

# CS 583: PROBABILISTIC GRAPHICAL MODELS

## TOPIC: PARAMETER ESTIMATION CHAPTER: 17



**Mustafa Bilgic**

 <http://www.cs.iit.edu/~mbilgic>

 <https://twitter.com/bilgicm>

# PARAMETER ESTIMATION FOR BNs

- Assume the network structure is given
- The data  $\mathcal{D}$  consists of fully observed instances of the network variables
  - $\mathcal{D} = \{x[1], x[2], \dots, x[n]\}$
- Estimate the network parameters, i.e., learn the CPDs
- Two approaches
  1. Maximum likelihood estimation
  2. Bayesian estimation

# SIMPLEST CASE – ONE VARIABLE

- Imagine we have a thumbtack
- Flip it, and it comes as heads or tails

heads



tails



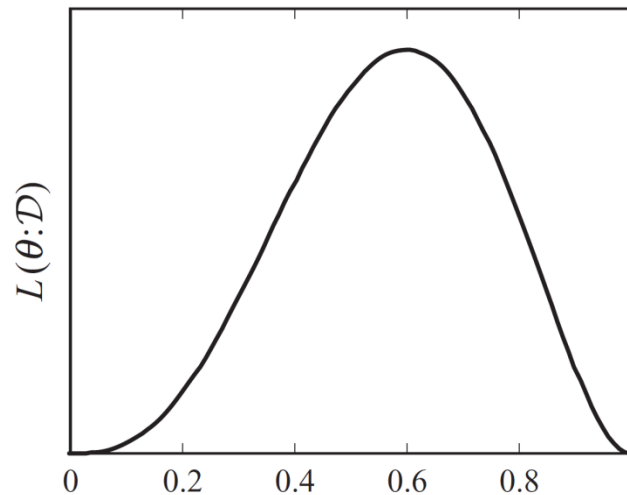
- Assume we flip it 100 times and it comes head 30 times
- What is  $\theta$ ?

# THUMBTRACK TOSSES

- Assume we have a set of thumbtrack tosses
  - $\mathcal{D} = \{\mathbf{x}[1], \dots, \mathbf{x}[n]\}$
- Also assume each toss,  $\mathbf{x}[i]$ , is IID
- We define a *hypothesis space*  $\Theta$ 
  - $\Theta$  is the set of all parameters  $\theta \in [0, 1]$
- We formulate an *objective function*
  - The objective function tells us how good a given hypothesis (in this case  $\theta$ ) is

# LIKELIHOOD

- What is the probability, or *likelihood*, of seeing the sequence H, T, T, H, H?
  - $\theta * (1 - \theta) * (1 - \theta) * \theta * \theta = \theta^3(1 - \theta)^2$



When is  $L(\theta; \mathcal{D})$  maximum?

# LIKELIHOOD/LOG-LIKELIHOOD

- Number of heads =  $h$ , number of tails =  $t$
- Likelihood:  $L(\theta; \mathcal{D}) = \theta^h(1-\theta)^t$
- Log-likelihood:  $l(\theta; \mathcal{D}) = h\ln\theta + t\ln(1-\theta)$
- Find  $\theta$  that maximizes the log-likelihood
- Take derivate of  $l(\theta; \mathcal{D})$  with respect to  $\theta$  and set it to zero

# MAXIMUM LIKELIHOOD FOR A MULTINOMIAL

- Domain of  $X$  is  $\{A, B, C\}$
- We see  $A$   $a$  times,  $B$   $b$  times, and  $C$   $c$  times.
- $P(X=A)$  is  $p$ ,  $P(X=B)$  is  $q$ , and  $P(C) = 1 - p - q$
- What are  $p$  and  $q$ ?
- Proof?

# CONSTRAINED OPTIMIZATION

- Assume  $X$  can take  $k$  values
- $P(X=x_i) = \theta_i$
- Find  $\theta$  that maximizes the entropy
  - $H(X) = -\sum_i \theta_i \log_2 \theta_i$
- If we take the partial derivative w.r.t.  $\theta_i$ 
  - ...



