

## Probability &amp; Machine Learning

## 1. Overview

a. Maximum Likelihood Estimation  $\rightarrow$  parameter model (pdf, pmf)i.  $\text{Max } L(D|\theta) \rightarrow$  example  $\theta = (\mu, \sigma^2)$  when it is normalii.  $L(D; \theta) =$  multiply from  $i=1$  to  $i=n$   $(p^{x_i}) * (1-p)^{(1-x_i)}$ 1.  $X_i = 1$  if the  $i$ th toss is H, 0 otherwise2.  $P$ : Heads:  $\sum x_i$ 3.  $1-P$ : Tails:  $\sum 1-x_i$ 

b. In other words, we assume iid samples from some distribution

## 2. The Thumbtack problem (cont.)

a. We shall denote  $\text{Pr}(\text{HTHHTHT})$  as  $\text{Pr}(\text{HTHHTHT}; p)$ i.  $P$  is not a random variableb.  $D$  data,  $\theta$  parameterc. We refer to  $\text{Pr}(D|\theta)$  as the likelihood of the data under the modeld. In our problem,  $\text{Pr}(D;p) = p^4 (1-p)^3$ 

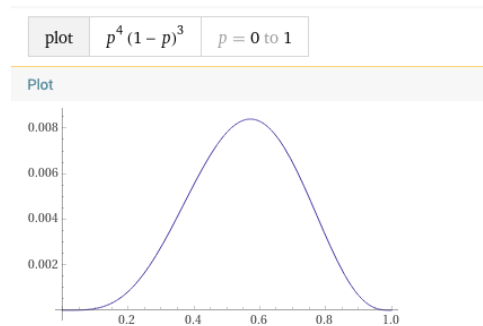
## 3. Maximum Likelihood Estimate (MLE)

a. Data: thumbtack tosses

b. Hypothesis: A flip is a Bernoulli distributed variable. Independence of flips

c. Learning: Find  $p^*$  that maximizes the data likelihood

$$\text{Pr}(\mathcal{D}; p) = p^4 (1 - p)^3$$



d.

4. Optimize to find  $p^*$

- $f(p) = p^4(1-p)^3$
- Max  $f(p)$
- Take the logarithm of  $L(D;p)$  to make life easier  $\rightarrow \log(L(D;p))$  since the product becomes the sum

$$\begin{aligned}
 p^* &= \arg \max_p \log(p^{a_H} (1-p)^{a_T}) \\
 \frac{d}{dp} (\log) &= \frac{d}{dp} (a_H \log p + a_T \log (1-p)) = \\
 &= a_H \frac{1}{p} + a_T \left(-\frac{1}{1-p}\right) = 0 \Rightarrow \\
 &\Rightarrow a_H/p = a_T/(1-p) \Rightarrow p^{a_T} = a_H(1-p) \Rightarrow \\
 &\Rightarrow p^{a_H+a_T} = a_H \Rightarrow \boxed{p_{MLE} = \frac{a_H}{n}}
 \end{aligned}$$

d.

e.  $P = 4/7$

5. Confidence intervals (reminder)

- Boston billionaire says: I want to know the true  $p$  within 0.01 accuracy, with confidence at least 95%
- Sampling theorem: Given  $n$  independent 0-1 RVs  $X_i$  such that  $\Pr(X_i = 1) = p$

$p(i=1, \dots, n)$  where  $n \geq \frac{3}{\epsilon^2} \ln\left(\frac{2}{\delta}\right)$  then the following holds:

$$\Pr\left(\left|\frac{\sum_{i=1}^n X_i}{n} - p\right| \leq \epsilon\right) \geq 1 - \delta$$

6. Two important properties of the MLE

- Consistent

$$\theta_{MLE} \rightarrow \theta_{true} \text{ in probability}$$

- Equivariant

$$\text{i. If } \theta_{MLE} \text{ is the MLE of } \theta_{true} \Rightarrow g(\theta_{MLE}) \text{ is the MLE of } g(\theta_{true})$$

