

# Estimating Probabilities from data

## - MLE Recap:

- Seek an estimate of  $\theta$  that maximizes the probability of the observed data

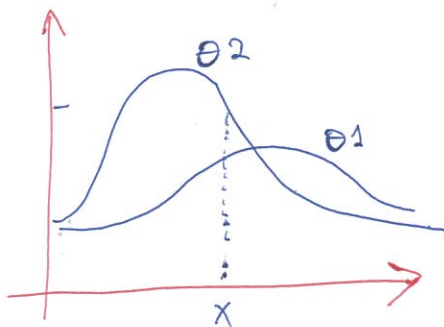
A function of  $x$  (m.s./d.f.)  $\rightarrow f(x; \theta)$

R.V.  $\rightarrow x$       fixed  $\rightarrow \theta$

- To do so, we define the likelihood function

$$L(\theta) = L(\theta|x) = f(x; \theta)$$

- MLE principle find  $\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} \log(f(x; \theta))$



Remark: Given that we observe the data, which parameter(s) would make it most likely that we observe what we observe

## - Extension to the data set

$$L(\theta) = L(\theta|D) = f(D; \theta) = f(x_1; \theta) \cdot f(x_2; \theta) \cdots f(x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$$

$$\Rightarrow \theta_{MLE} = \underset{\theta}{\operatorname{argmax}} \log \left( \prod_{i=1}^n f(x_i; \theta) \right)$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n (\log f(x_i; \theta))$$

## Summary:

- MLE gives the explanation of the data we observed.
- If  $n$  is large and  $f(x|\theta)$  is chosen correctly, then MLE will find the true parameter.
- MLE can overfit the data if  $n$  is small.
- If  $f(x|\theta)$  is chosen incorrectly, then MLE can be terribly wrong.

## MAP (Another way of estimating probabilities from data)

MLE: Estimate  $\theta_{MLE}$  that makes  $P(D|\theta)$  maximized.

R.V.      parameter (Fixed)

- MLE is frequentist statistics, meaning that  $\theta$  is just some constant.
- In Bayesian statistics,  $\theta$  can be r.v.

⇒ There is the distribution  $P(\theta)$

↑ Encode your belief of  $\theta$

Bayes rule:

$$P(\theta|D) = \frac{P(D|\theta) P(\theta)}{P(D)}$$

likelihood      prior on  $\theta$

Posterior

normalized form

$$\Rightarrow P(\theta|D) \propto P(D|\theta) P(\theta)$$

$P(D|\theta)$  ← which parameter makes our data the most likely

$P(\theta|D)$  ← Given that we have data, what is the most likely parameter

MAP Principle: find  $\theta_{\text{MAP}} = \underset{\theta}{\text{argmax}} \log P(\theta | D)$

$$= \underset{\theta}{\text{argmax}} \log P(D|\theta) P(\theta)$$

$$= \underset{\theta}{\text{argmax}} (\log P(D|\theta) + (\log P(\theta)))$$

MAP (Maximum a Posteriori Probability Estimation)

Simple Scenario: coin toss w/ prior knowledge

- Ex, suppose you toss a coin and observe  $D = \{H, H, H, H, H\}$

$$\text{MLE} \Rightarrow \theta_{\text{MLE}} = \frac{n_H}{n} = 1$$

- If you don't trust your estimate, then you can fix it by

$$\theta = \frac{n_H + m}{n + 2m}$$

- Here,  $m$  is the number of imaginary throws that would result in  $\frac{m}{2m} = 0.5$  (your hunch is close to 0.5)

- For large  $n$ ,  $\theta \rightarrow \theta_{\text{MLE}}$

- For small  $n$ , this incorporates your "prior belief" about what  $\theta$  should be

- Let's formalize this using MAP principle

Natural choice for the prior  $P(\theta)$  is the Beta distribution,

$$P(\theta) = \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)} \quad \leftarrow \text{constant}$$

$$\log P(\theta) \propto \log \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$\begin{aligned} P(\theta|D) &\propto P(D|\theta) P(\theta) \propto \binom{n}{n_H} \theta^{n_H} (1-\theta)^{n-n_H} \theta^{\alpha-1} (1-\theta)^{\beta-1} \\ &= \binom{n}{n_H} \theta^{n_H+\alpha-1} (1-\theta)^{n-n_H+\beta-1} \end{aligned}$$

$$\Rightarrow \theta_{\text{MAP}} = \frac{n_H + (\alpha - 1)}{n + (\alpha + \beta - 2)}$$

---