# CONSTRAINED OPTIMIZATION

Find  $\boldsymbol{\theta}$   
maximizing  $f(\boldsymbol{\theta})$   
subject to

$$c_1(\boldsymbol{\theta}) = 0$$

...

$$c_m(\boldsymbol{\theta}) = 0$$

Form the Lagrangian:

$$F(\boldsymbol{\theta}, \boldsymbol{\lambda}) = f(\boldsymbol{\theta}) - \sum_{j=1}^m \lambda_j c_j(\boldsymbol{\theta})$$

# LAGRANGE MULTIPLIERS EXAMPLES

1. Maximize  $x*y$  st.  $x+y = 10$
2. Maximize  $x+y$  st.  $x^2+y^2 = 1$
3. Entropy
4. Maximum likelihood estimate for a multinomial

# ML FOR BNs

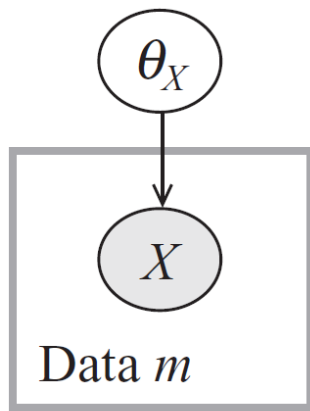
- Simple structure
  - $X \rightarrow Y$
- General structure
  - The key is that the parameters for each variable can be optimized independently
  - Examples

# BAYESIAN ESTIMATION

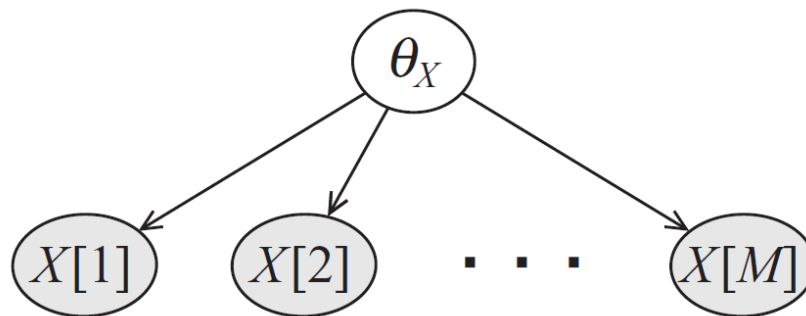
- Assume we flip a coin 10 times and we get 4 Heads, 6 Tails
  - What is  $P(C=H)$ ?
- Assume we flip a thumbtack 10 times and we get 4 Heads, 6 Tails
  - What is  $P(T=H)$ ?
- What if we repeat the flips 10M times and we get 4M Heads and 6M Tails?
- Bayesian estimation will let us encode our *prior knowledge*

# INDEPENDENCE?

- Earlier, we assumed the tosses are independent
- This is true if we know  $\theta$
- If we don't know  $\theta$ , then each toss tells us something about  $\theta$ , thus the next toss



(a)

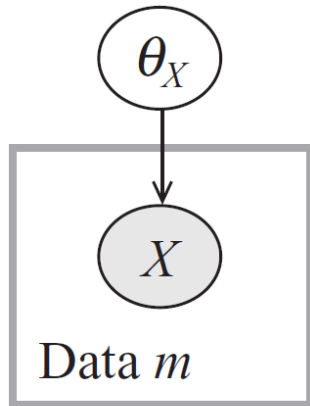


(b)

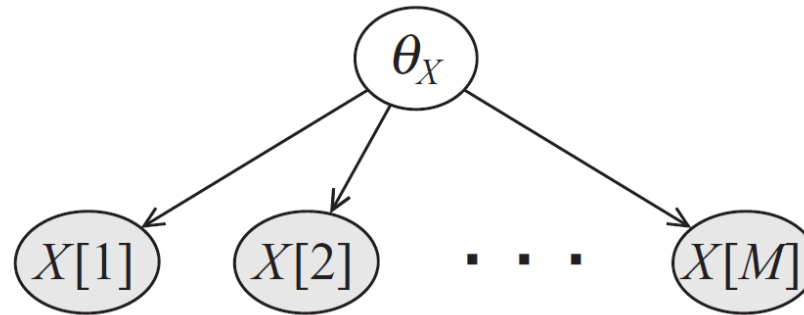
# BAYESIAN ESTIMATION

- Rather than a single  $\theta$ , we will instead have a probability distribution,  $P(\theta)$ , over  $\theta$

