

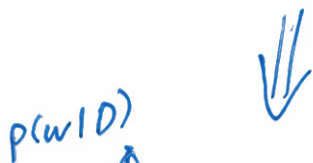
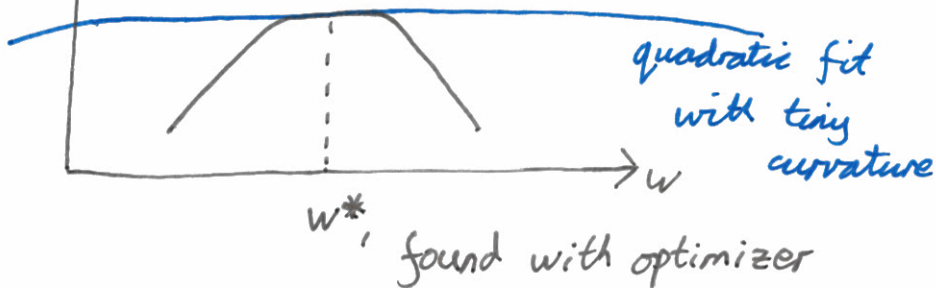
Laplace Approx

$$p(\underline{w} | D) \approx N(\underline{w}; \underline{w}^*, H^{-1})$$

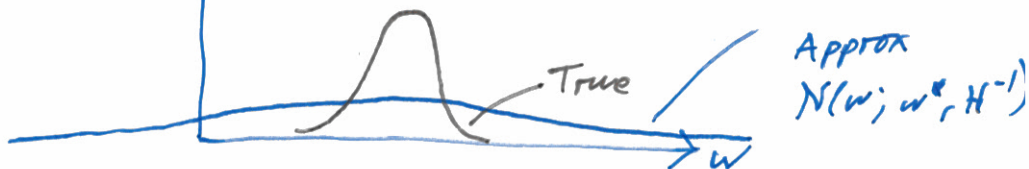
$$\log p(\underline{w}, D) = \log p(D | \underline{w}) p(\underline{w})$$

Hessian

$$H_{ij} = \frac{\partial^2 -\log p(\underline{w}, D)}{\partial w_i \partial w_j}$$



$$P(D) \approx \frac{p(w^*, D)}{N(w; w^*, H^{-1})}$$



Approx: Areas under both curves are 1.

$P(D)$? A) Too big; B) Too small; c) \approx Right; z) ?

Computing Predictions

$$\begin{aligned} p(y=1 | \underline{x}, D) &\approx \int p(y=1 | \underline{x}, \underline{w}) N(\underline{w}; \underline{w}^*, H^{-1}) d\underline{w} \\ &= \int \sigma(\underline{w}^T \underline{x}) N(\underline{w}; \underline{w}^*, H^{-1}) d\underline{w} \\ &= \mathbb{E}_{N(\underline{w}; \underline{w}^*, H^{-1})} [\sigma(\underline{w}^T \underline{x})] \end{aligned}$$

Average under an activation

$$a = \underline{w}^T \underline{x}$$

$$= \mathbb{E}_{N(a; \underline{w}^* \underline{x}, \underline{x}^T H^{-1} \underline{x})} [\sigma(a)]$$

mean a : $\underline{w}^* \underline{x}$

variance a : $\underline{x}^T H^{-1} \underline{x}$] TODO check.

$$= \int \sigma(a) N(a; \underline{w}^* \underline{x}, \underbrace{\underline{x}^T H^{-1} \underline{x}}_{\text{scalar}}) da$$

Could solve numerically.

Murphy § 8.4.4.2:

$$P(y=1 | \underline{x}, D) \approx \sigma(K \underline{w}^T \underline{x})$$

↑

$$K = \frac{1}{\sqrt{1 + \frac{\pi}{8} \underline{x}^T H^{-1} \underline{x}}}$$

Variational Methods

Another way to fit approx. to posterior

$$P(\underline{w} | D) \approx q(\underline{w}; \alpha)$$

For us $q(\underline{w}; \alpha) = N(\underline{w} | \underline{m}, V)$

Variational Params: $\alpha = \{\underline{m}, V\}$

Set up optimization problem.

