

# Bayesian Model Selection

①  $D \quad M_1, M_2, \dots, M_k$

$\mathcal{M} \stackrel{\text{def}}{=} p(\theta), p(D|\theta)$

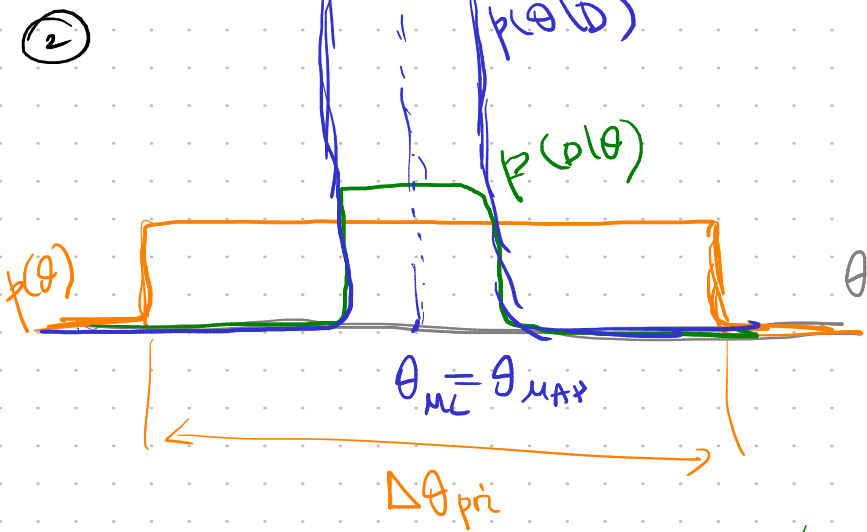
$p(\bar{\theta} | M_k), p(D|\bar{\theta}, M_k)$

$p(\underline{M}_k | D) = \frac{p(M_k) p(D|M_k)}{1/k \cdot p(D)}$

$\text{argmax}_k p(D|M_k)$

$p(\bar{\theta}_k | D, M_k) = \frac{p(\bar{\theta}_k | M_k) \cdot p(D|\bar{\theta}_k, M_k)}{p(D|M_k)}$

$p(D|M_k) = \int p(\bar{\theta}_k | M_k) p(D|\bar{\theta}_k, M_k) d\bar{\theta}_k$



$p(D) = \int p(\bar{\theta}) p(D|\bar{\theta}) d\bar{\theta} = \int_{(\Delta\theta_{pri})} \left( \frac{1}{\Delta\theta_{pri}} \right) p(D|\bar{\theta}) d\bar{\theta} =$

$= \int_{(\Delta\theta_{post})} \frac{1}{(\Delta\theta_{pri})} p(D|\bar{\theta}_{MAP}) d\bar{\theta} = p(D|\bar{\theta}_{MAP}) \cdot \left( \frac{\Delta\theta_{post}}{\Delta\theta_{pri}} \right)$

information criterion

$\log p(D) = \log p(D|\theta_{MAP}) - \log \frac{\Delta\theta_{pri}}{\Delta\theta_{post}}$

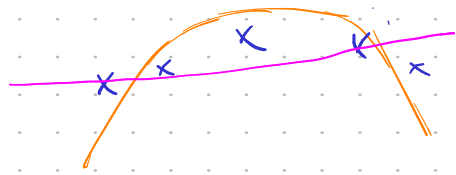


③ Sanity check

$D \sim \mathcal{M}_{true}$

$p(D|\mathcal{M}_{true}) \geq p(D|\mathcal{M})$

$E_{D \sim p(D|\mathcal{M}_{true})} \left[ \ln \frac{p(D|\mathcal{M}_{true})}{p(D|\mathcal{M})} \right] =$



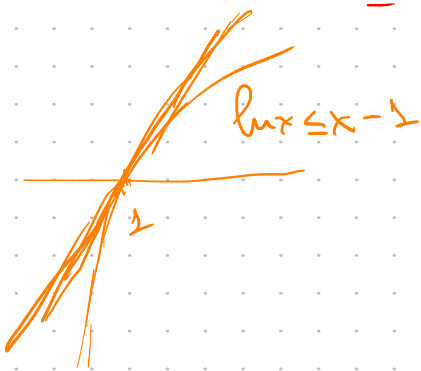
$= \int p(D|\mathcal{M}_{true}) \ln \frac{p(D|\mathcal{M}_{true})}{p(D|\mathcal{M})} dD = \text{KL}(p(D|\mathcal{M}_{true}) || p(D|\mathcal{M})) \geq 0$

$\text{KL}(p || q) = \int p(x) \ln \frac{p(x)}{q(x)} dx =$  Kullback-Leibler divergence

