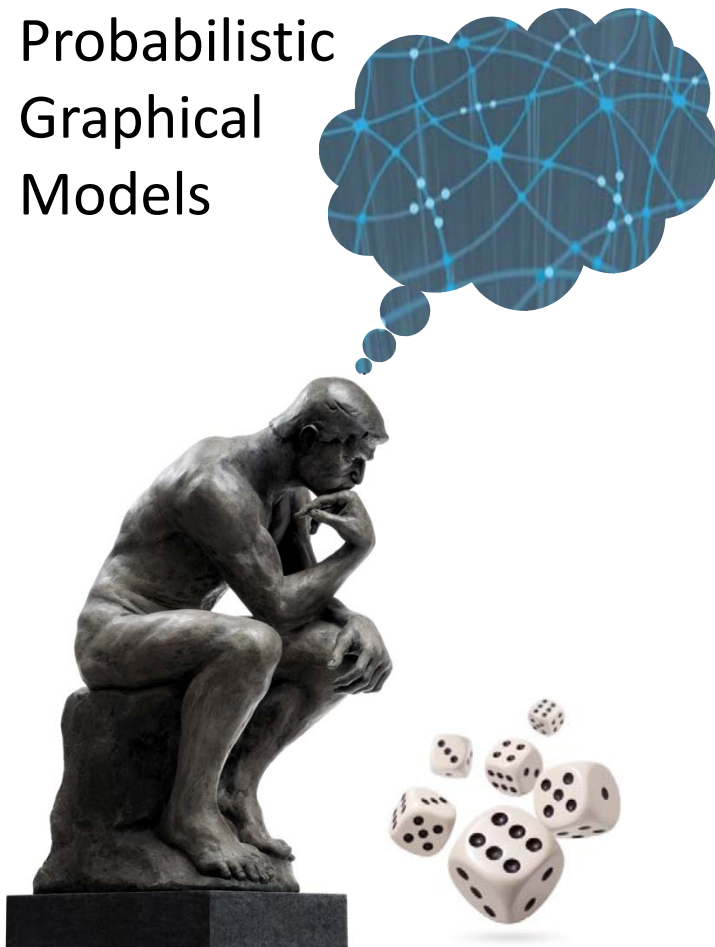


Probabilistic  
Graphical  
Models



Learning

---

Parameter Estimation

---

Maximum  
Likelihood  
Estimation

# Biased Coin Example

$P$  is a Bernoulli distribution:

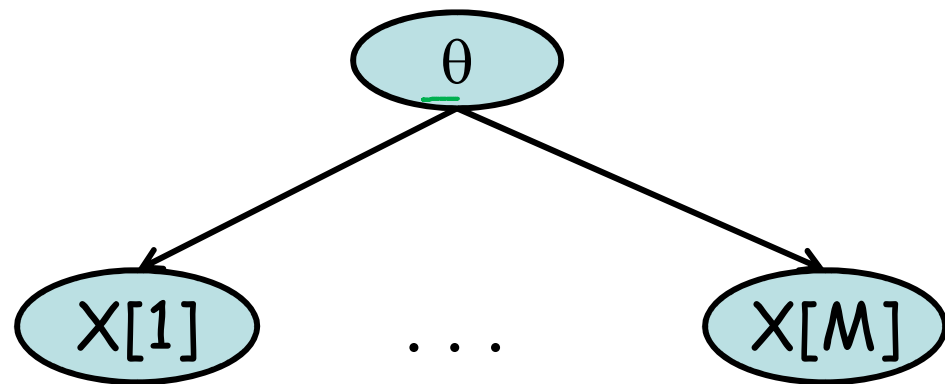
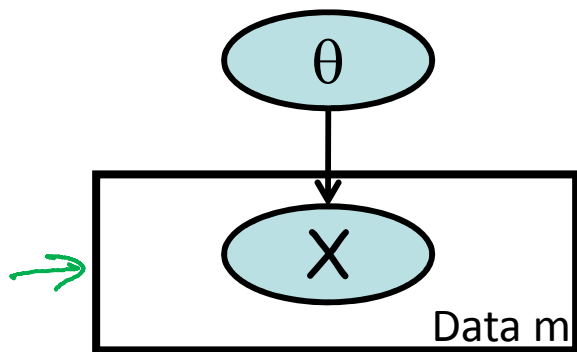
$$P(\underline{X}=1) = \underline{\theta}, P(\underline{X}=0) = \underline{1-\theta}$$



$\underline{\mathcal{D}} = \{x[1], \dots, x[M]\}$  sampled IID from  $P$

- Tosses are independent of each other
- Tosses are sampled from the same distribution (identically distributed)

# IID as a PGM



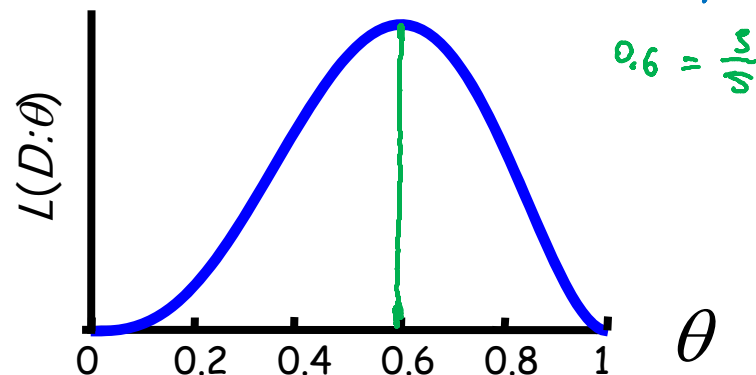
$$\underline{P(x[m] | \theta)} = \begin{cases} \underline{\theta} & x[m] = \underline{x^1} \\ \underline{1 - \theta} & x[m] = \underline{x^0} \end{cases}$$

