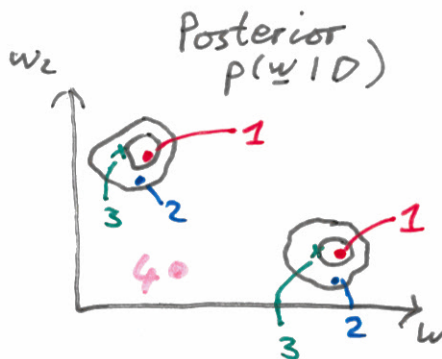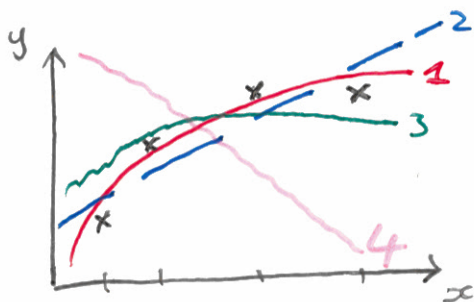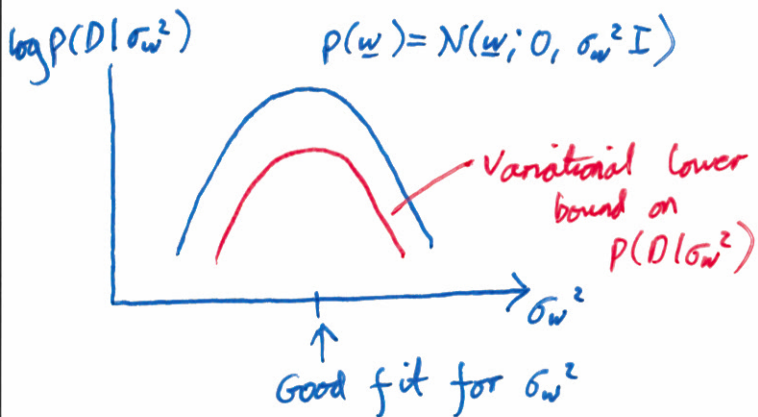# Variational Methods



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

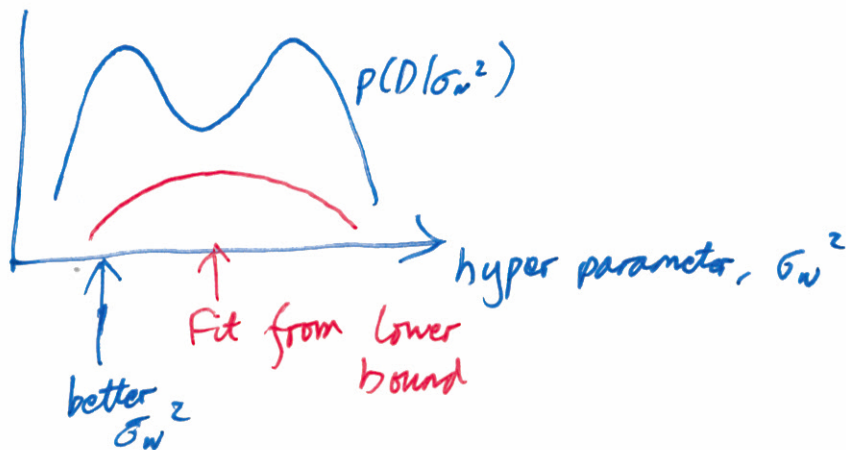$\underbrace{\qquad}_{\alpha}$

- For prediction 1 mode
                    might be ok.

- We use $q$ to approx. $p(D|M)$
    The number itself will be wrong by a
                                large factor.

