# 4 Estimating Probabilities of Data: MLE

Assumption: There is a distribution for data $p(x, y)$

If we have the distribution, then we can predict label using $p(y|x)$

Two ways to obtain $p(x, y)$

$p(x, y) = p(y|x) \, p(x)$      Discriminative learning    (Logistic Regression, SVM)

$\qquad\quad = p(x|y) \, p(y)$      Generative learning    (VAE, GAN)

## Maximum Likelihood Estimation

Binomial Distr:    $H, T, T, H, H, H, T, T, T, T$

What is $\mathbb{P}\{H\}$?    $\mathbb{P}\{H\} \approx \dfrac{n_H}{n_H + n_T}$

MLE: $\mathbb{P}\{D; \theta\}$    Let $\hat{\theta} = \underset{\theta}{\arg\max} \; \mathbb{P}\{D; \theta\}$

$\quad$ $D$ is data. Want to find $\theta$ that maximizes the prob of observing the data

$\quad \mathbb{P}\{D; \theta\} = \dbinom{n_H + n_T}{n_H} \theta^{n_H} (1-\theta)^{n_T}$

$\quad$ Find log-likelihood $\log$ monotonic increasing

$\quad \log(\mathbb{P}\{D; \theta\}) = \log\dbinom{n_H + n_T}{n_H} + n_H \log(\theta) + n_T \log(1-\theta)$

$\quad$ Maximize the log by taking derivative

$\quad \frac{\partial}{\partial \theta} \log(\mathbb{P}\{D; \theta\}) = \dfrac{n_H}{\theta} + \dfrac{n_T}{1-\theta} = 0$

$\qquad\qquad\qquad\qquad \theta = \dfrac{n_H}{n_H + n_T}$

Shortfalls: With small sample size, estimate unstable

$\qquad\qquad$ With 1 toss, predict H/T all the time

Fix: Smoothing. Hallucinate samples of

$\qquad\qquad \dfrac{n_H + \alpha}{n_H + n_T + \beta}$    $\alpha$ H in $\beta$ toss

## Bayesian vs. Frequentist

Frequentist $\mathbb{P}\{D; \theta\}$      $\theta$ is an unknown constant

$\qquad\qquad\qquad\qquad$ $\mathbb{P}\{\theta\}$ ill defined

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ MLE maximizes $\mathbb{P}\{D; \theta\}$

Bayesian $\mathbb{P}\{D|\theta\}$      $\theta$ is a distribution

$\qquad\qquad\qquad\qquad$ $\mathbb{P}\{\theta\}$ encodes your belief of what $\theta$ should be

# Baye's Rule

$P\{D|\theta\}$      Likelihood

$P\{\theta\}$      Prior

$P\{\theta|D\}$      Posterior

$$P\{\theta|D\} = \frac{P\{D|\theta\}\, P\{\theta\}}{P\{D\}}$$
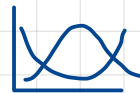
← Likelihood × Prior

← Normalization

## Back to coin toss example

Let prior be beta distribution

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

$B(\alpha,\beta)$ is normalizing constant

Depends on, $\alpha, \beta$,

Beta distr looks like 

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

$$= \binom{n_H + n_T}{n_H} \theta^{n_H}(1-\theta)^{n_T} \frac{\theta^{\alpha-1}(1-\theta)^{\beta-1}}{B(\alpha,\beta)}$$

$$= \binom{n_H + n_T}{n_H} \theta^{n_H + \alpha - 1}(1-\theta)^{n_T + \beta - 1}$$

If you were to do MLE on $P(\theta|D)$

$$\hat{\theta}_{MLE} = \frac{n_H + \alpha - 1}{n_H + n_T + \alpha + \beta - 1}$$

Same as smoothing