

Last lecture, we discussed the problem of Max Likelihood estimation, where given a dataset  $D$ , and a prob model whose parameters are specified

by  $\theta$ , we want to find the best parameter which makes observations  $D$  as likely / probable as possible

$$L(\theta) \triangleq p(D|\theta) \quad \max_{\theta} L(\theta) \quad \uparrow \text{likelihood}$$

$$\text{log-likelihood} \leftarrow \ell(\theta) = \log p(D|\theta) \stackrel{\{x_1, \dots, x_n\}}{\stackrel{\text{iid}}{=}} \sum_{i=1}^n \log p(x_i|\theta)$$

### Distributions over parameters

In the MLE, we were looking for best  $\theta$  that maximizes probability of observations,  $D$  what if we have some distribution over  $\theta$  (coming from prior information, e.g., the coin is fair)?

We want to find the most likely  $\theta$ , given observations,  $D$  i.e., want to maximize

$$\arg \max_{\theta} p(\theta|D) \quad \text{Maximum A Posteriori (MAP) estimation}$$

This is really a Bayesian framework, where data and parameters have distributions

$$\text{To do MAP,} \quad p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)} \propto p(D|\theta)p(\theta)$$

$$\arg \max_{\theta} p(\theta|D) = \arg \max_{\theta} p(D|\theta)p(\theta)$$

$$p(\theta|D) \propto p(D|\theta)p(\theta)$$

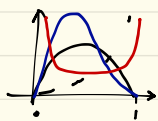
posterior
likelihood
prior

IE, we have some belief about  $\theta$  ( $p(\theta)$ ), we observe data  $D$ , how our belief on  $\theta$  changes, given  $D$  ( $p(\theta|D)$ )

Let's assume  $\theta$  has a Beta distribution (good dist for variable  $\in [0,1]$ )

$$p(\theta) = \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)}, \quad \theta \in [0,1]$$

$$\left\{ \begin{array}{l} \Gamma(x) \triangleq \int_0^\infty t^{x-1} e^{-t} dt \\ \Gamma(n) = (n-1)! \end{array} \right. \quad \rightarrow \quad \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$$



$E[\theta] = \frac{\alpha}{\alpha+\beta}$   
 $\alpha, \beta > 0$   
 $Mode = \frac{\alpha-1}{\alpha+\beta-2} \quad \alpha, \beta > 1$

We know  $p(D|\theta) = \theta^{\sum_n x_n} (1-\theta)^{N - \sum_n x_n}$

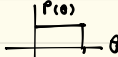
$$\rightarrow p(\theta|D) \propto \theta^{\alpha + \sum_n x_n - 1} (1-\theta)^{\beta + N - \sum_n x_n - 1} \propto_{\text{Beta}} (\alpha + \sum_n x_n, \beta + N - \sum_n x_n)$$

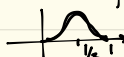
$\rightarrow$  Also a Beta distribution with mode

$$\begin{aligned} \arg\max_{\theta} p(\theta|D) &= \hat{\theta}_{MAP} = \frac{\sum_n x_n + \alpha - 1}{N + \alpha + \beta - 2} \\ &= \left( \frac{N}{N + \alpha + \beta - 2} \right) \frac{\sum_n x_n}{N} + \left( \frac{\alpha + \beta - 2}{N + \alpha + \beta - 2} \right) \frac{\alpha - 1}{\alpha + \beta - 2} \\ &= \text{Convex Comb} \left( \frac{\sum_n x_n}{N}, \frac{\alpha - 1}{\alpha + \beta - 2} \right) \end{aligned}$$

$N \rightarrow 0 \Rightarrow$  only prior

$N \rightarrow \infty \Rightarrow$  MLE (prior forgotten as  $N$  increases)

If  $p(\theta) = \text{uniform}$    $\rightarrow \alpha, \beta = 1 \rightarrow MAP = MLE$

If  $p(\theta)$   fair coin  $\rightarrow \alpha = \beta = 2 \rightarrow \frac{1}{N+2} + \frac{\sum_n x_n}{N+2}$