$\log p(D|\sigma_w^2)$          $p(\underline{w}) = N(\underline{w}; 0, \sigma_w^2 I)$



- Variational lower
    bound on
        $p(D|\sigma_w^2)$

Good fit for $\sigma_w^2$

L25

Bad case :



better $\sigma_w^2$

Fit from lower bound

hyper parameter, $\sigma_w^2$

$p(D|\sigma_w^2)$

$$D_{KL}\left( q(\underline{w};\alpha) \| p(\underline{w}|D) \right) = \mathbb{E}_q\left[ \log \frac{q(\underline{w};\alpha)}{p(\underline{w}|D)} \right]$$

$$= \int q(\underline{w};\alpha) \log \frac{q(\underline{w};\alpha)}{p(\underline{w}|D)} \, d\underline{w} \geqslant 0$$

$$= -\underbrace{\mathbb{E}_q\left[ \log p(\underline{w}|D) \right]}_{\text{Minimize this term with } q(\underline{w};\alpha)=N(\underline{w};\underline{w}^*, \underline{0})} + \overbrace{\mathbb{E}_q\left[ \log q(\underline{w};\alpha) \right]}^{-\text{Entropy of } q}$$

Substitute $\quad p(\underline{w}\,|\,D) = \dfrac{p(D\,|\,\underline{w})\,p(\underline{w})}{P(D)}$

$$D_{KL} = \underbrace{\mathbb{E}_q\left[\log q(\underline{w};\alpha)\right] - \mathbb{E}_q\left[\log p(\underline{w})\right]}_{}$$

$$\underbrace{-\mathbb{E}_q\left[\log p(D\,|\,w)\right]}_{} + \underbrace{\mathbb{E}_q\left[\log p(D)\right]}_{\log p(D)}$$

$$\hookrightarrow = J(\alpha) = J(\underline{m}, V)$$

Fit $\quad q(\underline{w}) = N(\underline{w};\, m, V) \quad$ by minimizing $J$

Also want to fit model hyperparameters.

We'd like to maximize $\log p(D\,|\,\text{hypers})$

$$\log p(D\,|\,\text{hypers}) \geqslant -J$$

$$(\text{Because } D_{KL} \geqslant 0)$$

$\rightarrow$ Jointly minimize $J$

wrt $\{\underline{m}, V\}$ and model hypers
Variational params $\quad$ eg $\sigma_w^2$

# Optimizing $D_{KL}(q(w) \| p(\underline{w} | D))$

Gradient-based optimization

    Particularly stochastic gradient descent.
                                (S.G.D.)

Not on $\underline{w}$, weights of logistic regression
                        or a n.n.

On beliefs about $\underline{w}$, $q(\underline{w}) = N(\underline{w}; \underline{m}, V)$
                                             Opt. these.

Also optimize hyper-parameters.

## Unconstrained optimization (Trick #1)

If we optimize $\sigma_w^2$ S.G.D.
                    might make it -ve.

Also $V$ has to be positive definite

Instead we optimize $\log \sigma_w$

Also transform $V$:            ← — $L$ is Cholesky decomp of $V$.

$$V = LL^T, \quad L \text{ lower triangular}$$
matrix with +ve diagonal

We create another matrix

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

Optimize $\hat{L} \longrightarrow L \longrightarrow V = LL^T \rightarrow$ est.
     exp.                                                          cost
     diagonal                                                  backprop.
                                          $\longleftarrow$
                                                               gradients

"Entropy Terms"

We can evaluate $\mathbb{E}_{N(\underline{w};\underline{m},V)}\left[\log N(\underline{w};M,\Sigma)\right]$

For any $\underline{m}, V, M, \Sigma \dots$

Likelihood Term

$$\mathbb{E}_q\left[\log p(D|\underline{w})\right]$$
$$= \mathbb{E}_q\left[\sum_{n=1}^{N} \log p(y^{(n)}|\underline{x}^{(n)}, \underline{w})\right]$$

Could do by numerical integration.

## Stochastic estimation

$$\mathbb{E}_{N(\underline{w};\underline{m},v)}\left[f(\underline{w})\right]$$

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

$$\left[\begin{array}{l} \text{Sample } \underline{w}, \quad \text{by } \underline{v}\sim N(0,I) \\ \qquad \underline{w}=\underline{m}+L\underline{v} \end{array}\right]$$

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

Monte Carlo estimate. Unbiased est.

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

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

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

$$\dots \text{chain rule only } \nabla_{\underline{w}}\,f$$