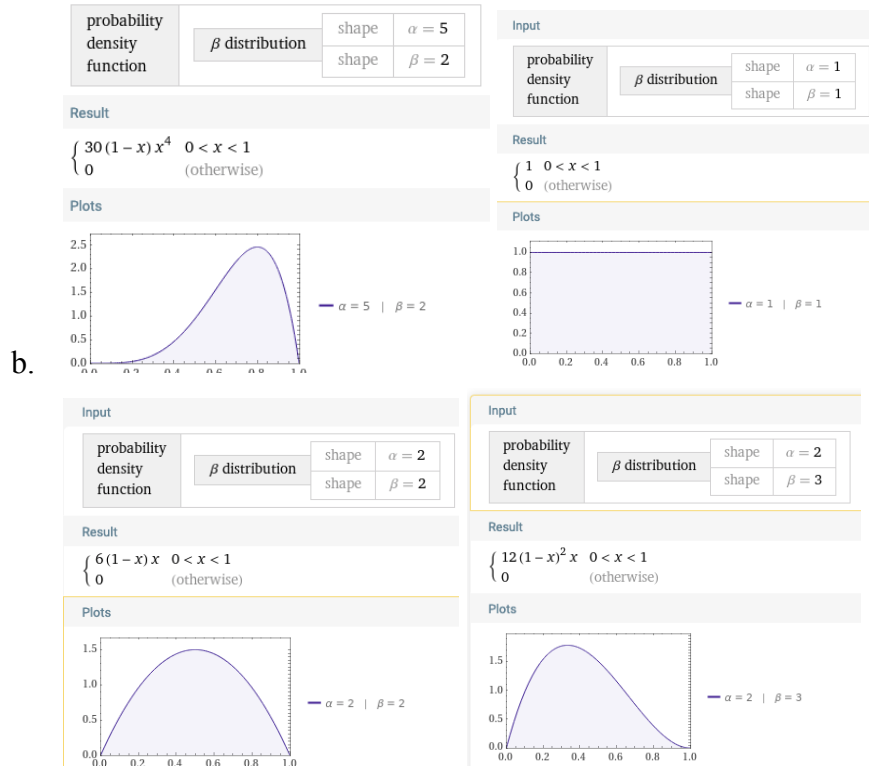
- ii. This means that if we want to the function it (example:  $p^2$ ), we simply put the function
  - iii. For example,  $\hat{p}^2$  for  $p^2$
- 7. Billionaire with prior beliefs
  - a. He says: Wait! I know that the thumbtack should be close to 50-50
  - b. You say: let's be Bayesian
  - c. Rather than learn a single value for  $p$ , we learn a probability distribution
    - i.  $P$  now becomes a random variable
- 8. Inference using Bayes' rule
  - a. Notice the notation  $\Pr(D|p)$  instead of  $\Pr(D;p)$
  - $$\Pr(p|D) = \frac{\Pr(D|p)\Pr(p)}{\Pr(D)}$$
  - b.
  - c.  $\Pr(D) \rightarrow$  does not depend on  $p$
- 9. Bayesian inference - summary
  - a. We choose the prior distribution  $f(\theta)$ . This distribution expresses our prior beliefs on the parameter  $\theta$ .
  - b. We choose the statistical model for the likelihood function  $f(D|\theta)$
  - c. After observing the data  $D=X_1, \dots, X_n$ , we update our beliefs and calculate the posterior distribution  $f(\theta|D)$
  - d. Maximum a posteriori (MAP) estimate is the mode of the posterior distribution
- 10. Important observation
  - a. If we impose a uniform prior on  $p$ , then
 