$= \int p(x) \ln p(x) dx - \int p(x) \ln q(x) dx$

$= H(p) - H(p, q)$

$\text{KL}(q || p) = \int q \ln \frac{q}{p} dx$



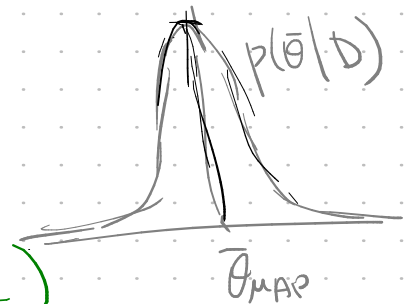
$= - \int p \ln \frac{q}{p} dx \geq - \int p \left( \frac{q}{p} - 1 \right) dx$

$= - \left( \int q dx - \int p dx \right) = 0$

④ BIC - Bayesian information criterion

$p(D) = \int p(\bar{\theta}) p(D|\bar{\theta}) d\bar{\theta}$

$\approx \mathcal{N}(\bar{\theta} | \bar{\theta}_{MAP}, \Sigma)$



$p(D) = \int \underbrace{f(\bar{\theta})}_{\text{const}} d\bar{\theta} \approx \int f(\bar{\theta}_{MAP}) \cdot e^{-\frac{1}{2}(\bar{\theta} - \bar{\theta}_{MAP})^T A (\bar{\theta} - \bar{\theta}_{MAP})} d\bar{\theta}$

$= \underline{f(\bar{\theta}_{MAP})} \cdot \sqrt{\frac{(2\pi)^d}{\det A}}$  , use  $A = -\nabla \nabla \ln f(\bar{\theta}) |_{\bar{\theta}_{MAP}}$

$\ln f(\bar{\theta}) \approx \ln f(\bar{\theta}_{MAP}) + (\bar{\theta} - \bar{\theta}_{MAP})^T \nabla \ln f + \frac{1}{2} (\bar{\theta} - \bar{\theta}_{MAP})^T \left( -\frac{\partial^2 \ln f}{\partial \theta_i \partial \theta_j} \right) (\bar{\theta} - \bar{\theta}_{MAP})$

$$\ln p(D) \approx \ln p(D|\bar{\theta}_{MAP}) + \ln p(\bar{\theta}_{MAP}) + \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln \det A$$

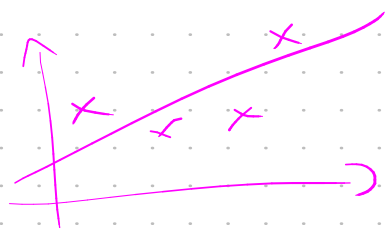
ATC - Akaike IC

Occam's factor  $\downarrow N = |D|$

BIC

$$\ln p(D) \approx \ln p(D|\bar{\theta}_{MAP}) - \frac{1}{2} \cdot d \cdot \ln N$$

5 Empirical Bayes



hyperparameters

$$p(y|\bar{x}, \bar{w}) = \mathcal{N}(y|\bar{w}^T \bar{x}, \sigma^2)$$

$$p(\bar{w}) = \mathcal{N}(\bar{w}|\bar{0}, \sigma_0^2 \cdot \mathbf{I})$$

$$\log p(\bar{w}|\bar{y}, X) = \text{const} - \frac{1}{2\sigma^2} \sum_n (y_n - \bar{x}_n^T \bar{w})^2 - \frac{\alpha}{2\sigma_0^2} \bar{w}^T \bar{w}$$

$$\log p(\bar{w}|\bar{y}, X) = \text{const} - \frac{\beta}{2} \sum_n (-)^2 - \frac{\alpha}{2} \bar{w}^T \bar{w}$$

MLE-II

$$p(\bar{w}|D, \alpha, \beta) = \frac{p(\bar{w}|\alpha) p(D|\bar{w}, \beta)}{p(D|\alpha, \beta)}$$

evidence function  
marginal likelihood  $\xrightarrow{\alpha, \beta} \max$

$$p(\bar{y}|X, \alpha, \beta) = \int p(\bar{w}|\alpha) p(\bar{y}|X, \bar{w}, \beta) d\bar{w}$$

$$= \int \left(\frac{\alpha}{2\pi}\right)^{d/2} e^{-\frac{\alpha}{2} \bar{w}^T \bar{w}} \cdot \left(\frac{\beta}{2\pi}\right)^{N/2} e^{-\frac{\beta}{2} \sum_n (y_n - \bar{w}^T \bar{x}_n)^2} d\bar{w}$$

