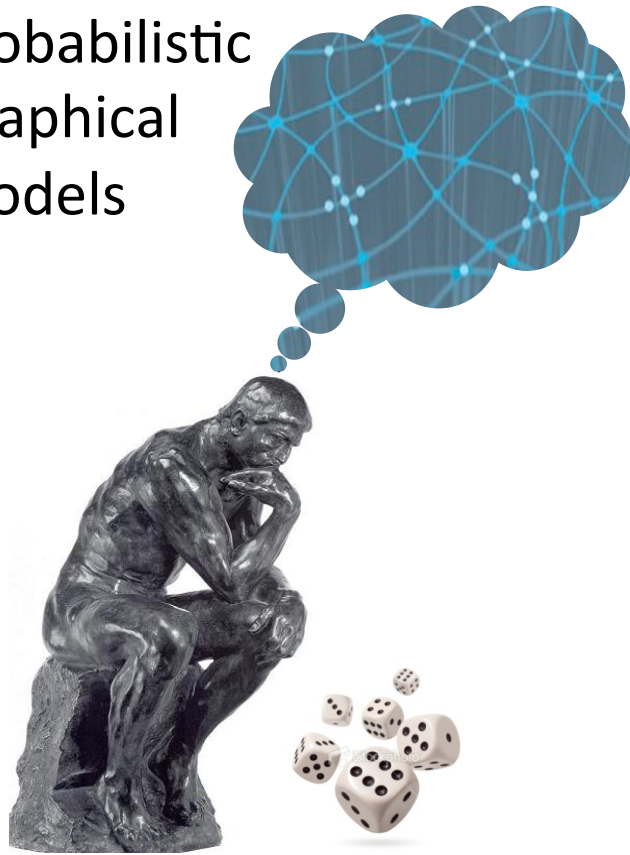


Probabilistic
Graphical
Models



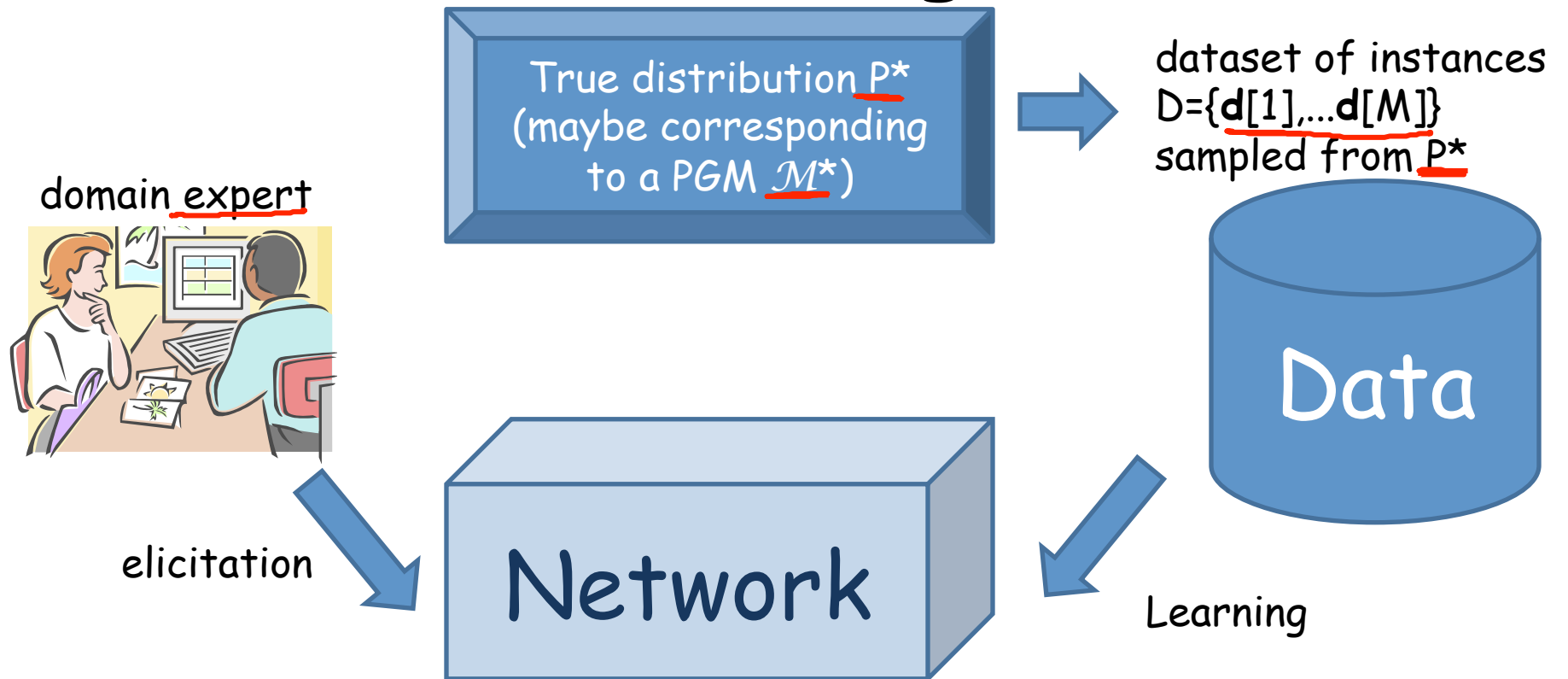
Learning

Overview

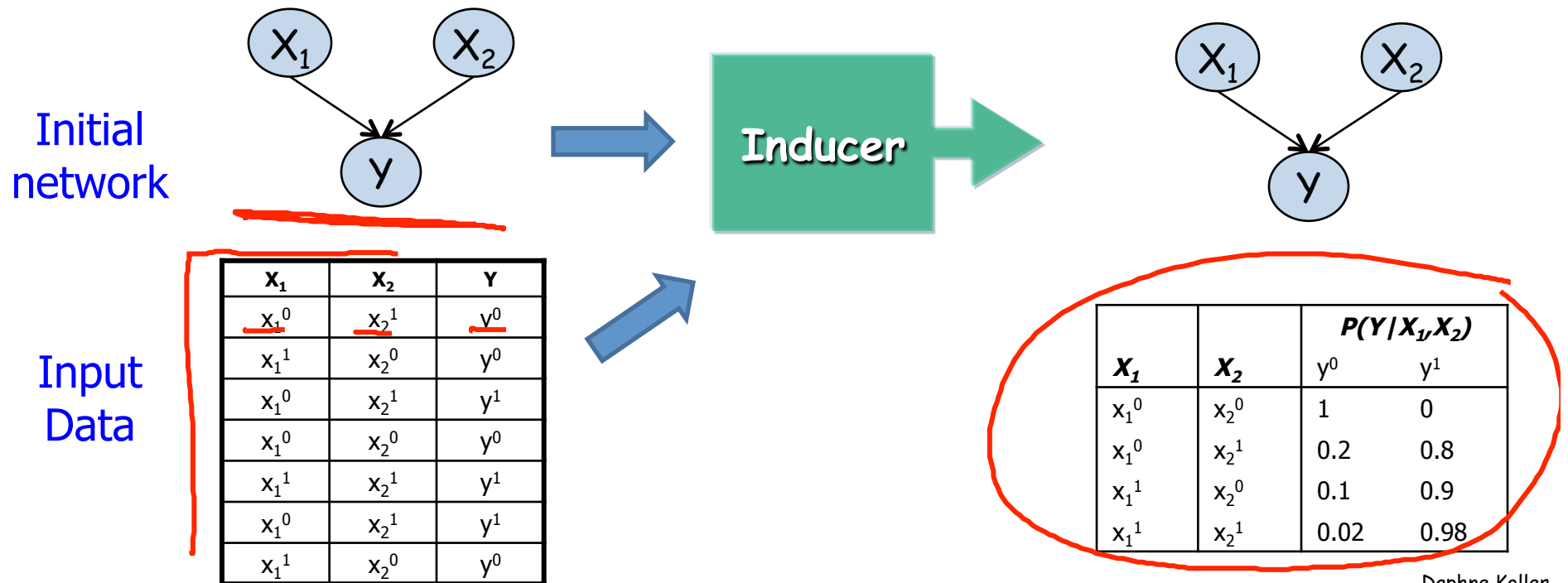
PGM

Learning Tasks
& Metrics

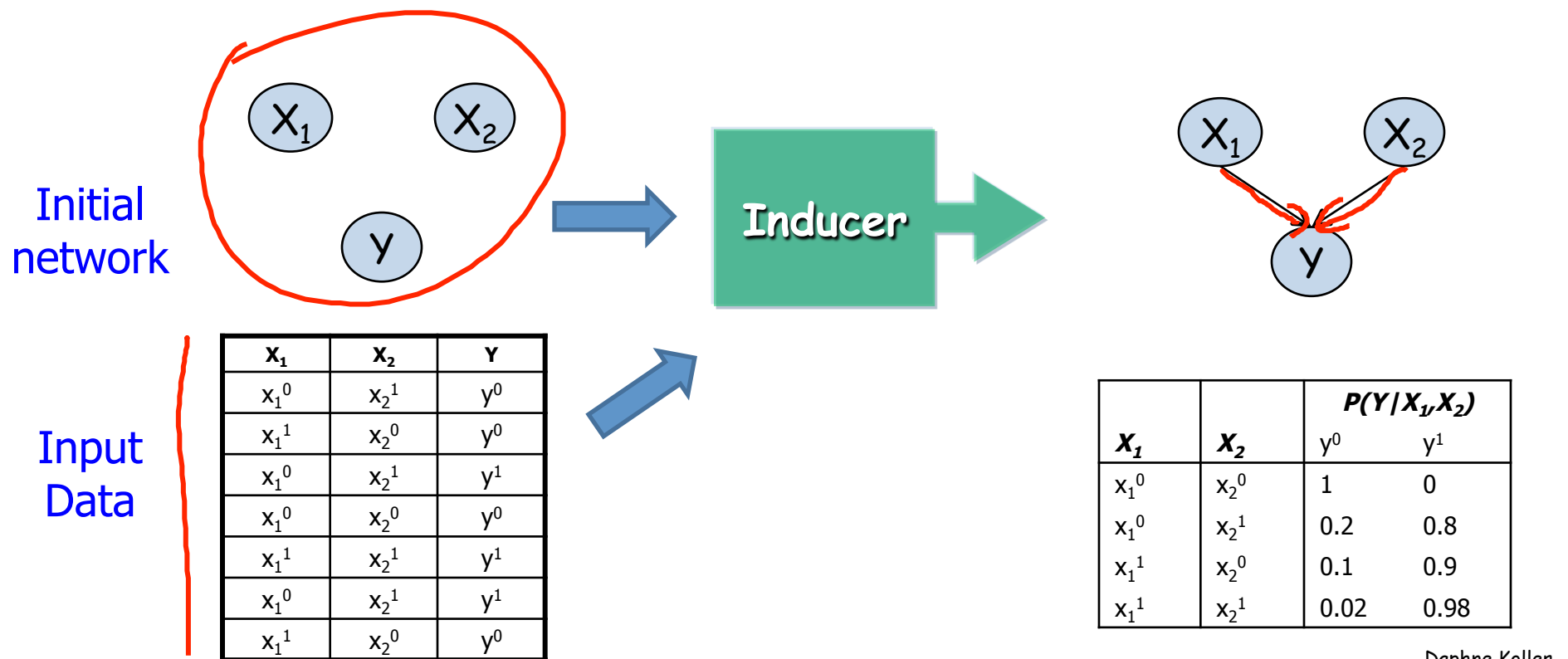
Learning



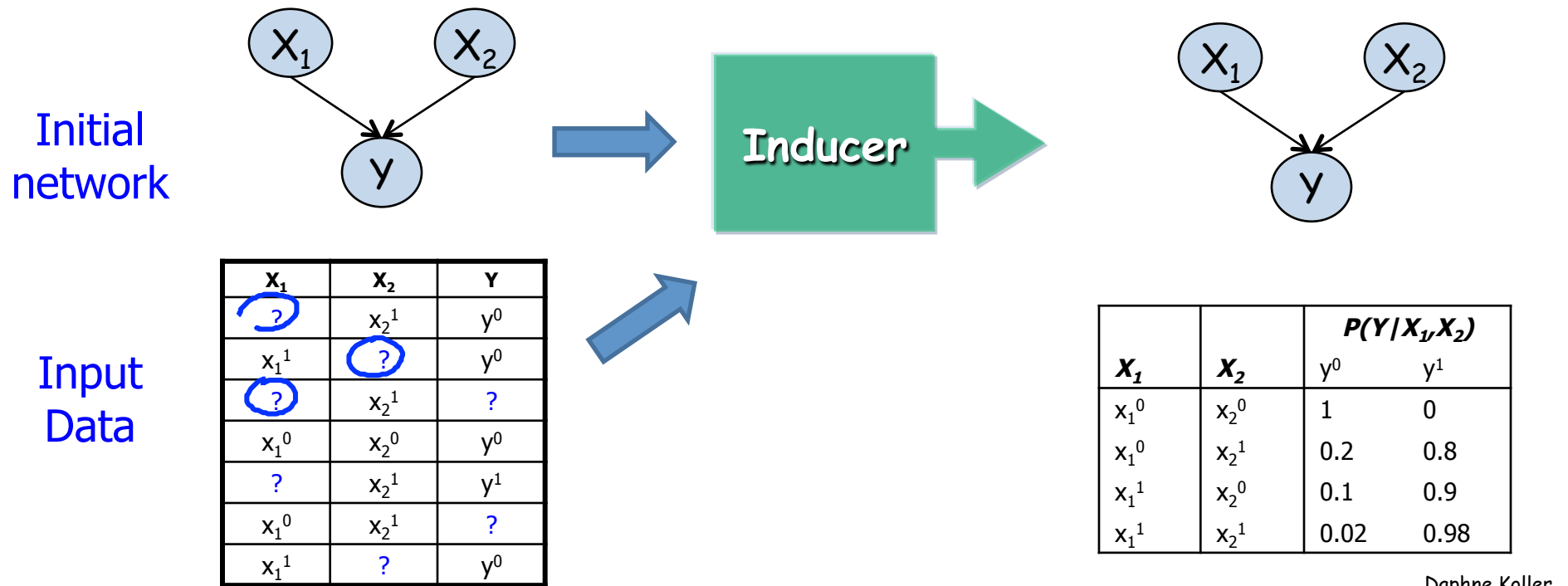
Known Structure, Complete Data



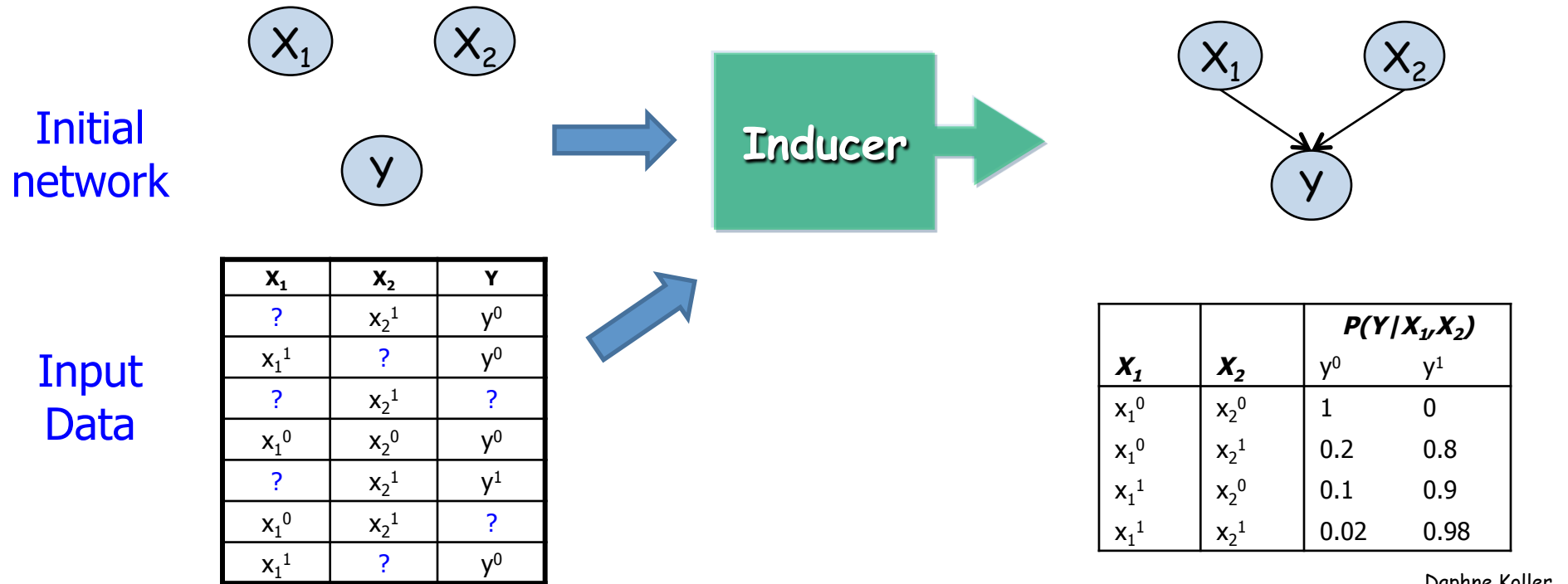
