# Probability, Maximum likelihood, and MAP Estimators

Tom Mitchell
Machine Learning 10-601
Jan 26 2009

part of this material is based on slides from Carlos Guestrin, and from Eric Xing

### You should know

- Events
  - discrete random variables, continuous random variables, compound events
- Axioms of probability
  - What defines a reasonable theory of uncertainty
- Independent events
- Conditional probabilities
- Bayes rule and beliefs
- Joint probability distribution

### **Expected values**

Given discrete random variable X, the expected value of X, written E[X] is

$$E[X] = \sum_{x \in \mathcal{X}} x P(X = x)$$

We also can talk about the expected value of functions of X

$$E[f(X)] = \sum_{x \in \mathcal{X}} f(x)P(X = x)$$

### Covariance

Given two discrete r.v.'s X and Y, we define the covariance of X and Y as

$$Cov(X,Y) = E[(X - E(X))(Y - E(Y))]$$

e.g., X=gender, Y=playsFootball or X=gender, Y=leftHanded

Remember: 
$$E[X] = \sum_{x \in \mathcal{X}} x P(X = x)$$

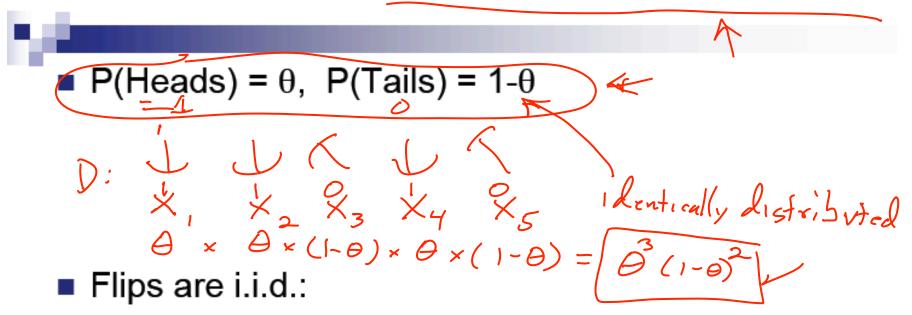
# Your first consulting job (PCYIX)

- A billionaire from the suburbs of Seattle asks you a question:
  - He says: I have thumbtack, if I flip it, what's the probability it will fall with the nail up?
  - ☐ You say: Please flip it a few times:



- □ You say: The probability is: 6
- ■He says: Why???
- ☐ You say: Because...

### Thumbtack – Binomial Distribution



- □ Independent events
- Identically distributed according to Binomial distribution
- Sequence *D* of  $\alpha_H$  Heads and  $\alpha_T$  Tails

data 
$$P(\mathcal{D} \mid \theta) = \theta^{lpha_H} (1- heta)^{lpha_T}$$

#### Maximum Likelihood Estimation



- **Data:** Observed set *D* of  $\alpha_H$  Heads and  $\alpha_T$  Tails
- Hypothesis: Binomial distribution
- Learning θ is an optimization problem
  - □ What's the objective function?
- MLE: Choose θ that maximizes the probability of observed data:

$$\widehat{\theta} = \arg \max_{\theta} P(\mathcal{D} \mid \theta)$$

$$= \arg \max_{\theta} \ln P(\mathcal{D} \mid \theta)$$

### Maximum Likelihood Estimate for $\Theta$



$$\widehat{\theta} = \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta)$$

$$= \arg\max_{\theta} \ln \theta^{\alpha_H} (1-\theta)^{\alpha_T}$$

Set derivative to zero:

$$rac{d}{d heta}$$
 In  $P(\mathcal{D} \mid heta) = 0$ 

100

Set derivative to zero:

$$rac{d}{d heta} \, \ln P(\mathcal{D} \mid heta) = 0$$

$$\hat{\theta} = \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

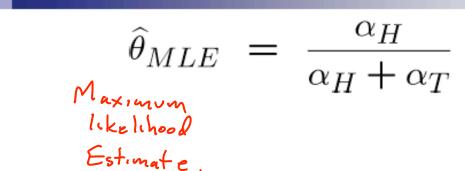
$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

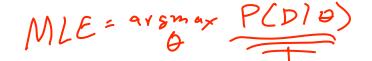
$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative to Zero.}$$

$$= \arg\min_{\theta} \ln P(\mathcal{D} \mid \theta) \qquad \text{set derivative$$

### How many flips do I need?



### Bayesian Learning





arsmar 
$$P(\theta \mid \mathcal{D}) = \frac{P(\mathcal{D} \mid \theta)P(\theta)}{P(\mathcal{D})}$$