$E(\bar{w})$

$$-\frac{\alpha}{2} \bar{w}^T \bar{w} - \frac{\beta}{2} (\bar{y} - X\bar{w})^T (\bar{y} - X\bar{w}) =$$

$$= -\frac{\alpha}{2} \bar{w}^T \bar{w} - \frac{\beta}{2} (\bar{y}^T \bar{y} - 2\bar{w}^T X^T \bar{y} + \bar{w}^T X^T X \bar{w}) =$$

$$= -\frac{1}{2} \bar{w}^T (\underbrace{\beta X^T X + \alpha \mathbf{I}}_A) \bar{w} + \underbrace{\beta \bar{w}^T X^T \bar{y}}_{A\bar{w}} - \frac{\beta}{2} \bar{y}^T \bar{y} =$$

$$= -\frac{1}{2}(\bar{w} - \bar{\mu})^T A (\bar{w} - \bar{\mu}) - \frac{\beta}{2} \bar{y}^T \bar{y} + \frac{1}{2} \bar{\mu}^T A \bar{\mu} \quad \text{use } A = \beta X^T X + \alpha I = Z_N^{-1}, \quad \bar{\mu} = \bar{\mu}_N = \beta \cdot A^{-1} X^T \bar{y}$$

$$E(\bar{\mu}) = -\frac{\alpha}{2} \bar{\mu}^T \bar{\mu} - \frac{\beta}{2} (\bar{y}^T \bar{y} - 2 \bar{\mu}^T X^T \bar{y} + \bar{\mu}^T X^T X \bar{\mu}) =$$

$$- \frac{\beta}{2} \bar{y}^T \bar{y} + \frac{\beta}{2} \bar{\mu}^T X^T X \bar{\mu} + \frac{\alpha}{2} \bar{\mu}^T \bar{\mu} =$$

$$= -\frac{\beta}{2} (\bar{\mu}^T X^T X \bar{\mu} + \bar{y}^T \bar{y} - 2 \bar{\mu}^T X^T \bar{y}) - \beta \bar{\mu}^T X^T \bar{y} + \frac{\alpha}{2} \bar{\mu}^T \bar{\mu} + \beta \bar{\mu}^T X^T X \bar{\mu}$$

$$= \bar{\mu}^T A \bar{\mu} + \frac{\alpha}{2} \bar{\mu}^T \bar{\mu} + \beta \bar{\mu}^T X^T X \bar{\mu}$$

$$= \bar{\mu}^T (\beta X^T X + \alpha I) \bar{\mu} + \frac{\alpha}{2} \bar{\mu}^T \bar{\mu} + \beta \bar{\mu}^T X^T X \bar{\mu}$$

$$= \frac{\alpha}{2} \bar{\mu}^T \bar{\mu}$$

$$E(\bar{w}) = E(\bar{\mu}) - \frac{1}{2}(\bar{w} - \bar{\mu})^T A (\bar{w} - \bar{\mu}) \quad \alpha, \beta \rightarrow \max$$

$$\log p(D | \alpha, \beta) = E(\bar{\mu}) + \frac{d}{2} \log \alpha + \frac{N}{2} \log \beta - \frac{N}{2} \log Z_N - \frac{1}{2} \log \det A$$

log det A

$$A = \alpha I + \beta (X^T X)$$

$\lambda_i$  - eigne values

$$\beta (X^T X)$$

$\lambda_i + \alpha$  - ch. A

$$\log \det A = \sum_{i=1}^d \log(\lambda_i + \alpha)$$

$$\frac{\partial \lambda_i}{\partial \beta} = \frac{\partial}{\partial \beta} (\beta (X^T X) + \alpha I) \cdot \bar{u}_i = \lambda_i \cdot \bar{u}_i$$

$$\lambda_i \propto \beta$$

$$\lambda_i = \beta \cdot \lambda_i$$

$$\frac{\partial \log p(D | \alpha, \beta)}{\partial \alpha} = \frac{d}{2\alpha} - \frac{1}{2} \bar{\mu}^T \bar{\mu} - \frac{1}{2} \sum_i \frac{1}{\lambda_i + \alpha} = 0$$

$$\alpha \cdot \bar{\mu}^T \bar{\mu} + \alpha \cdot \sum \frac{1}{\lambda_i + \alpha} = d \quad \bar{\mu} = \bar{\mu}(\alpha, \beta)$$

$$\alpha = \frac{1}{\bar{\mu}^T \bar{\mu}} \cdot \sum_{i=1}^d \left( 1 - \frac{\alpha}{\lambda_i + \alpha} \right) = \frac{1}{\bar{\mu}^T \bar{\mu}} \cdot \sum_{i=1}^d \frac{\lambda_i}{\lambda_i + \alpha}$$

$$\frac{\partial \log p(D|\alpha, \beta)}{\partial \beta} = \frac{N}{2\beta} - \frac{1}{2} \sum_n (y_n - \bar{x}_n^T \bar{\mu})^2 - \frac{1}{2} \sum_i \frac{\lambda_i(\beta)}{\lambda_i + \alpha} = 0$$

$$\beta = \frac{N - \sum_i \frac{\lambda_i}{\lambda_i + \alpha}}{\sum_n (y_n - \bar{\mu}^T \bar{x}_n)^2}$$