

Hyperparameter tuning using Evidence maximization

~~P(data)~~
Finding evidence directly would be complex. It will involve solving of the integral which is not trivial i.e.

$$\int_{\mathbf{w}} \text{Prior}(\mathbf{w}) \underbrace{P(\text{data}|\mathbf{w})}_{\text{Likelihood}} d\mathbf{w}$$

This likelihood ~~has a standard form~~ contains terms involving $\sigma(\mathbf{w}^T \mathbf{x})$ ^{because of} which will be very hard to compute.

→ Evidence approximation

We can approximate the posterior distribution with normal dist to get approximate posterior. This also gives us a way to calculate the evidence

$$\text{Evidence} \approx f(\mathbf{w}_{\text{MAP}}) (2\pi)^{d/2} |\Sigma|^{-1/2}$$

where $f(\mathbf{w}_{\text{MAP}})$ is posterior at \mathbf{w}_{MAP}

This equation comes from solving ~~approximate~~ approximate posterior which is distributed as ~~N(w|w_MAP, Sigma^-1)~~ $N(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \Sigma^{-1})$

$$\text{where } \Sigma = (\mathbf{Q}^T \mathbf{R} \mathbf{Q} + \alpha \mathbf{I})^{-1}$$

$$\Rightarrow \text{Therefore evidence} \approx f(\mathbf{w}_{\text{MAP}}) (2\pi)^{d/2} |(\mathbf{Q}^T \mathbf{R} \mathbf{Q} + \alpha \mathbf{I})^{-1}|^{1/2}$$

$$\log \text{Ev} = \log f(\mathbf{w}_{\text{MAP}}) + \frac{d}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma|$$

lets look at $f(\mathbf{w}_{\text{MAP}}) \rightarrow$

$$f(w_{MAP}) \propto \text{prior} | w_{MAP} * \text{likelihood} | w_{MAP}$$

$$\log(f(w_{MAP})) = C + \log(\text{prior})_{w_{MAP}} + \log \text{likelihood} | w_{MAP}$$

Now ~~log likelihood~~ ~~likelihood~~ ~~is~~

$$\text{likelihood for logistic regression is } \prod_i y_i^{t_i} (1 - y_i)^{1 - t_i}$$

$$\text{where } y_i = \sigma(w^T \phi(x_i))$$

→ this likelihood function doesn't contain α .

Similarly likelihood for poisson regression → does not contain α .
and ordinal regression

$$\rightarrow \boxed{\frac{\partial \log \text{likelihood}}{\partial \alpha} = 0}$$

Now using this function, let's maximise log evidence.

$$\log E = \log f(w_{MAP}) + \frac{d}{2} \log(2\pi) + \frac{1}{2} \log(|\Sigma|)$$

$$\frac{\partial \log E}{\partial \alpha} = \frac{\partial}{\partial \alpha} (\log f(w_{MAP})) + \frac{1}{2} \frac{\partial}{\partial \alpha} \log(|\Sigma|)$$

$$= \frac{\partial}{\partial \alpha} (C + \log(\text{prior})_{w_{MAP}} + \log \text{likelihood}) + \frac{1}{2} \frac{\partial}{\partial \alpha} (\log |\Sigma|)$$

$$= \frac{\partial}{\partial \alpha} (C + \frac{d}{2} \log(1/2\pi) + \frac{d}{2} \log(\alpha) - \frac{1}{2} \alpha \|w_{MAP}\|^2 + \log \text{likelihood}) + \frac{1}{2} \frac{\partial}{\partial \alpha} \log(|\Sigma|)$$

$$= \frac{d}{2\alpha} - \frac{1}{2} \|w_{MAP}\|^2 + \frac{1}{2} \frac{\partial}{\partial \alpha} (\log |\Sigma|)$$

