

# Notes on Learning from Data, MLE and MAP

## ① Chain Rule

$$P(x_1, \dots, x_D) = P(x_1) \prod_{i=2}^D P(x_i | x_1, \dots, x_{i-1})$$

$$P(x, y, z) = P(x | y, z) P(y | z) P(z) \quad (1)$$

$$= P(x | y, z) P(z | y) P(y) \quad (2)$$

$$= P(y | x, z) P(x | z) P(z) \quad (3)$$

$$= P(y | x, z) P(z | x) P(x) \quad (4)$$

$$= P(z | x, y) P(x | y) P(y) \quad (5)$$

$$= P(z | x, y) P(y | x) P(x) \quad (6)$$

In general, suppose that we have  $N$  Random Variables  $X_1, \dots, X_N$ . The total number  $T$  of possible joint distribution  $P(X_1, \dots, X_N)$  factorizations using chain-rule can be written as:

$$T = N \cdot (N-1) (N-2) \dots 2 \cdot 1 \\ = N!$$

## ② Marginal Independence



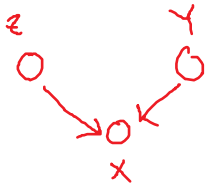
Choose the factorization of  $P(x, y, z)$  that fits the graphical model

$$\textcircled{6} P(x, y, z) = P(z | x, y) P(y | x) P(x) \\ = P(z | y, x) P(y) P(x)$$

$$\textcircled{5} P(x, y, z) = P(z | x, y) P(x | y) P(y)$$

$$= p(z|x, y) p(x) p(y)$$

### ③ Conditional Independence



To prove  $Z \perp Y | X$ ,  
Just show that either:

①  $p(z|x, y) = p(z|x)$   
 $p(y|x, z) = p(y|x)$

②  $p(y, z|x) = p(y|x) p(z|x)$

Ex. show that  $Z \not\perp Y | X$

$$p(z|x, y) = \frac{p(x, y, z)}{p(x, y)} = \frac{p(x|y, z) \cdot \cancel{p(y)} \cdot p(z)}{p(x|y) \cancel{p(y)}} \neq p(z|x)$$

Further Reading:

1. Chapter 10 Probabilistic Machine Learning, Murphy 2012
2. Chapter 20 Information Theory, Inference, and Learning Algorithm, Mackay 2003

### ④ Mean and Variance

Mean  $\Rightarrow \mathbb{E}[X]_{x \sim p}$ . Suppose that  $x \sim \text{Bernoulli}(\theta)$ , that is  $p(x) = \theta^x (1-\theta)^{(1-x)}$  with  $x \in \{0, 1\}$ . Then:

$$\mathbb{E}[X]_{x \sim p} = \sum_{x=0}^1 p(x) \cdot x = 0 \cdot p(x=0) + 1 \cdot p(x=1)$$

$$= \theta$$

$$\begin{aligned} \text{Variance} \Rightarrow \mathbb{E}[(X - \mathbb{E}[X])^2]_{x \sim p} &= \sum_{x=0}^1 (x - \theta)^2 \cdot p(x) \\ &= (0 - \theta)^2 \cdot p(x=0) + (1 - \theta)^2 \cdot p(x=1) \\ &= \theta^2 (1 - \theta) + (1 - \theta)^2 \cdot \theta \\ &= \theta (1 - \theta) (\theta + 1 - \theta) \\ &= \theta (1 - \theta) \end{aligned}$$

## ⑤ Entropy and KL-divergence

Entropy:  $H(X) = \mathbb{E}_{X \sim P} [I(X)] = \mathbb{E}_{X \sim P} [-\log P(X)]$