$$\Pr(p|D) \propto \Pr(D|p)$$
  - b. Image denoising lecture
    - i. Had we imposed a uniform prior on images  $x$ , then our MAP inference would be the same as the MLE
    - ii. Choosing a good prior is important in applications of Bayesian inference
- 11. Conjugate priors
  - a. Definition
    - i. "If the posterior distribution  $p(\theta|x)$  is in the same probability distribution family as the prior probability distribution  $p(\theta)$ , the prior and posterior are then called conjugate distributions, and the prior is called a conjugate prior for the likelihood function  $p(x|\theta)$ "  $\rightarrow$  assume that prior is same as likelihood
  - b. Conjugate prior: create fictitious data
    - i. For example: billionaire is sure that the probability of head to tails is 8:2
    - ii. Append  $800^H, 200^T$  (the more confident you are, the more data you collect) from the data  $H^4, T^3$ , you append the data
    - iii. Therefore,  $H \rightarrow 804, T \rightarrow 203$ , so  $P(H) = 804/1007$

## 12. Bayesian inference for the thumbtack problem

- a. The probability density for beta distribution is

$$f(x; a, b) = \frac{\Gamma(a+b)x^{a-1}(1-x)^{b-1}}{\Gamma(a)\Gamma(b)}, 0 \leq x \leq 1, a, b > 0$$

where  $\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx$  is the gamma function



- c. As you increase beta and theta, you increase the confidence

$$P(p) \propto p^{\beta_+ - 1} (1-p)^{\beta_- - 1} \text{ PRIOR}$$

$$P(D|p) = p^{a_+} (1-p)^{a_-} \text{ LIKELIHOOD.}$$

POSTERIOR

$$P(p|D) \propto P(D|p) P(p) \propto p^{a_+ + \beta_+ - 1} (1-p)^{a_- + \beta_- - 1}$$

$$\text{Beta}(a_+ + \beta_+, a_- + \beta_-)$$

fictitious coin tosses reflecting our prior belief.

- d.

### 13. Method of moments

- a. Suppose our model has parameters  $\theta = (\theta_1, \dots, \theta_k)$ 
  - i. Recall that the  $j$ -th moment of a RV  $X$  is  $E[X^j]$ . To denote that this is a function of the model, we write  $E_\theta[X^j]$ .
    1. Compute analytically,  $\alpha_j = E_\theta[X^j], j = 1, \dots, k$
  - ii. Consider the  $j$ -th sample moment for data  $x_1, \dots, x_n$  is
 
$$\hat{\alpha}_j = \frac{1}{n} \sum_{i=1}^n x_i^j, j = 1, \dots, k = \mu_{\text{MOM}}$$
  - iii.  $\sum x_i^2/n = \sigma^2_{\text{MOM}}$
  - iv. Equate the analytical moment expressions with the sample moments, and solve a system of  $k$  equations with unknowns to learn  $\theta = (\theta_1, \dots, \theta_k)$

### 14. Method of methods: example 1

- a. Let  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ 
  - i. We have one parameter, so  $k=1$ .
  - ii. The first moment is the mean  $\alpha_1 = E_p[X] = p$ .
  - iii. The first sample moment is the sample mean.
 
$$p_{\text{MOM}} = \frac{1}{n} \sum_{i=1}^n x_i$$
  - iv. Thus we directly get
  - v. Remark: Here, MOM is same as MLE, but this is not always the case

### 15. MLE vs MAP

- a. MLE
  - i. Goal: Find  $\theta$  maximizing the log-likelihood  $\Pr(x; \theta) \rightarrow$  also denoted as  $\Pr(x|\theta)$
- b. MAP
  - i. Goal: Find  $\theta$  maximizing the posterior  $\Pr(\theta|x)$
- c. In some cases, solving analytically for  $\theta$  is hard
- d. EM algorithm is an iterative approach to solving hard parameter learning problems