# BAYESIAN ESTIMATION



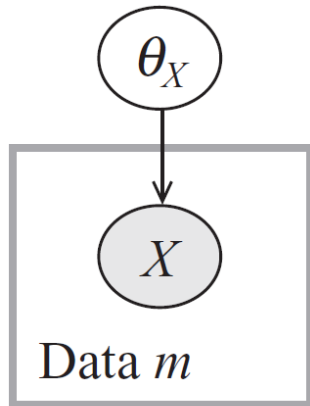
(a)



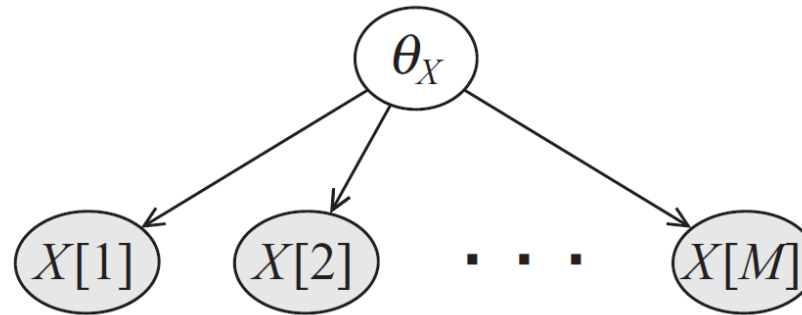
(b)

- We treat the parameter  $\theta$  as a random variable
- We ascribe a prior probability to  $\theta$ ,  $P(\theta)$ , encoding our prior knowledge

# PARAMETERS



(a)



(b)

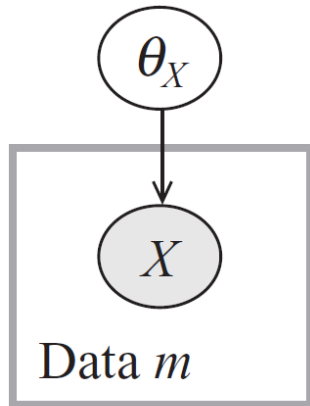
- $P(X[i] = x^1 \mid \theta_x) = \theta; P(X[i] = x^0 \mid \theta_x) = (1 - \theta)$
- $P(\theta_x)$ ?
  - A continuous distribution over the interval  $[0,1]$



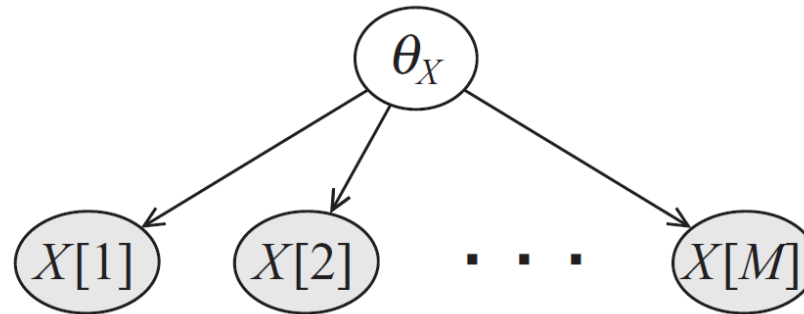
# POSTERIOR AND PREDICTION

- We are interested in
  - The probability of the next instance, given data
    - $P(x[M+1] \mid D)$
  - The posterior distribution of  $\theta$  given data
    - $P(\theta \mid D)$

# FACTORIZATION



(a)