Unknown Structure, Complete Data



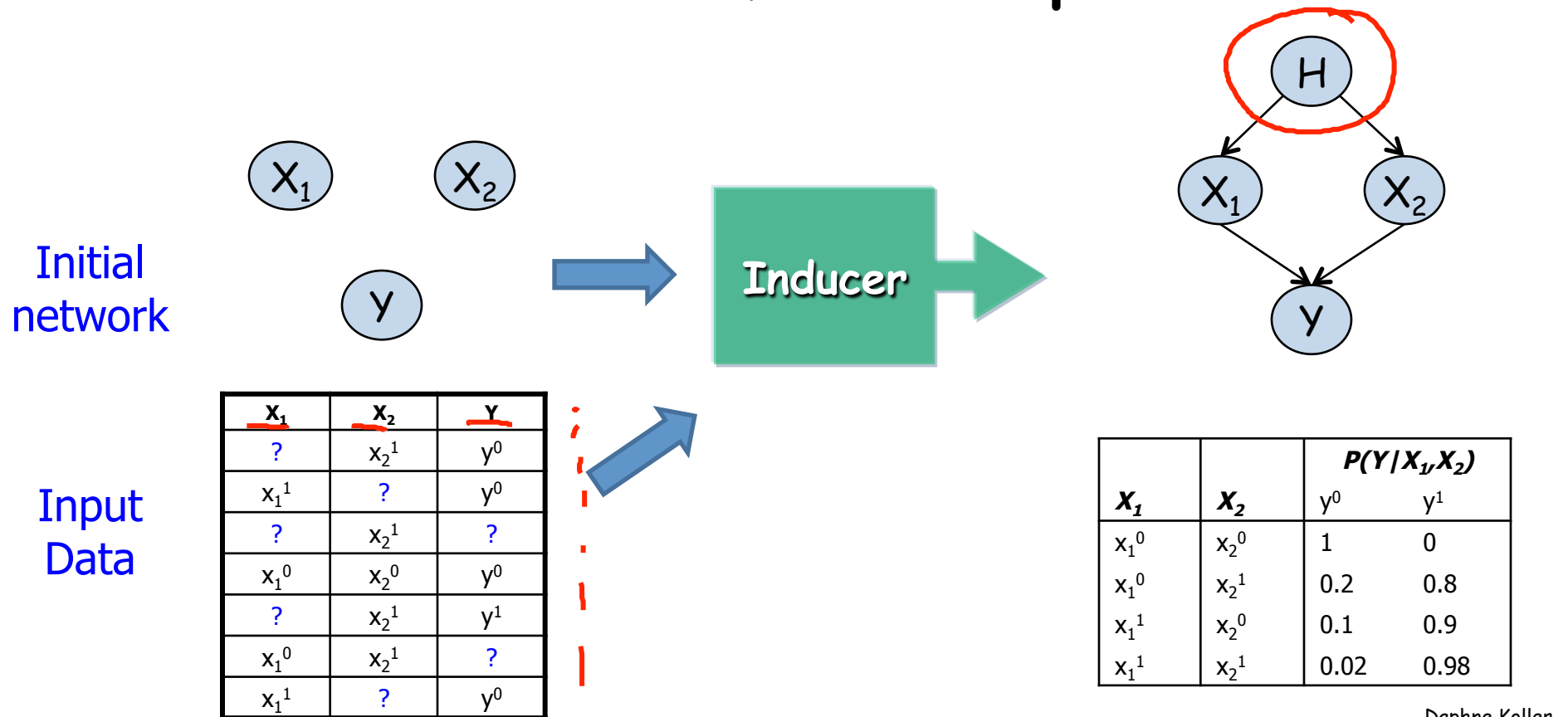
Known Structure, Incomplete Data



Unknown Structure, Incomplete Data



Latent Variables, Incomplete Data



PGM Learning Tasks I

- Goal: Answer general probabilistic queries about new instances
- Simple metric: Training set likelihood