### 16. Two coin problem



$$\Pr(H) = \theta_A$$



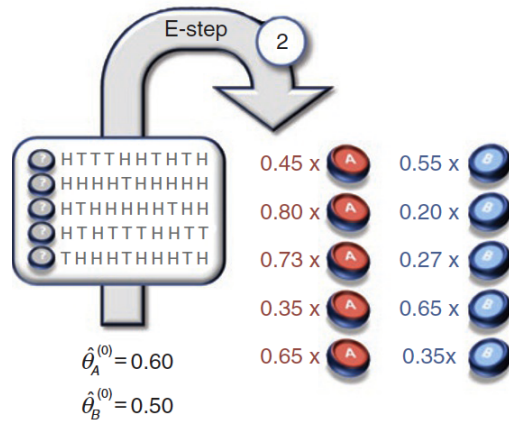
$$\Pr(H) = \theta_B$$

- a.
- b. Process
  - i. Suppose we choose a coin (A or B) uniformly at random
  - ii. We toss the coin  $n$  times, and record the total number of heads

- c. We repeat the process  $k$  times
- d. We have two unknown parameters:  $\theta_A, \theta_B$
- e. Suppose  $k = 5, n = 10$
- f. The data are two vectors
  - i.  $x = X_H = (x_1, x_2, x_3, x_4, x_5)$  #heads per round
  - ii.  $z = Z_c = (z_1, z_2, z_3, z_4, z_5)$  #coin id per round
- g. E.g.,  $x = (5, 9, 8, 4, 7), z = (B, A, A, B, A)$
- h. Solve two separate maximum posterior

Coin A	Coin B	
	5 H, 5 T	
9 H, 1 T		$\hat{\theta}_A = \frac{24}{24 + 6} = 0.80$
8 H, 2 T		
	4 H, 6 T	$\hat{\theta}_B = \frac{9}{9 + 11} = 0.45$
7 H, 3 T		
24 H, 6 T	9 H, 11 T	

- i.
- j. Use EM if there is missing data
- k. Suppose we only see the number of heads per round
  - i. In other words, we do not have access to  $z$
- l. We refer to  $z$  as hidden variables or latent factors
- m. Remarks
  - i. Not uncommon/common setting in data science applications  $\rightarrow$  missing data
  - ii. Clear that maximizing  $\Pr(x|\theta)$  is much harder than  $\Pr(x, z|\theta)$  in the presence of missing data  $z$
- n. We will proceed iteratively by updating
- o. Each iteration starts with a guess of the unknown parameters
 
$$\hat{\theta}^{(t)} = \left( \hat{\theta}_A^{(t)}, \hat{\theta}_B^{(t)} \right)$$
- p. E-step: a probability distribution over possible completions is computed using the current parameters  $\hat{\theta}^{(t)}$
- q. M-step: the new parameters are determined using the current completions
- r. Suppose our initial guess for the unknown variables are 0.6, 0.5
- s. E-step: what is the probability that round  $i$  comes from coin A/ coin B?



t.

*Likelihoods*

Coin A:  $\propto 0.6^5 \cdot 0.4^5 \approx 0.000796$

Coin B:  $\propto 0.5^5 \cdot 0.5^5 = \left(\frac{1}{2}\right)^{10} \approx 0.000976$

$$Pr(Z_1=A) = \frac{0.6^5 \cdot 0.4^5}{0.6^5 \cdot 0.4^5 + (0.5)^{10}} = 0.45, \quad Pr(Z_1=B) = 0.55$$

u.

- v. M-step: in order to learn  $\hat{\theta}^{(t+1)} = \left( \hat{\theta}_A^{(t+1)}, \hat{\theta}_B^{(t+1)} \right)$ , we first need to estimate the number of heads/tails from coins A/B given our estimate of the latent variables
- i. Notice that instead of being 100% certain whether a round was due to coin A or B, we have a probability distribution

**Two coin problem**

Round 1: 5H, 5T

HTTTHHTHTH 0.45 x A 0.55 x B

	Coin A	Coin B
Heads	0.45x5	0.55x5
Tails	0.45x5	0.55x5

Coin A	Coin B
$\approx 2.2$ H, $2.2$ T	$\approx 2.8$ H, $2.8$ T

We repeat this for all five rounds

w.

