

# Lecture 9

## Bayesian information criterion (BIC)

Recall that under the Laplace approximation,

$$p(\vec{\theta} | \mathcal{D}, M) = \frac{p(\vec{\theta}, \mathcal{D} | M)}{p(\mathcal{D} | M)}, \text{ where}$$

"2 from before

$$p(\mathcal{D} | M) = \ell \log p(\vec{\theta} | \mathcal{D} | M) (2\pi)^{D/2} |\mathcal{H}|^{-1/2}.$$

$D = \dim \{\vec{\theta}\}$

Using  $\hat{\vec{\theta}} = \vec{\theta}_{MAP}$ , we obtain:

$$\begin{aligned} \log p(\mathcal{D} | M) &= \underbrace{\log p(\mathcal{D} | \vec{\theta}_{MAP}, M)}_{\substack{\text{model type} \\ \text{log-likelihood}}} + \underbrace{\log p(\vec{\theta}_{MAP} | M)}_{\substack{\text{penalty terms}}} - \frac{1}{2} \log |\mathcal{H}| + \underbrace{\text{const}}_{\substack{\text{indep. of } M, \\ \text{indep. of } N}} \end{aligned}$$

In the large- $N$  limit,  $\vec{\theta}_{MAP} \xrightarrow{\text{MLE}} \vec{\theta}^*$   
and the effects of the prior  
are negligible, such that

$$\log p(\mathcal{D} | M) = \log p(\mathcal{D} | \vec{\theta}^*, M) - \frac{1}{2} \log |\mathcal{H}|, \text{ where}$$

Occam factor

$$H_{ij} = \frac{\partial^2}{\partial \theta_i \partial \theta_j} \left[ -\log p(\mathcal{D} | \vec{\theta}, M) - \log p(\vec{\theta} | M) \right] \Big|_{\vec{\theta}^*}$$

discard the prior

$$\textcircled{2} - \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log p(\mathcal{D} | \bar{\theta}, M) \Big|_{\bar{\theta}^*}$$

$\sum_{k=1}^N \log p(D_k | \bar{\theta}, M)$ , additive over data points

sample size

We can rewrite

$$H = \sum_{k=1}^N H_k = N \langle H \rangle, \text{ then}$$

$$\log |H| = \log |N \langle H \rangle| = \log(N^D \langle H \rangle) \textcircled{3}$$

if  $H$  is full rank

$\textcircled{3} \Rightarrow \log N + \underbrace{\log \langle H \rangle}_{\text{discard, indep. of } N}$

$$\text{Thus, } \underbrace{\log p(\mathcal{D} | M)}_{\text{BIC score}} \approx \log p(\mathcal{D} | \bar{\theta}^*, M) - \underbrace{\frac{D_M}{2} \log N}_{\text{model complexity penalty}}$$

$D_M = \# \text{ parameters in model } M$

The issue of rank: suppose that ( $D_M = 3$ )

$$H_2 = \begin{pmatrix} H_{11} & H_{12} & 0 \\ H_{21} & H_{22} & 0 \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} N \langle H_{11} \rangle & N \langle H_{12} \rangle & 0 \\ N \langle H_{21} \rangle & N \langle H_{22} \rangle & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

due to likelihood only

and that  $-\log p(\bar{\theta} | M) = \lambda |\bar{\theta}|^2$ , then

$$H = H_2 + \lambda \mathbb{I} = \begin{pmatrix} N\langle H_{11} \rangle + \lambda & N\langle H_{12} \rangle & 0 \\ N\langle H_{21} \rangle & \text{neglect} & N\langle H_{22} \rangle + \lambda \\ 0 & 0 & \text{neglect } \lambda \end{pmatrix}.$$

Now,  $\det H = N^2 \langle H_{11} \rangle \langle H_{22} \rangle \lambda - N^2 \langle H_{12} \rangle \langle H_{21} \rangle \lambda = N^2 \lambda [\langle H_{11} \rangle \langle H_{22} \rangle - \langle H_{12} \rangle \langle H_{21} \rangle].$