 - $P(\overset{\text{data}}{\underline{D}} : \underline{\mathcal{M}}) = \underline{\prod_m P(\underline{d[m]} : \underline{\mathcal{M}})}$ (IID)
- But we really care about new data
 - Evaluate on test set likelihood - $P(\underline{D'} : \underline{\mathcal{M}})$
generalization performance

PGM Learning Tasks II

- Goal: Specific prediction task on new instances
 - Predict target variables y from observed variables x
 - E.g., image segmentation, speech recognition
- Often care about specialized objective
 - E.g., pixel-level segmentation accuracy
- Often convenient to select model to optimize
 - likelihood $\prod_m P(\mathbf{d}[m] : \mathcal{M})$ or
 - conditional likelihood $\prod_m P(\mathbf{y}[m] \mid \mathbf{x}[m] : \mathcal{M})$
- Model evaluated on "true" objective over test data

PGM Learning Tasks III



- Goal: Knowledge discovery of \mathcal{M}^*
 - Distinguish direct vs indirect dependencies
 - Possibly directionality of edges
 - Presence and location of hidden variables
- Often train using likelihood
 - Poor surrogate for structural accuracy
- Evaluate by comparing to prior knowledge

Avoiding Overfitting

- Selecting \mathcal{M} to optimize training set likelihood overfits to statistical noise
- Parameter overfitting
 - Parameters fit random noise in training data
 - Use regularization / parameter priors
- Structure overfitting
 - Training likelihood always increases for more complex structures
 - Bound or penalize model complexity