Ex:  $X \sim \text{Bernoulli}(\theta) \Rightarrow H(X) = -\sum_{x=0}^1 P(x) \cdot \log(P(x))$

$$= -[P(x=0) \log P(x=0) + P(x=1) \log P(x=1)]$$

$$= -[(1-\theta) \cdot \log(1-\theta) + \theta \cdot \log \theta]$$

$$= -[\log(1-\theta)^{(1-\theta)} + \log \theta^{\theta}]$$

$$= -[\log \cdot \theta^{\theta} \cdot (1-\theta)^{(1-\theta)}]$$

KL-divergence:  $KL[P \parallel Q] = \mathbb{E}_{X \sim P(X)} \left[ \log \frac{P(X)}{Q(X)} \right]$

Properties of KL-divergence:

①  $KL[P \parallel Q] \geq 0$  (Gibb's inequality)

Proof:

'Jensen's inequality states that if  $f$  is a convex function and  $X$  is a random variable, then  $f$  satisfies:

$$\mathbb{E}_{X \sim P} [f(X)] \geq f(\mathbb{E}_{X \sim P} [X])$$

Since  $\log$  is a concave function, then  $-\log$  is a convex function.

Setting  $f(x) = -\log \frac{Q(x)}{P(x)}$  gives us:

$$KL[P(X) \parallel Q(X)] = \mathbb{E}_{X \sim P} \left[ \log \frac{P(X)}{Q(X)} \right]$$

$$= \mathbb{E}_{X \sim P} \left[ -\log \frac{Q(X)}{P(X)} \right]$$

$$\geq -\log \mathbb{E}_{X \sim P} \left[ \frac{Q(X)}{P(X)} \right]$$

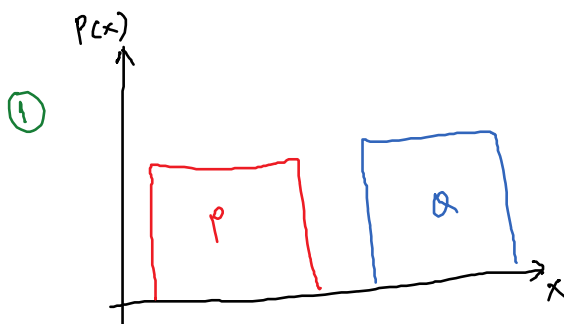
$$= -\log \int_{\mathcal{X}} P(x) \cdot \frac{Q(x)}{P(x)} dx$$

$$= -\log \int_{\mathcal{X}} Q(x) dx = -\log 1 = 0$$

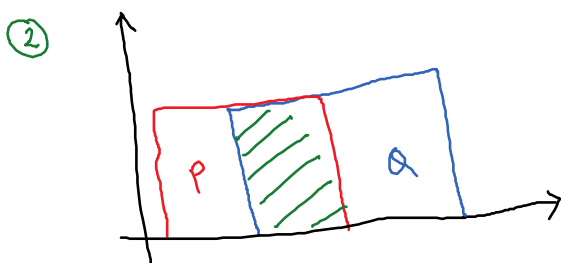
$P = Q$  iff  $KL[P \parallel Q] = 0$

②  $KL[P \parallel Q] \neq KL[Q \parallel P]$

$$\int p(x) \cdot \log \frac{p(x)}{q(x)} dx \neq - \int q(x) \log \frac{p(x)}{q(x)} dx$$

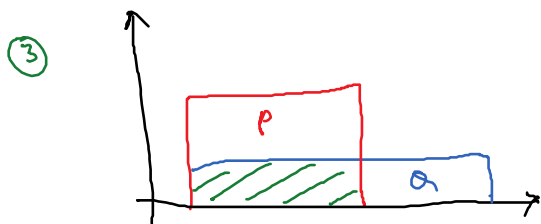


Support (P) and support (Q) is disjoint  
Therefore  $KL[P \parallel Q] = KL[Q \parallel P] = \infty$