Or equivalently:

$$P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta) P(\theta)$$

## Beta prior distribution – $P(\theta)$

$$P(\theta) = \frac{\theta^{\beta_H - 1}(1-\theta)^{\beta_T - 1}}{B(\beta_H, \beta_T)} \sim Beta(\beta_H, \beta_T) \qquad \text{Mean:}$$

$$Beta(\beta_H, \beta_T) \qquad beta(\beta_H, \beta_T) \qquad beta(\beta_H,$$

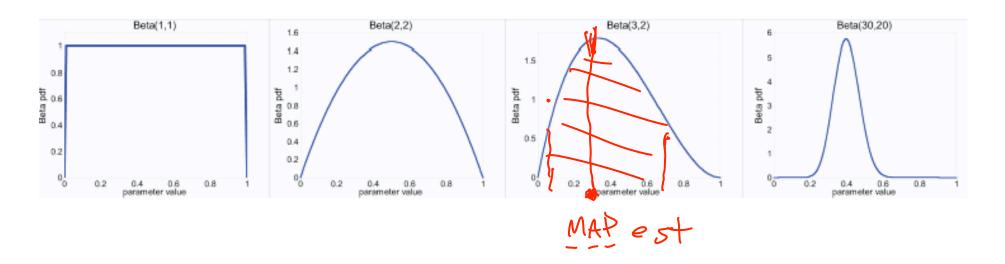
- Likelihood function:  $P(\mathcal{D} \mid \theta) = \theta^{\alpha_H} (1 \theta)^{\alpha_T}$
- Posterior:  $P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta)P(\theta)$

### Posterior distribution

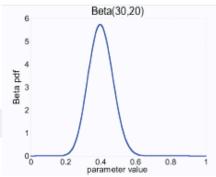


- Prior:  $Beta(\beta_H, \beta_T)$
- Data:  $\alpha_H$  heads and  $\alpha_T$  tails
- Posterior distribution:

$$P(\theta \mid \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$



## MAP for Beta distribution





$$P(\theta \mid \mathcal{D}) = \frac{\theta^{\beta_H + \alpha_H - 1} (1 - \theta)^{\beta_T + \alpha_T - 1}}{B(\beta_H + \alpha_H, \beta_T + \alpha_T)} \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

MAP: use most likely parameter:

$$\widehat{\theta} = \arg\max_{\theta} P(\theta \mid \mathcal{D}) = \underbrace{\begin{array}{c} \mathcal{L}_{\mathcal{H}} + \mathcal{B}_{\mathcal{H}} \\ \mathcal{L}_{\mathcal{H}} + \mathcal{L}_{\mathcal{T}} + \mathcal{B}_{\mathcal{H}} + \mathcal{B}_{\mathcal{T}} \end{array}}_{\mathcal{L}_{\mathcal{H}} + \mathcal{L}_{\mathcal{T}}}$$

- Beta prior equivalent to extra thumbtack flips
- As  $N \to \infty$ , prior is "forgotten"
- But, for small sample size, prior is important!

### Dirichlet distribution

- number of heads in N flips of a two-sided coin
  - follows a binomial distribution
  - Beta is a good prior (conjugate prior for binomial)
- what it's not two-sided, but k-sided?
  - follows a multinomial distribution
  - Dirichlet distribution is the conjugate prior

$$P( heta_1, heta_2,... heta_K) = rac{1}{B(lpha)} \prod_i^K heta_i^{(lpha_1-1)}$$



### **Estimating Parameters**

• Maximum Likelihood Estimate (MLE): choose  $\theta$  that maximizes probability of observed data  $\mathcal{D}$ 

$$\widehat{\theta} = \arg \max_{\theta} P(\mathcal{D} \mid \theta)$$

 Maximum a Posteriori (MAP) estimate: choose θ that is most probable given prior probability and the data

$$\widehat{\theta} = \arg\max_{\theta} P(\theta \mid \mathcal{D})$$

$$= \arg\max_{\theta} = \frac{P(\mathcal{D} \mid \theta)P(\theta)}{P(\mathcal{D})}$$

### You should know

#### Probability basics

- random variables, events, sample space, conditional probs, ...
- independence of random variables
- Bayes rule
- Joint probability distributions
- calculating probabilities from the joint distribution

#### Point estimation

- maximum likelihood estimates
- maximum a posteriori estimates
- distributions binomial, Beta, Dirichlet, ...





- Data:
  - We observed Niid coin tossing: D={1, 0, 1, ..., 0}
- Representation:



$$x_n = \{0,1\}$$

$$P(x) = \begin{cases} 1 - \theta & \text{for } x = 0 \\ \theta & \text{for } x = 1 \end{cases} \Rightarrow P(x) = \theta^{x} (1 - \theta)^{1 - x}$$

• How to write the likelihood of a single observation  $x_i$ ?

$$P(x_i) = \theta^{x_i} (1 - \theta)^{1 - x_i}$$

The likelihood of dataset D={x<sub>1</sub>, ...,x<sub>N</sub>}:

$$P(x_{1}, x_{2}, ..., x_{N} \mid \theta) = \prod_{i=1}^{N} P(x_{i} \mid \theta) = \prod_{i=1}^{N} \left(\theta^{x_{i}} (1 - \theta)^{1 - x_{i}}\right) = \theta^{\sum_{i=1}^{N} x_{i}} (1 - \theta)^{\sum_{i=1}^{N} 1 - x_{i}} = \theta^{\text{\#head}} (1 - \theta)^{\text{\#tails}}$$

### Conditional Independence

Definition: X is <u>conditionally independent</u> of Y given Z, if the probability distribution governing X is independent of the value of Y, given the value of Z

$$(\forall i, j, k) P(X = x_i | Y = y_j, Z = z_k) = P(X = x_i | Z = z_k)$$

Which we often write

$$P(X|Y,Z) = P(X|Z)$$

E.g., 
$$P(Thunder|Rain, Lightning) = P(Thunder|Lightning)$$