(b)

$$\begin{aligned} P(x[1], x[2], \dots, x[M], \theta_X) &= P(\theta_X) \prod_{m=1}^M P(x[m] | \theta_X) \\ &= P(\theta_X) \theta^{M^{[1]}} (1 - \theta)^{M^{[0]}} \end{aligned}$$

# POSTERIOR AND $P(x[M+1] | D)$

Posterior distribution

$$P(\theta_x | D) = \frac{P(x[1], \dots, x[M] | \theta_x) P(\theta_x)}{P(x[1], \dots, x[M])}$$

$$\begin{aligned} P(x[M+1] | D) &= \int_0^1 P(x[M+1] | \theta_x, x[1], \dots, x[M]) P(\theta_x | x[1], \dots, x[M]) d\theta \\ &= \int_0^1 \underbrace{P(x[M+1] | \theta_x)}_{\substack{\theta \text{ or } 1-\theta \\ \text{(if binary)}}} \underbrace{P(\theta_x | x[1], \dots, x[M])}_{\text{Posterior}} d\theta \end{aligned}$$

Think of taking a weighted average

## $P(x[M+1] \mid D)$

$$\begin{aligned} P(x[M+1] \mid x[1], \dots, x[M]) &= \int_0^1 P(x[M+1] \mid \theta_x) P(\theta_x \mid x[1], \dots, x[M]) d\theta \\ &= \int_0^1 P(x[M+1] \mid \theta_x) \frac{P(\theta_x) P(x[1], \dots, x[M] \mid \theta_x)}{P(x[1], \dots, x[M])} d\theta \end{aligned}$$

$P(x[1], \dots, x[M])$  is a constant

$$P(x[M+1] \mid x[1], \dots, x[M]) \propto \int_0^1 P(x[M+1] \mid \theta_x) P(\theta_x) P(x[1], \dots, x[M] \mid \theta_x) d\theta$$

# UNIFORM PRIOR

- We have a uniform prior over  $\theta_x$ . That is,  $p(\theta_x)=1$
- $P(X[M+1]=x^1 \mid x[1], \dots, x[M])?$

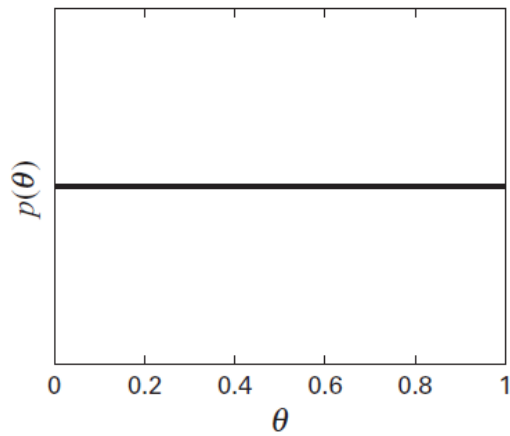
# UNIFORM PRIOR

- We have a uniform prior over  $\theta_x$ . That is,  $p(\theta_x)=1$
- $P(X[M+1]=x^1 \mid x[1], \dots, x[M])$ ? That is,  $P(X[M+1]=x^1 \mid D)$ ?
- For the binary case,  $P(X[M+1]=x^1 \mid D) = (t + 1) / (t + f + 2)$ , where  $t$  is the number of True cases and  $f$  is the number of False cases in  $D$
- This is also called *Laplace smoothing*
- What about the posterior,  $P(\theta \mid D)$ , if the prior  $P(\theta)$  is uniform?