# Maximum Likelihood Estimation

- **Goal:** find  $\theta \in [0,1]$  that predicts  $D$  well
- **Prediction quality = likelihood of  $D$  given  $\theta$**

$$L(\theta : D) = P(D | \theta) = \prod_{m=1}^M P(x[m] | \theta)$$

$L(\theta : \langle H, T, T, H, H \rangle) = \underbrace{P(H|\theta)}_{\theta} \cdot \underbrace{P(T|\theta)}_{(1-\theta)} \cdot \underbrace{P(T|\theta)}_{(1-\theta)} \cdot \underbrace{P(H|\theta)}_{\theta} \cdot \underbrace{P(H|\theta)}_{\theta} = \theta^3 (1-\theta)^2$



# Maximum Likelihood Estimator

- Observations:  $M_H$  heads and  $M_T$  tails
- Find  $\theta$  maximizing likelihood

- Equivalent to maximizing log-likelihood

$$l(\theta : M_H, M_T) = M_H \log \theta + M_T \log(1 - \theta)$$

- Differentiating the log-likelihood and solving for  $\theta$ :

$$\hat{\theta} = \frac{M_H}{M_H + M_T}$$

# Sufficient Statistics

- For computing  $\theta$  in the coin toss example, we only needed  $M_H$  and  $M_T$  since

$$L(\theta : D) = \theta^{M_H} (1 - \theta)^{M_T}$$

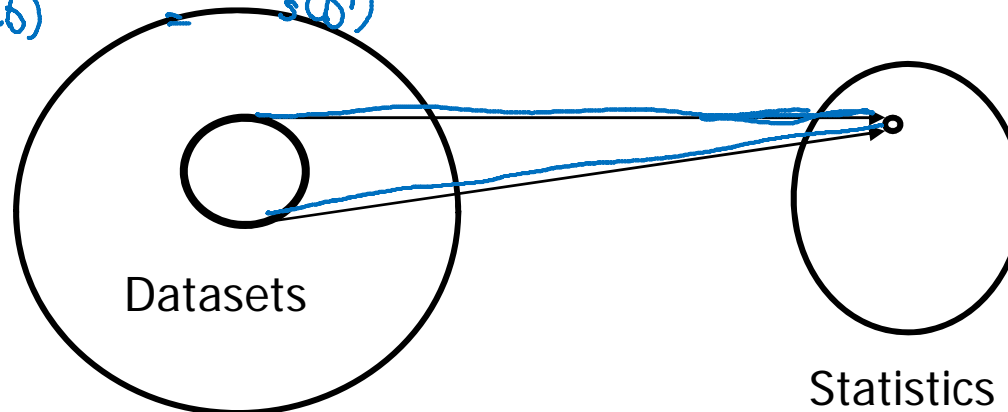
- $\rightarrow M_H$  and  $M_T$  are sufficient statistics

# Sufficient Statistics

- A function  $s(D)$  is a sufficient statistic from instances to a vector in  $\mathbb{R}^k$  if for any two datasets  $D$  and  $D'$  and any  $\theta \in \Theta$  we have

$$\sum_{x[i] \in D} s(x[i]) = \sum_{x[i] \in D'} s(x[i]) \Rightarrow L(\theta : D) = L(\theta : D')$$

$\underbrace{\sum_{x[i] \in D} s(x[i])}_{s(D)} = \underbrace{\sum_{x[i] \in D'} s(x[i])}_{s(D')}$



# Sufficient Statistic for Multinomial

- For a dataset  $D$  over variable  $X$  with  $k$  values, the sufficient statistics are counts  $\langle \bar{M}_1, \dots, \bar{M}_k \rangle$  where  $M_i$  is the # of times that  $X[m]=x^i$  in  $D$
- Sufficient statistic  $s(x)$  is a tuple of dimension  $k$ 
  - $s(x^i) = (0, \dots, 0, 1, 0, \dots, 0)$   $\sum_n s(x[m]) = \{M_1, M_2, \dots, M_k\}$

$$L(\theta : D) = \prod_{i=1}^k \theta_i^{M_i} \quad \text{where } \theta_i \text{ is param for } x=x^i$$



# Sufficient Statistic for Gaussian

- Gaussian distribution:

$$P(X) \sim N(\underline{\mu}, \underline{\sigma^2}) \quad \text{if} \quad p(X) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

- Rewrite as

$$p(X) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-x^2 \frac{1}{2\sigma^2} + x \frac{\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2}\right)$$

- Sufficient statistics for Gaussian:

$$s(x) = \langle 1, x, x^2 \rangle \quad s(D) = \left( \sum_n x[n]^2, \sum_n x[n], n \right)$$

# Maximum Likelihood Estimation

- MLE Principle: Choose  $\theta$  to maximize  $L(D;\Theta)$

- Multinomial MLE:  $\hat{\theta}^i = \frac{M_i}{\sum_{i=1}^m M_i}$  *fraction of  $x_i$  in data*

- Gaussian MLE:  
 $\hat{\mu} = \frac{1}{M} \sum_m x[m]$  *empirical mean*  
 $\hat{\sigma} = \sqrt{\frac{1}{M} \sum_m (x[m] - \hat{\mu})^2}$  *empirical st dev*

# Summary

- Maximum likelihood estimation is a simple principle for parameter selection given  $D$
- Likelihood function uniquely determined by sufficient statistics that summarize  $D$
- MLE has closed form solution for many parametric distributions