

## 5.5 Maximum Likelihood Estimation.

### Motivation:

We want some principle from which we can derive specific functions that are good estimators for different models.

The most common: maximum likelihood principle

$X = \{x^{(1)}, \dots, x^{(m)}\}$ ,  $m$  examples drawn iid w.r.t. some unknown true data generating dist.  $P_{\text{data}}(x)$ .

$\{P_{\text{model}}(x; \theta)\}_{\theta}$ : a family of prob. dist. indexed by  $\theta$ .

Def

The maximum likelihood estimator for  $\theta$  is defined as

$$\begin{aligned} \theta_{ML} &= \arg \max_{\theta} P_{\text{model}}(X; \theta) \\ &= \arg \max_{\theta} \prod_{i=1}^m P_{\text{model}}(x^{(i)}; \theta) \quad (\text{due to iid}) \end{aligned}$$

$$\begin{aligned} &= \arg \max_{\theta} \sum_{i=1}^m \log(P_{\text{model}}(x^{(i)}; \theta)) \\ &= \arg \max_{\theta} \frac{1}{m} \sum_{i=1}^m \log(P_{\text{model}}(x^{(i)}; \theta)) \end{aligned}$$

$$= \arg \max_{\theta} \mathbb{E}_{x \sim \hat{P}_{\text{data}}} \log(P_{\text{model}}(x; \theta))$$

(where  $\hat{P}_{\text{data}}$  is the empirical distribution)

$$= \arg \min_{\theta} -\mathbb{E}_{x \sim \hat{P}_{\text{data}}} \log(P_{\text{model}}(x; \theta))$$

↑ Cross entropy b/w  $\hat{P}_{\text{data}}$  &  $P_{\text{model}}$

$$= \arg \min_{\theta} \mathbb{E}_{x \sim \hat{P}_{\text{data}}} [\log \hat{P}_{\text{data}}(x; \theta) - \log P_{\text{model}}(x; \theta)]$$

↑  $D_{KL}(\hat{P}_{\text{data}} \parallel P_{\text{model}})$ , the KL divergence

### Remark:

① In view of KL divergence (as measuring difference),

maximum likelihood  $\Leftrightarrow$  trying to make  $P_{\text{model}}$  match  $\hat{P}_{\text{data}}$ .

$\Leftrightarrow$  minimization of negative log-likelihood (NLL)

$\Leftrightarrow$  ..... cross entropy / KL divergence.

## 5.5.1 Conditional Log-Likelihood and Mean Square Error.

Def

$X$ : all our inputs,  $Y$ : all our observed targets. DL  
Supp.  
PI

The conditional maximum likelihood estimator is

$$\theta_{ML} = \arg \max_{\theta} P(Y|X; \theta).$$

Remark:

For iid drawn  $X, Y$ ,

$$\theta_{ML} = \arg \max_{\theta} \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}; \theta)$$

Prop

If the model

$P(y|x; \theta) = \mathcal{N}(y; \hat{y}(x; \omega), \sigma^2)$ , w/  $\sigma^2$  fixed and  $\omega$  as parameter, (so  $\theta = \omega$ ).

Gaussian  
mean variance

then the maximum likelihood estimator is the same as minimization of MSE on  $\{\hat{y}(x; \omega)\}_{\omega}$ .

<PP>

$$\begin{aligned} \theta_{ML} (= \omega_{ML}) &= \arg \max_{\omega} \sum_{i=1}^m \log \left( \frac{1}{\sqrt{2\pi}\sigma^2} \cdot e^{-\frac{(y^{(i)} - \hat{y}(x^{(i)}; \omega))^2}{2\sigma^2}} \right) \\ &= \arg \max_{\omega} \left( -m \log \sigma - \frac{m}{2} \log(2\pi) - \sum_{i=1}^m \frac{|y^{(i)} - \hat{y}(x^{(i)}; \omega)|^2}{2\sigma^2} \right) \\ &= \arg \min_{\omega} \left( \sum_{i=1}^m \frac{|y^{(i)} - \hat{y}(x^{(i)}; \omega)|^2}{m} \right) = \arg \min_{\omega} \text{MSE}. \quad (\#) \end{aligned}$$

denoted  $\hat{y}^{(i)}$

Remark:

In short, <sup>para.</sup> modelling the mean of Gaussians w/ fixed variance together maximum likelihood

$\Leftrightarrow$  para. modelling w/ MSE as cost function.