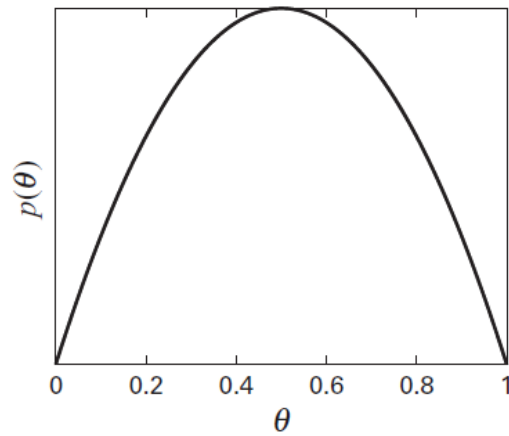
# BETA DISTRIBUTION

- $\theta \sim \text{Beta}(\alpha, \beta)$  if  $P(\theta) = \gamma \theta^{\alpha-1} (1-\theta)^{\beta-1}$  where  $\gamma$  is a normalizing constant
- Mean:  $\alpha/(\alpha+\beta)$
- Mode:  $(\alpha-1)/(\alpha+\beta-2)$
- Note that the mode is closer to the mean when  $\alpha$  and  $\beta$  are large
- Read more at
  - [https://en.wikipedia.org/wiki/Beta\\_distribution](https://en.wikipedia.org/wiki/Beta_distribution)