Thus,

$$\log |H| = \underbrace{2 \log N}_{\text{not } 3!} + \underbrace{\text{const}(N)}_{\substack{\text{discard} \\ \text{as before}}}$$


---

### Information theory

Entropy of a discrete random variable  $X$ :

$$H(X) = - \sum_{k=1}^K p(X=k) \log_2 p(X=k) = - E[\log p(X)].$$

If  $\log_2$  is used, units of entropy are bits.

Max entropy: ~~p<sub>k</sub>~~  $p(X=k) = \frac{1}{K}, \forall k$ .

Then  $H(X) = \log_2 K$ .

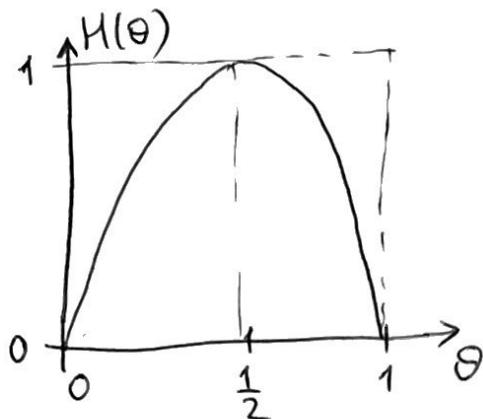
Min entropy:  $p(X=k') = 1$ , all others = 0.

Then  $H(X) = 0$ .

If  $K=2$  (binary data), we have:

$$\begin{cases} p(X=1) = \theta, \\ p(X=0) = 1-\theta. \end{cases}$$

Then  $H(X) = -\theta \log_2 \theta - (1-\theta) \log_2 (1-\theta) = H(\theta)$



$$\max\{H(\theta)\} = 1 \text{ bit}$$

Cross-entropy: between discrete distributions

$$P \& Q: H(P, Q) = - \sum_{k=1}^K p_k \log_2 q_k$$

Joint entropy:

$$H(X, Y) = - \sum_{x,y} p(x, y) \log_2 p(x, y)$$

If  $X \& Y$  are indep.,  $p(x, y) = p(x)p(y)$

and  $H(X, Y) = H(X) + H(Y)$  upper bound  
on entropy

If  $X \& Y$  are correlated, we expect

$$H(X, Y) < H(X) + H(Y).$$

Lower bound:  $H(X, Y) \geq \max\{H(X), H(Y)\} \geq 0.$  (\*)

## Conditional entropy:

$$\begin{aligned}
 H(Y|X) &= - \sum_x p(x) \sum_y p(y|x) \log p(y|x) = \\
 &= - \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)} = - \sum_{x,y} p(x,y) \log p(x,y) + \\
 &\quad + \sum_x p(x) \log p(x) = H(X,Y) - H(X).
 \end{aligned}$$

$$\text{Thus, } H(X,Y) = H(X) + H(Y|X).$$

$$\text{Likewise, } H(X,Y) = H(Y) + H(X|Y).$$

Since  $H(X|Y) \geq 0$ ,  $H(Y|X) \geq 0$ , we recover Eq. (\*).

moreover,

$$\underbrace{H(X,Y)}_{H(X)+H(Y|X)} \leq H(X) + H(Y) \Rightarrow H(Y|X) \leq H(Y)$$

- If  $X$  completely determines  $Y$ :  $H(Y|X) = 0$ .
- If  $X$  &  $Y$  are indep.:  $H(Y|X) = H(Y)$ .

"Conditioning on data does not increase uncertainty, on average".