- x. Suppose we do one more round with different number of heads and tails

## One more round

Round 2: 9H, 1T

From the E-round we have

$$0.80 \times \text{A} \quad 0.20 \times \text{B}$$

Likelihoods Round 2

$$\text{Coin A: } \propto 0.6^9 \cdot 0.4^1 \approx 0.00403$$

$$\text{Coin B: } \propto 0.5^9 \cdot 0.5^1 = \left(\frac{1}{2}\right)^{10} \approx 0.000976$$

$$P_r(Z=A) = \frac{0.6^9 \cdot 0.4^1}{0.6^9 \cdot 0.4^1 + (0.5)^{10}} = 0.8049 \approx 0.8, P_r(Z=B) = 0.2$$

	Coin A	Coin B
Heads	0.8x9	0.2x9
Tails	0.8x1	0.2x1

$$\approx 7.2 \text{ H}, 0.8 \text{ T}$$

$$\approx 1.8 \text{ H}, 0.2 \text{ T}$$

y.

z. Having done this for all 5 rounds, we obtain the following

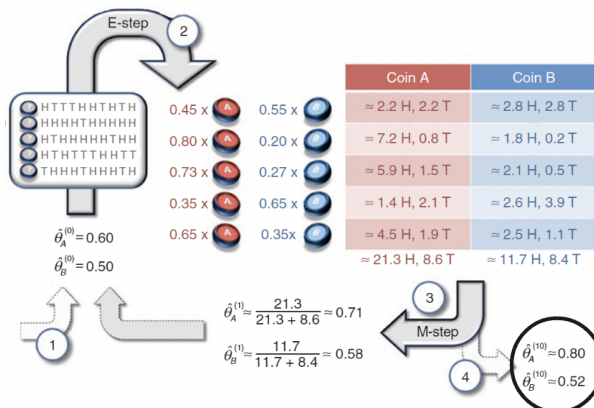
Coin A	Coin B
$\approx 2.2 \text{ H}, 2.2 \text{ T}$	$\approx 2.8 \text{ H}, 2.8 \text{ T}$
$\approx 7.2 \text{ H}, 0.8 \text{ T}$	$\approx 1.8 \text{ H}, 0.2 \text{ T}$
$\approx 5.9 \text{ H}, 1.5 \text{ T}$	$\approx 2.1 \text{ H}, 0.5 \text{ T}$
$\approx 1.4 \text{ H}, 2.1 \text{ T}$	$\approx 2.6 \text{ H}, 3.9 \text{ T}$
$\approx 4.5 \text{ H}, 1.9 \text{ T}$	$\approx 2.5 \text{ H}, 1.1 \text{ T}$
$\approx 21.3 \text{ H}, 8.6 \text{ T}$	$\approx 11.7 \text{ H}, 8.4 \text{ T}$

aa. The M-step now simply becomes the MLEs according to this data

$$\hat{\theta}_A^{(1)} \approx \frac{21.3}{21.3 + 8.6} \approx 0.71$$

$$\hat{\theta}_B^{(1)} \approx \frac{11.7}{11.7 + 8.4} \approx 0.58$$

bb. Expectation maximization





17. EM algorithm

- a. Suppose maximizing  $\Pr(x|\theta)$  has no closed form solution/intractable
- b. Key idea: latent variables  $z$  that make likelihood computations tractable
- c. Intuition:  $\Pr(x,z|\theta)$ ,  $\Pr(z|x,\theta)$  should be easy to compute after introducing the “right” latent variables
- d. EM guaranteed to converge to a local maximum
- e. Define “expected log”  $Q(\theta|\theta')$  where  $\theta'$  is the current estimate of  $\theta$ :

$$Q(\theta | \theta') = \sum_z \Pr(z | x, \theta') \log \Pr(x, z | \theta)$$

- f. The EM algorithm is an iterative method consisting of two steps
  - i. E-step: Find  $Q(\theta|\theta')$  in terms of the latent variables  $z$
  - ii. M-step: find  $\theta^*$  maximizing  $Q(\theta|\theta')$