha) \\ \vec{w}_{\text{MAP}} = \vec{w}_{\text{MAP}}(\alpha) \end{cases}$

Similarly,

$$\frac{\partial}{\partial \beta} \log P(\mathcal{D} | \alpha, \beta) = \frac{N}{2\beta} - \frac{1}{2} \sum_n (y(\vec{x}_n, \vec{w}_{MAP}) - t_n)^2 -$$

Note that

$$\lambda_i \sim \beta \Rightarrow \frac{d\lambda_i}{d\beta} = \frac{\lambda_i}{\beta}$$

$$- \frac{1}{2} \frac{\partial}{\partial \beta} = 0, \text{ or}$$

$$\text{Then } \frac{d}{d\beta} \log |A| = \sum_i \frac{d\lambda_i/d\beta}{\lambda_i + \alpha} = \frac{1}{\beta} \underbrace{\sum_i \frac{\lambda_i}{\lambda_i + \alpha}}_{\gamma}$$

$$\frac{N - \gamma}{\beta} = \sum_n (y - t_n)^2,$$

$$\hat{\beta}^{-1} = \frac{1}{N - \gamma} \sum_n (y - t_n)^2$$

[ alternate between estimating  $\vec{w}_{MAP}$  &  $\hat{\alpha}, \hat{\beta}$ . ]

We can use  $p(\mathcal{D} | \hat{\alpha}, \hat{\beta})$  to compare e.g. two-layer networks with different  $M$  (# hidden units). Recall the  $M! 2^M$  symmetry factor  $\Rightarrow$  compare  $M! 2^M p(\mathcal{D} | \hat{\alpha}, \hat{\beta})$  when  $M$  is changed.

## Classification

Consider a single sigmoid output in the  $K=2$  classification problem:

$$\log \mathcal{L} = \log P(\mathcal{D} | \vec{w}) = \sum_n [t_n \log y_n + (1-t_n) \log(1-y_n)],$$

where  $t \in (0,1)$  &  $y_n = y(\vec{x}_n, \vec{w})$

Then we need to maximize the posterior:

$$\log P(\vec{w} | \mathcal{D}, \alpha) = -\frac{\alpha}{2} \vec{w}^T \vec{w} + \log P(\mathcal{D} | \vec{w}) + \text{const}(\vec{w})$$

↓  
for a given  $\alpha$ , find  $\vec{w}_{\text{MAP}}$  by error backpropagation + e.g. conjugate grad.

Next, find  $H$  at  $\vec{w}_{\text{MAP}}$ , use

$$A = \alpha \mathbf{I} + H \quad \text{as before to get}$$

$$P(\vec{w} | \mathcal{D}, \alpha) \Rightarrow q(\vec{w} | \mathcal{D}) = \mathcal{N}(\vec{w} | \vec{w}_{\text{MAP}}, A^{-1})$$

gaussian approx'n to posterior

as before, under Laplace approx'n:

$$P(\mathcal{D} | \alpha) = \int d\vec{w} p(\mathcal{D} | \vec{w}) p(\vec{w} | \alpha)$$

is given by

$$\log P(\mathcal{D} | \alpha) \approx -\frac{\alpha}{2} \vec{w}_{\text{MAP}}^T \vec{w}_{\text{MAP}} + \sum_n [t_n \log y_n + (1-t_n) \log(1-y_n)] - \frac{W}{2} \log(2\pi\alpha^{-1}) \oplus$$

↑  
here,  $y_n = y(\vec{x}_n, \vec{w}_{\text{MAP}})$

$$\oplus \frac{W}{2} \log(2\pi) - \frac{1}{2} \log |A|, \text{ s.t.}$$

$$\frac{\partial}{\partial \lambda} \log P(\mathcal{D}|\lambda) \approx -\frac{1}{2} \vec{w}_{\text{MAP}}^T \vec{w}_{\text{MAP}} + \frac{W}{2\lambda} -$$

$$-\frac{1}{2} \sum_i \frac{1}{\lambda_i + \lambda} = 0 \quad \text{gives}$$

$$\hat{\lambda} = \frac{\sigma}{\vec{w}_{\text{MAP}}^T \vec{w}_{\text{MAP}}} \quad \text{as before}$$

We need to alternate  $\vec{w}_{\text{MAP}}$  &  $\hat{\lambda}$  estimates.

Finally, we need the predictive distr'n:

$$p(t=1|\vec{x}, \mathcal{D}) = \int d\vec{w} \underbrace{p(t=1|\vec{x}, \vec{w})}_{\sigma(a(\vec{x}, \vec{w}))} \underbrace{q(\vec{w}|\mathcal{D})}_{\text{Gaussian}}$$

Consider

$$a(\vec{x}, \vec{w}) = a(\vec{x}, \vec{w}_{\text{MAP}}) + \vec{b}^T \cdot (\vec{w} - \vec{w}_{\text{MAP}}),$$

output unit activation  $\vec{b} = \nabla a(\vec{x}, \vec{w}) \Big|_{\vec{w}_{\text{MAP}}}$

(cannot expand  $y(\vec{x}, \vec{w})$  as in regression due to an extra sigmoid)

Use results of (4.5.2): linear in  $\vec{w}$

$$p(t=1|\vec{x}, \mathcal{D}) \approx \int d\vec{w} \sigma(a(\vec{x}, \vec{w}_{\text{MAP}}) + \vec{b}^T \cdot (\vec{w} - \vec{w}_{\text{MAP}})) \times q(\vec{w}|\mathcal{D}) \quad \ominus$$



$$\textcircled{=} \int d\tilde{a} \delta(\tilde{a}) \left[ \int d\vec{w} \overbrace{q_g(\vec{w} | \mathcal{D})}^{\text{gaussian}} \delta(\tilde{a} - a(\vec{x}, \vec{w}_{\text{MAP}}) - \underbrace{\vec{b}^T(\vec{x}) \cdot (\vec{w} - \vec{w}_{\text{MAP}})}_{\text{}}) \right]$$

" $p(\tilde{a} | \vec{x}, \mathcal{D})$ ", can be computed by expanding the  $\delta$ -f'n in plane waves:

$$p(\tilde{a} | \vec{x}, \mathcal{D}) = \mathcal{N}(\tilde{a} | a(\vec{x}, \vec{w}_{\text{MAP}}), \underbrace{\vec{b}^T A^{-1} \vec{b}}_{\text{"}} \underbrace{\sigma_a^2(\vec{x})}_{\sigma^2})$$

Finally,

$$p(t=1 | \vec{x}, \mathcal{D}) = \int d\tilde{a} \delta(\tilde{a}) \mathcal{N}(\tilde{a} | a(\vec{x}, \vec{w}_{\text{MAP}}), \underbrace{\sigma_a^2(\vec{x})}_{\sigma^2}) \quad \textcircled{\approx}$$

approximate  $\delta(\tilde{a})$  with  $\phi(\tilde{a})$  and use (4.153) to get:

$$\textcircled{\approx} \delta(K(\sigma_a^2) a(\vec{x}, \vec{w}_{\text{MAP}})), \text{ where } \left. \begin{array}{l} \text{double-} \\ \text{-check:} \\ (5.190) \\ \text{wrong?} \end{array} \right\}$$

$$\rightarrow K(\sigma_a^2) = \frac{1}{\sqrt{1 + \frac{\pi \sigma_a^2}{8}}}$$

(4.154)