⇒ Need cost f^n , measures discrepancy between $P(\underline{w} | D)$ and q .

A common way to compare dists:

Kullback-Leibler Divergence

→ KL Divergence

$$D_{KL}(p \parallel q) = \int p(z) \log \frac{p(z)}{q(z)} dz$$

≥ 0 (Can show Gibbs' inequality.)

Isn't a distance:

- Not symmetric $D_{KL}(p \parallel q) \neq D_{KL}(q \parallel p)$

Logistic Regression eq.

Minimize

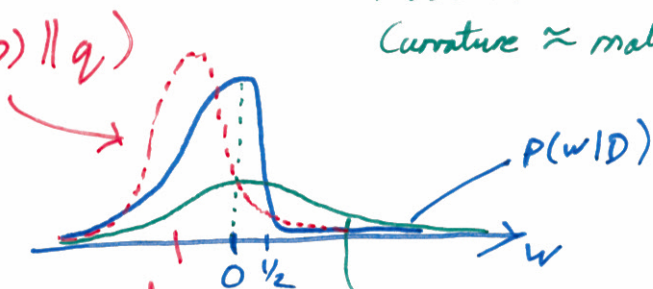
$$D_{KL}(p(w|D) \parallel q)$$

Mode at $\approx w=0$

Curvature \approx matches prior



Matches mean and variance of posterior

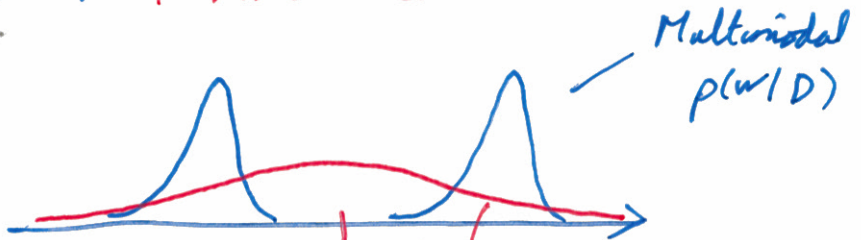


Laplace approx
 $p(w|D) \approx p(w)$

We don't minimize $D_{KL}(p||q)$

- 1) We don't know how.
- 2) [↑] Not a good idea
often

Example: match mean and variance



I would predict $D_{KL}(p||q)$
focussing on implausible parameters.

Minimizing $D_{KL}(q||p)$

$$D_{KL}(q(\underline{w}; \alpha) || p(\underline{w}|D))$$

$$= \int q(\underline{w}; \alpha) \log \frac{q(\underline{w}; \alpha)}{p(\underline{w}|D)} d\underline{w}$$

$$= \underbrace{- \int q(\underline{w}; \alpha) \log p(\underline{w}|D) d\underline{w}}_{\text{It's good if } q(\underline{w}; \alpha) \text{ is big when } p(\underline{w}|D) \text{ is.}} + \underbrace{\int q(\underline{w}; \alpha) \log q(\underline{w}; \alpha) d\underline{w}}_{\substack{-H[q(\underline{w}; \alpha)] \\ \uparrow \\ \text{Entropy}}}}$$

Really bad if $q(\underline{w}; \alpha)$ is big when $p(\underline{w}|D)$ is small.

Encourages q to be spread out

$$D_{KL}(q(w; \alpha) \| p(w|D))$$

↑
From Bayes' Rule

$$= \underbrace{\mathbb{E}_q[\log q]}_{\text{Can evaluate, } J} - \mathbb{E}_q[\log p(D/w)] - \mathbb{E}_q[\log p(w)]$$

~~$\mathbb{E}_q[\log p(D)]$~~
Don't know
but it's a
const.

$$D_{KL} \geq 0$$

$$J + \log p(D) \geq 0$$

$$\log p(D) \geq -J$$

$$\mathbb{E}_q[f(w)] = \int f(w) q(w; \alpha) dw$$