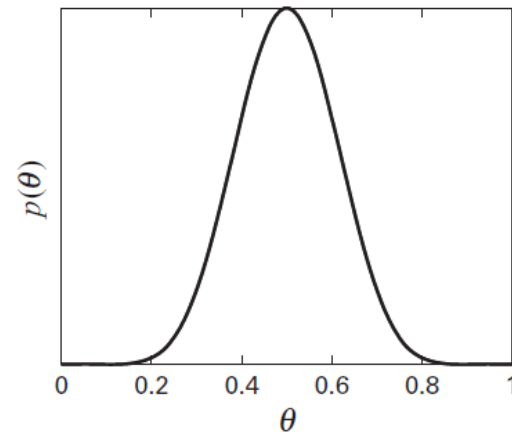
# BETA DISTRIBUTION



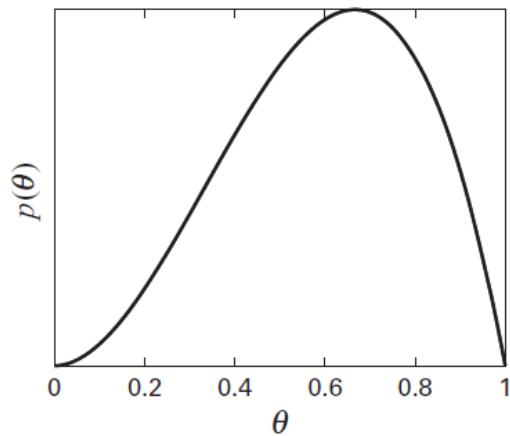
*Beta(1,1)*



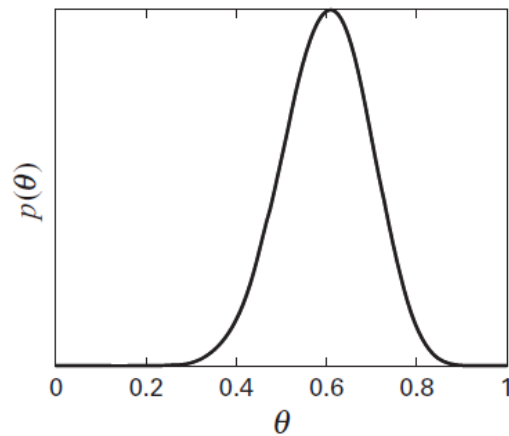
*Beta(2,2)*



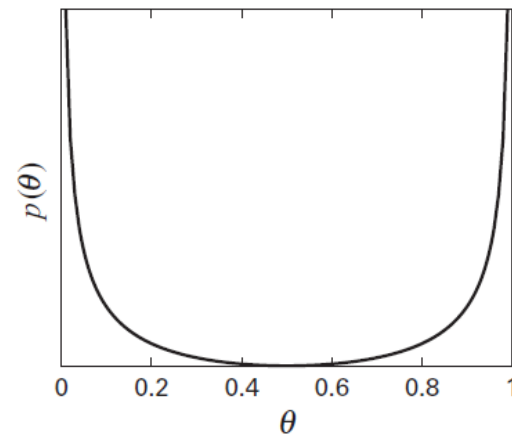
*Beta(10,10)*



*Beta(3,2)*



*Beta(15,10)*



*Beta(0.5,0.5)*



# BETA DISTRIBUTION

- What is  $P(X[M+1]=x^1 \mid D)$  if the prior is  $\text{Beta}(\alpha, \beta)$ ?
  - $P(X[M+1]=x^1 \mid D) = (p + \alpha) / (p + n + \alpha + \beta)$
- What is the posterior,  $P(\theta \mid D)$ , if the prior is  $\text{Beta}(\alpha, \beta)$ ?
  - $P(\theta \mid D) = \text{Beta}(p + \alpha, n + \beta)$
- $\alpha$  and  $\beta$  work like pseudo-counts for the positive and negative cases respectively
- What values to choose for  $\alpha$  and  $\beta$ ?
  - It depends on our belief and the strength of our belief

# DIRICHLET PRIORS

- Generalizes the Beta distribution for multinomials

$$\theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K) \text{ if } P(\theta) \propto \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

- What is  $P(X[M+1]=x^i \mid D)$  if the prior is Dirichlet?

- $P(X[M+1]=x^i \mid D) = (n_i + \alpha_i) / (|D| + \alpha)$  where  $n_i$  is the number of times the  $i^{\text{th}}$  case appears in  $D$  and  $\alpha = \alpha_1 + \alpha_2 + \dots + \alpha_K$

- What is the posterior,  $P(\theta \mid D)$ , if the prior is Dirichlet?

- $P(\theta \mid D) = \text{Dirichlet}(n_1 + \alpha_1, n_2 + \alpha_2, \dots, n_K + \alpha_K)$