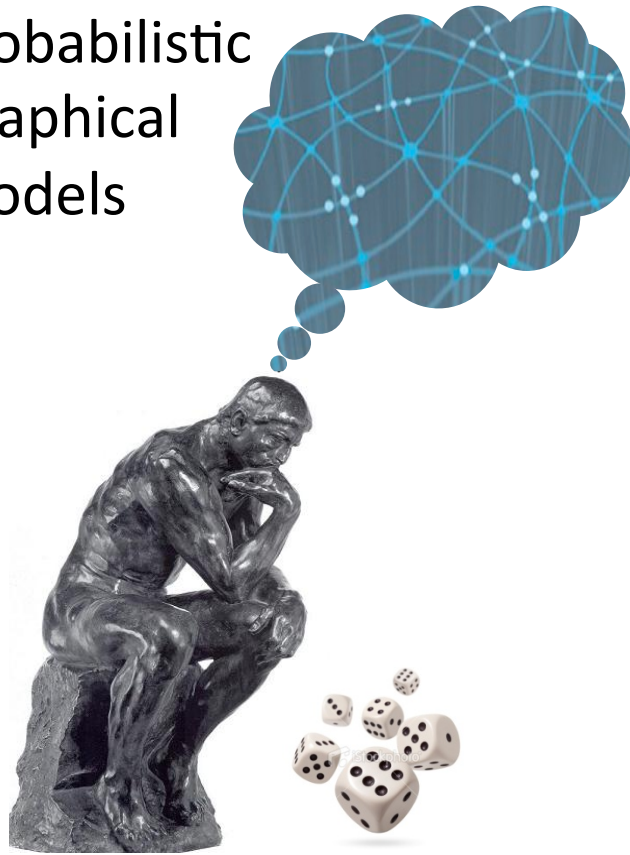
Selecting Hyperparameters

- Regularization for overfitting involves hyperparameters:
 - Parameter priors (regularization)
 - Complexity penalty
- Choice of hyperparameters makes a big difference to performance
- Must be selected on validation set ~~training set~~ ~~test~~ (cross-validation)

Why PGM Learning

- Predictions of structured objects
(sequences, graphs, trees)
 - Exploit correlations between several predicted variables
- Can incorporate prior knowledge into model
- Learning single model for multiple tasks
- Framework for knowledge discovery

Probabilistic
Graphical
Models



Learning

Parameter Estimation

Maximum
Likelihood
Estimation

Biased Coin Example

P is a Bernoulli distribution:

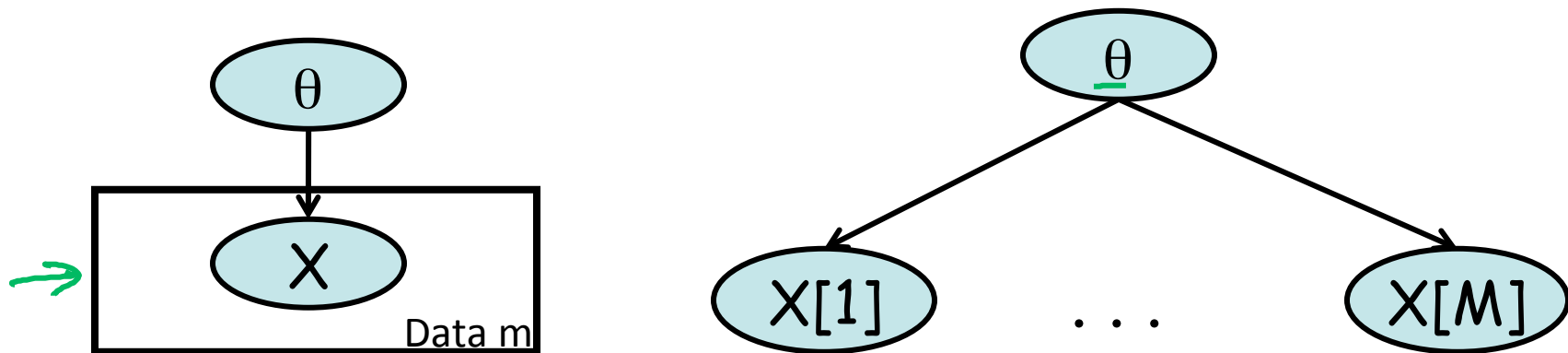
$$P(\underline{X=1}) = \underline{\theta}, P(X=0) = \underline{1-\theta}$$