→ chain rule for entropy:

$$\begin{aligned}
 H(X_1, X_2, \dots, X_n) &= \sum_{i=1}^n H(X_i | X_1, \dots, X_{i-1}) = \\
 &= H(X_1) + H(X_2 | X_1) + \dots + H(X_n | X_1, \dots, X_{n-1})
 \end{aligned}$$

## Continuous random variables:

Differential entropy:

$$h(x) = - \int dx p(x) \log p(x).$$

For example, if  $X \sim U(0, a)$ ,

$$h(x) = - \int_0^a dx \frac{1}{a} \log \frac{1}{a} = \log a,$$

can be  $< 0$  if  $a < 1$ .

Similarly, if  $X \sim N(\mu, \sigma^2)$ ,

$$\begin{aligned} h(x) &= - \int dx p(x) \left[ -(x-\mu)^2 / 2\sigma^2 - \frac{1}{2} \log (2\pi\sigma^2) \right] = \\ &= \frac{1}{2} + \frac{1}{2} \log (2\pi\sigma^2). \end{aligned}$$

## KL divergence:

$$D_{KL}(p||q_f) = \sum_{k=1}^K p_k \log \frac{p_k}{q_{fk}} \quad \text{or}$$

$$D_{KL}(p||q_f) = \int dx p(x) \log \frac{p(x)}{q_f(x)}.$$

Note that  $D_{KL}(p||q_f) = - \underbrace{\sum_k p_k \log q_{fk}}_{H(p, q_f), \text{ cross entropy}} + \underbrace{\sum_k p_k \log p_k}_{-H(p), \text{ entropy of } p}$

$$\text{Ex. } D_{KL}(p, q_\theta) = -\underbrace{H(p)}_{\log \frac{\log(2\pi\sigma_1^2) + 1}{2}} - \int dx p(x) \left[ -\frac{(x-\mu)^2}{2\sigma_2^2} - \frac{1}{2} \log(2\pi\sigma_2^2) \right] \quad (11)$$

$$\begin{cases} p = N(\mu, \sigma_1^2), \\ q_\theta = N(\mu, \sigma_2^2). \end{cases}$$

$$\textcircled{B} - \frac{1}{2} - \frac{1}{2} \log(2\pi\sigma_1^2) + \frac{\sigma_1^2}{2\sigma_2^2} + \frac{1}{2} \log(2\pi\sigma_2^2) =$$

$$= \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2}{2\sigma_2^2} - \frac{1}{2}.$$

$$\text{If } \sigma_2 \gg \sigma_1, \quad D_{KL} \approx \log \frac{\sigma_2}{\sigma_1} - \frac{1}{2} \geq 1.$$

$$\text{If } \sigma_2 \ll \sigma_1, \quad D_{KL} \approx \frac{\sigma_1^2}{2\sigma_2^2} \gg 1.$$

In fact,  $D_{KL} \geq 0$ , with  $D_{KL} = 0$  iff  $p = q$ .

Recall Jensen's inequality:

for any convex function  $f$ ,

$$f\left(\sum_{i=1}^n \lambda_i \bar{x}_i\right) \leq \sum_{i=1}^n \lambda_i f(\bar{x}_i)$$

$$0 \leq \lambda_i \leq 1, \quad \sum_{i=1}^n \lambda_i = 1.$$

" $f$  of the average is less than the average of  $f$ 's"

Now, consider  $-D_{KL}(p||q_\theta) = \sum_x p(x) \log \frac{q_\theta(x)}{p(x)} \leq \log \left( \sum_x p(x) \frac{q_\theta(x)}{p(x)} \right) = \log 1 = 0.$

Jensen;  $\log$  is concave

$$\boxed{D_{KL}(p||q_\theta) \geq 0}.$$

→ Uniform distribution maximizes the entropy:

$$u^{(x)} = \frac{1}{K}$$

↙ # outcomes

Then  $0 \leq D_{KL}(p || u) = \sum_x p(x) \log \frac{p(x)}{u(x)} =$

$$= \log K - H(p), \text{ or}$$

$$\underline{H(p) \leq \log K}.$$