↓

Now only this

$$\frac{\partial}{\partial \alpha} \log(|\Sigma|)$$

Now $|\Sigma| = \text{product of eigenvalues}$

\cdot product of eigenvalues of $(Q^T R Q + \alpha I)^{-1}$

If d_i are eigenvalues of $Q^T R Q$

$\Rightarrow d_i + \alpha$ are eigenvalues of $(Q^T R Q + \alpha I)$

$\Rightarrow \frac{1}{d_i + \alpha}$ are eigenvalues of $(Q^T R Q + \alpha I)^{-1}$

$$\rightarrow \frac{\partial}{\partial \alpha} \left[\log \left(\pi \frac{1}{d_i + \alpha} \right) \right]$$

$$= \frac{\partial}{\partial \alpha} \left(\sum_i \log \frac{1}{d_i + \alpha} \right)$$

$$= \frac{\partial}{\partial \alpha} \sum_i - \log(d_i + \alpha)$$

$$= - \frac{\partial}{\partial \alpha} \left[\sum_i \log(d_i + \alpha) \right]$$

$$= \sum_i \frac{-1}{d_i + \alpha}$$

$$\Rightarrow \frac{\partial \log E}{\partial \alpha} = \frac{d}{2\alpha} - \frac{1}{2} \|W_{MAP}\|^2 - \frac{1}{2} \sum_i \frac{1}{d_i + \alpha} = 0$$

$$= \frac{d}{\alpha} - \|W_{MAP}\|^2 = \sum_i \frac{1}{d_i + \alpha}$$

$$\Rightarrow d - \frac{\sum \alpha}{d_i + \alpha} = \alpha \|W_{MAP}\|^2$$

$$\Rightarrow \sum_i \frac{d_i}{d_i + \alpha} = \alpha \|W_{MAP}\|^2 \Rightarrow \sum_i \frac{d_i}{d_i + \alpha} = \alpha \|W_{MAP}\|^2$$

$$\Rightarrow \alpha =$$

$$\frac{\sum y_i}{\sum 1} = \frac{\sum \frac{d_i}{d_i + \alpha}}{\|W_{MAP}\|^2}$$

$$\alpha = \frac{Y}{\|W_{MAP}\|^2}$$

Now α can be optimised using the iterative process.

- 1) Initialise α
- 2) find $\|W_{MAP}\|^2$ using that α .
- 3) update α
- 4) Repeat until convergence.

Note that this α will be used to calculate gradient and Hessian α .
 \rightarrow Hence $\|W_{MAP}\|$ depends on α .

For calculating optimal S and Q_j in ordinal \rightarrow

$$\lg E = \lg(f(W_{MAP})) + d/2 \lg(2\pi) + 1/2 \lg(1/\epsilon)$$

$$\frac{\partial \lg E}{\partial S} = \frac{\partial}{\partial S} \left[c + \frac{d}{2} \lg\left(\frac{1}{2}\pi\right) + \frac{d}{2} \lg(\alpha) - \frac{1}{2} \|W_{MAP}\|^2 + \lg \text{likelihood} \right]$$

$$+ \frac{1}{2} \frac{\partial}{\partial S} \lg(1/\epsilon)$$

$$= \frac{\partial}{\partial S} \left[\sum_i \sum_j t_{ij} \left[\ln(y_{ij} - y_{(j-1)}) \right] + \frac{1}{2} \frac{\partial}{\partial S} \lg(1/\epsilon) \right]$$

$$= \frac{\partial}{\partial S} \left[\sum_i \sum_j (q_j - a) \left[\sigma(s^*(q_j - a)) \right] \left[1 - \sigma(s^*(q_{j-1} - a)) \right] \right. \\ \left. + \frac{1}{2} \frac{\partial}{\partial S} \lg(1/\epsilon) \right]$$

~~Now this can be solved using gradient descent to find optimal~~

Now $\frac{\partial}{\partial S} \lg(|\Sigma|)$

$$\Sigma^{-1} = Q^T R Q + \alpha I$$

$$\rightarrow \frac{\partial}{\partial S} \lg[|\Sigma|] = \frac{\partial}{\partial S} \lg \left[\prod_i \frac{1}{d_i + \alpha} \right] \quad \text{where } d_i \text{ are eigenvalues of } Q^T R Q$$

$$= - \frac{\partial}{\partial S} \sum_i \lg(d_i + \alpha)$$

$$= - \sum_i \frac{1}{d_i + \alpha} \frac{\partial d_i}{\partial S}$$

recalculate S where $\frac{\partial \lg |\Sigma|}{\partial S} = 0$ using gradient descent.

and then use that S to calculate W_{MAP} and repeat this process.