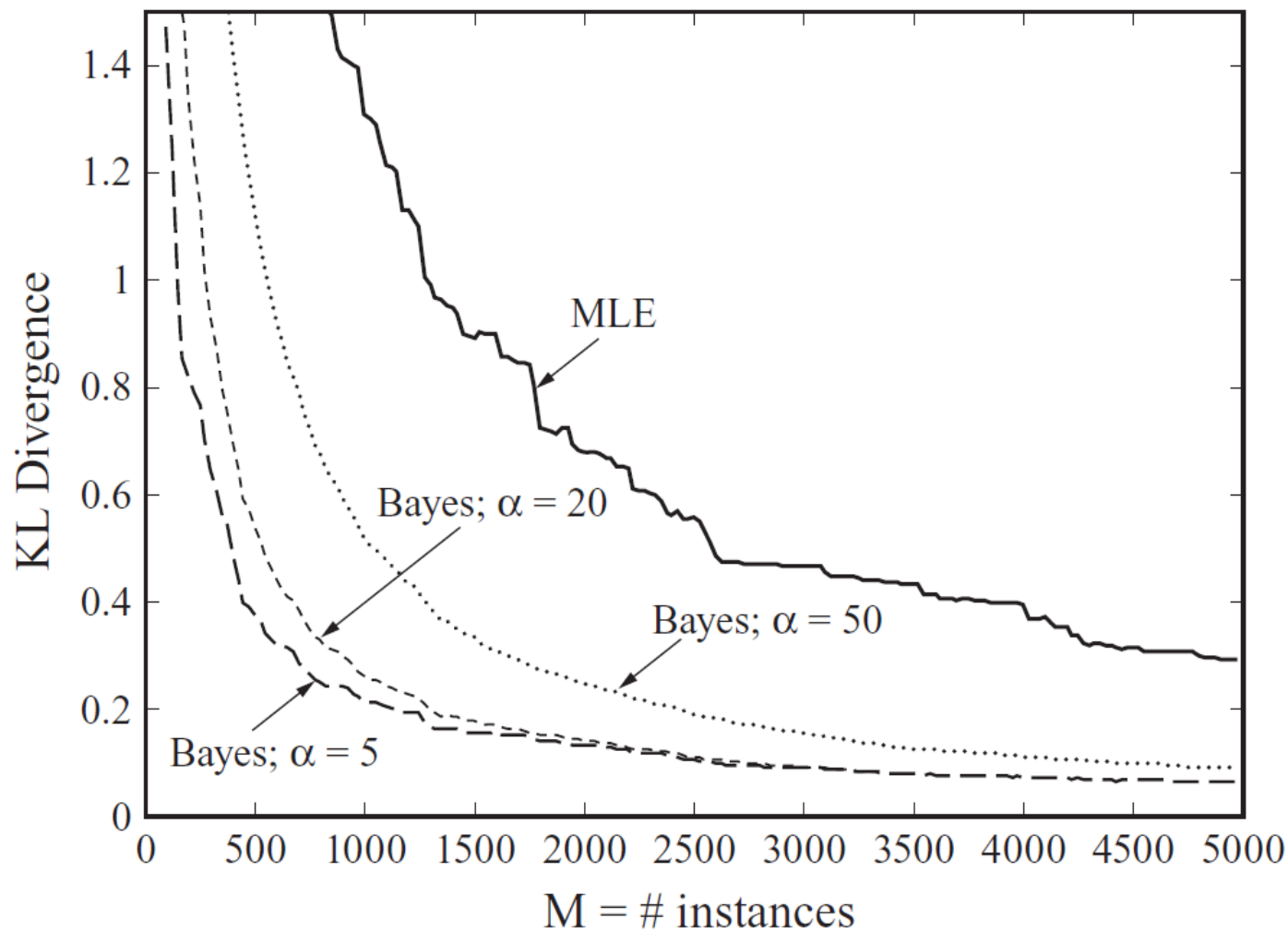
# BAYESIAN ESTIMATION

- In MLE for BNs, we optimized each parameter independently
- Can we do the same for Bayesian estimation for BNs?
  - Only if the prior also factorizes wrt the BN
- What about the priors? How do we choose them?
  1. Ask the prior for each variable to an expert
  2. Use the same prior for all variables
    - This is called the *K2 prior*
  3. Imagine a dataset  $D'$  of imaginary instances
    - The number of imaginary instances for  $x$  is  $|D'| * P'(x, \text{pa}(x))$
    - This is called the *BDe prior*
    - What is  $P'$ ?
      - Could be anything; e.g., a marginally independent distribution

# BAYESIAN ESTIMATION EXAMPLES

- Try a dataset using
  - MLE
  - Bayesian
    - K2
    - BDe

# ICU ALARM NETWORK – FIG 17.C.1



# MISSING DATA

- Note: the following information is from Chapter 19

# MISSING DATA

1. Accidental
  - E.g., sensor failure
2. Intentional
  - E.g., not all tests are ordered for medical diagnosis
3. Hidden variable
  - E.g., cluster assignment; hidden cause

# APPROACHES

1. Ignore the data points with missing values
  - Not economical (throwing data away), not necessarily accurate (if missing intentionally), and might not be possible (hidden variables)
2. Gradient optimization
3. Expectation Maximization



# EXPECTATION MAXIMIZATION (EM)

- Initialize  $\theta$
- Iterate
  - Expectation
    - $M[x, pa(x)] = \sum P(x, pa(x) \mid \text{observed})$  for each  $X$
  - Maximization
    - $\theta_x = M[x, pa(x)] / M[pa(x)]$  for each  $X$

# EM EXAMPLES

- Simple network

- $X \rightarrow Y$
- Let's see a few cases:
  - $X$  is missing in a small percentage; sometimes  $X$  and sometimes  $Y$  is missing;  $X$  is a hidden variable

- A more complicated network:

- A disease ( $D$ ) variable and three test variables ( $T_1, T_2$  and  $T_3$ )
- A naïve Bayes structure:  $D$  is the parent of  $T_1, T_2$  and  $T_3$
- $T_1$  is the fastest but least accurate test.  $T_2$  requires more time but is more accurate.  $T_3$  requires the most time but also is the most accurate test.
- The doctors order  $T_1$  for everyone first. Depending on the results, they might order  $T_2$ , and depending on its results, they might order  $T_3$
- $D$  might be totally hidden or might be observed for only a small subset of the data