Support (P) and support (Q) overlap, but neither is a subset of another. Therefore

$$KL[P \parallel Q] = KL[Q \parallel P] = \infty$$



support (Q)  $\subseteq$  support (P). Therefore

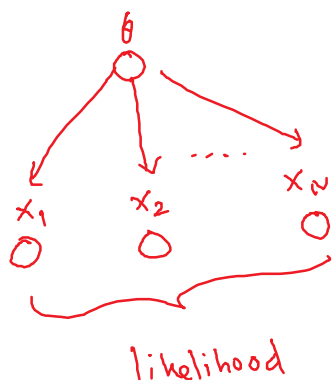
$$KL[P \parallel Q] \text{ is finite}$$

$$KL[Q \parallel P] = \infty$$

Further reading:

- Bishop 2006, chapter 1: Information theory

⑥ MLE



$$- x_i \perp x_j \mid \theta \quad \forall i, j: i \neq j \quad 1 \leq i, j \leq N$$

$$- \mathcal{D} = \{x_i\}_{i=1}^N$$

$$- \hat{\theta}_{MLE} = \max_{\theta} \log p(\mathcal{D} \mid \theta)$$

$$= \max_{\theta} \log \prod_{i=1}^N p(x_i \mid \theta)$$

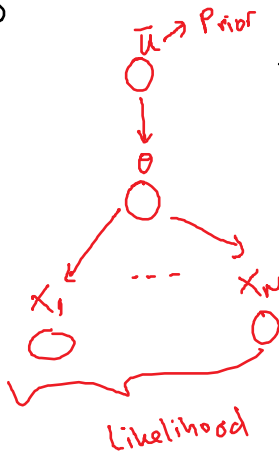
$$= \max_{\theta} \sum_{i=1}^N \log p(x_i \mid \theta)$$

likelihood

$$= \max_{\theta} \sum_{i=1}^N \log P(x_i | \theta)$$

$$\approx \frac{\partial \left[ \sum_{i=1}^N \log P(x_i | \theta) \right]}{\partial \theta} \bigg|_{\theta} = 0$$

⑦ MAP



The idea comes from Bayesian inference

① compute posterior

$$P(\theta | \mathcal{D}) = \frac{P(\mathcal{D} | \theta) \cdot P(\theta)}{P(\mathcal{D})} = \frac{P(\mathcal{D} | \theta) \cdot P(\theta)}{\int P(\mathcal{D} | \theta) \cdot P(\theta) d\theta}$$

② Given new data  $x^*$ , compute predictive posterior

$$P(x^* | \mathcal{D}) = \int P(x^* | \theta) \cdot P(\theta | \mathcal{D}) d\theta$$

Drawback of Bayesian  $\rightarrow$  intractability  $\rightarrow$  why?  $\rightarrow P(\mathcal{D})$

Semi-Bayes  $\rightarrow$  discard  $P(\mathcal{D}) \rightarrow$  MAP

$$P(\theta | \mathcal{D}) \propto P(\mathcal{D} | \theta) P(\theta)$$

$$\hat{\theta}_{MAP} = \max_{\theta} \log P(\mathcal{D} | \theta) P(\theta)$$

$$= \max_{\theta} \log P(\mathcal{D} | \theta) + \log P(\theta)$$

$$\lim_{n \rightarrow \infty} \hat{\theta}_{MAP} = \hat{\theta}_{MLE} \text{ (Bernstein-von Mises theorem)}$$

Simple observation; suppose that prior is uniform:  $P(\theta) = \text{constant}$ . Therefore:

$$\hat{\theta}_{MAP} = \max_{\theta} \log P(\mathcal{D} | \theta) + \log P(\theta)$$

$$= \max_{\theta} \log P(\mathcal{D} | \theta) + \text{constant}$$

$$\begin{aligned} &\approx \max_{\theta} \log p(\mathcal{D}|\theta) \\ &= \hat{\theta}_{MLE} \end{aligned}$$

Bonus: conjugacy

$$p(\theta) = \mathcal{U}[0, 1]$$

$$p(x|\theta) = \text{Bin}(x|n, \theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$$

$$\begin{aligned} p(\theta|x) &\propto p(x|\theta) \cdot p(\theta) \\ &\propto \theta^x (1-\theta)^{n-x} \\ &= \text{Beta}(x+1, n-x+1) \end{aligned}$$