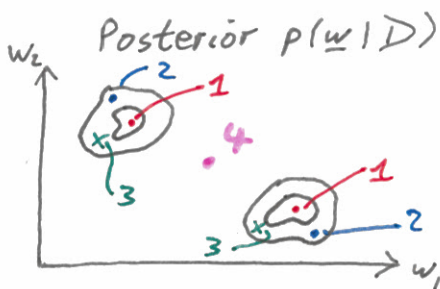
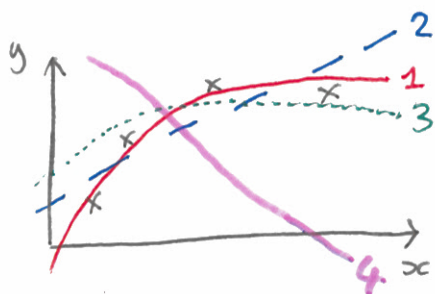


Variational Methods



Approximate posterior with $q(\underline{w}; \alpha) = N(\underline{w}; \underline{m}, \underline{V})$
 e.g. α

For prediction, fitting one mode might be ok.

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

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

$$= \underbrace{\mathbb{E}_q[\log q(\underline{w}; \alpha)]}_{- \text{Entropy}[q]} - \underbrace{\mathbb{E}_q[\log p(\underline{w} | D)]}_{\text{Could min. this term with } q = N(\underline{w}; \underline{w}^*, \underline{0})}$$

\Rightarrow Spread out distribution.

Substitute in

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

$$D_{KL} = \underbrace{\mathbb{E}_q[\log q] - \mathbb{E}_q[\log p(D|\underline{w})] - \mathbb{E}_q[\log p(\underline{w})]}_{\text{J, "can evaluate"}} + \cancel{\mathbb{E}_q[\log p(D)]}$$

Don't know,
log Marginal Likelihood

Minimize D_{KL} , by minimizing J

Gibbs' inequality $D_{KL} \geq 0$

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

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

\Rightarrow Lower bound on marginal likelihood.

To fit model choices or hyperparameters

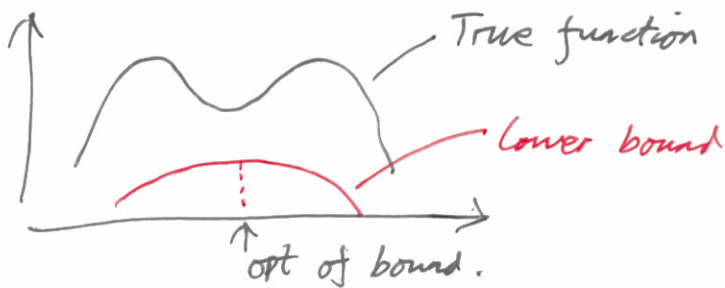
Jointly minimize J wrt $\{\underline{m}, V\}$

and wrt model hyperparameters σ_w^2

Might work...



Bad case:



Optimizing J

Gradient-based optimization.

Particularly stochastic gradient descent SGD

On $\alpha = \{\underline{m}, V\}$ and hypers... eg σ_w^2

Unconstrained optimization (Trick #1)

If we optimized σ_w^2 with SGD
we might make it -ve

Optimize $\log \sigma_w$ instead

V has to be positive definite, symmetric

$$V = LL^T, \quad L \text{ lower triangular} \\ \text{Diagonal is +ve.}$$

We create another matrix

$$\tilde{L}_{ij} = \begin{cases} L_{ij} & i \neq j \\ \log L_{ii} & i = j \end{cases}$$

Optimize $\tilde{L} \xrightarrow{\text{exp diagonal}} L \rightarrow V = LL^T \rightarrow \text{est cost}$
SGD \leftarrow backprop.

Evaluating the terms

"Entropy Terms" - we can compute...

For any $\underline{m}, V, \sigma_w^2 \dots$

Likelihood Term

$$\begin{aligned} & \mathbb{E}_q[\log p(D|\underline{w})] \\ &= \mathbb{E}_q\left[\underbrace{\sum_{n=1}^N \log p(y^{(n)} | \underline{x}^{(n)}, \underline{w})}_{\text{arrow}}\right] \end{aligned}$$

At least for logistic regression
can solve numerically.

Stochastic estimate ("Reparameterization trick")

$$\begin{aligned} & \mathbb{E}_{N(\underline{w}; \underline{m}, V)} [f(\underline{w})] \\ &= \mathbb{E}_{N(\underline{z}; 0, \mathbf{I})} [f(\underline{m} + L\underline{z})] \end{aligned}$$

Sample \underline{w} , by $\underline{z} \sim N(0, \mathbf{I})$

$$\underline{w} = \underline{m} + L\underline{z}$$

Monte Carlo estimate

$$\approx \frac{1}{S} \sum_{s=1}^S f(\underline{m} + L \underline{v}^{(s)}),$$

laziest
Simplest approx $S=1$

$$\underline{v}^{(s)} \sim N(0, \mathbf{I})$$

$$\approx f(\underline{m} + L \underline{v}), \quad \underline{v} \sim N(0, \mathbf{I})$$

Unbiased estimate.

$$\nabla_{\underline{m}} \mathbb{E}_{N(\underline{v}; 0, \mathbf{I})} [f(\underline{m} + L \underline{v})]$$

$$\approx \nabla_{\underline{m}} f(\underline{m} + L \underline{v}), \quad \underline{v} \sim N(\underline{v}; 0, \mathbf{I})$$

Unbiased

$$\nabla_L \mathbb{E}_{N(\underline{v}; 0, \mathbf{I})} [f(\underline{m} + L \underline{v})]$$

$$\dots \nabla_{\underline{w}} f(\underline{w}) \Big|_{\substack{\underline{v}^T \\ \underline{w} = \underline{m} + L \underline{v}}}$$

As long as you can

differentiate $f(\underline{w}) = \log P(y^{(n)} | \underline{x}^{(n)}, \underline{w})$

... apply chain rule.