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

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

IID as a PGM



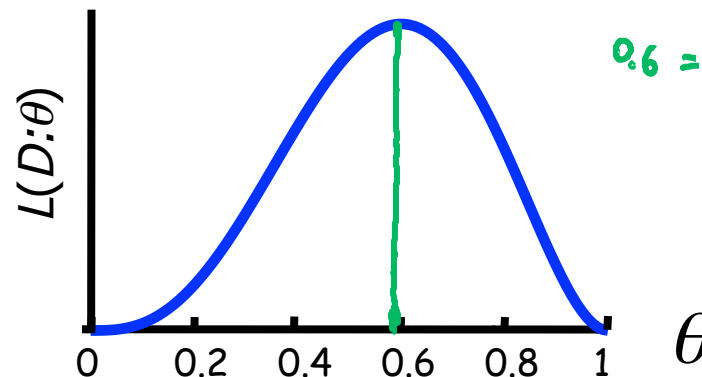
$$\underline{P(x[m] | \theta)} = \begin{cases} \underline{\theta} & x[m] = \underline{x^1} \\ \underline{1 - \theta} & x[m] = \underline{x^0} \end{cases}$$

Maximum Likelihood Estimation

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

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

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



$$0.6 = \frac{3}{5}$$

Maximum Likelihood Estimator

- Observations: M_H heads and M_T tails
- Find θ maximizing likelihood

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

- 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 θ :

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

Sufficient Statistics

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

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

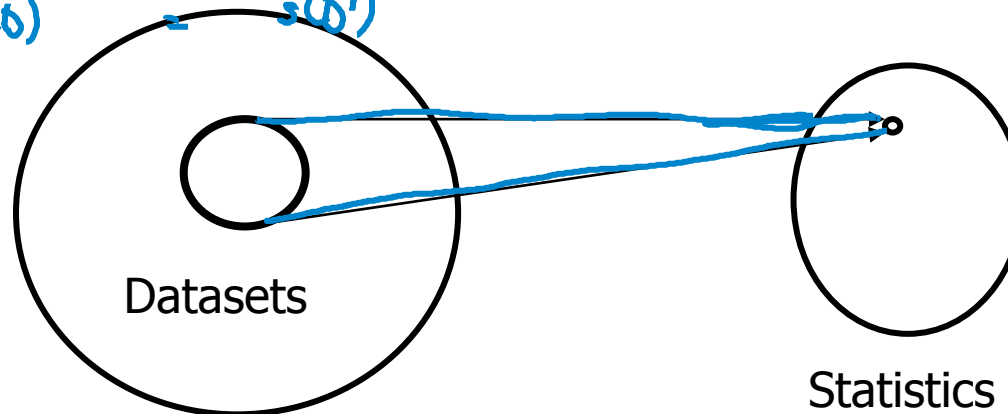
- $\rightarrow M_H$ and M_T are sufficient statistics

Sufficient Statistics

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

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

(Handwritten blue annotations: $s(D)$ under the first sum, $s(D')$ under the second sum, and a blue arrow pointing from the equality to the likelihood equation.)



Sufficient Statistic for Multinomial

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

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

Sufficient Statistic for Gaussian

- Gaussian distribution:

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

- Rewrite as

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

- 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 θ to maximize $L(D;\Theta)$

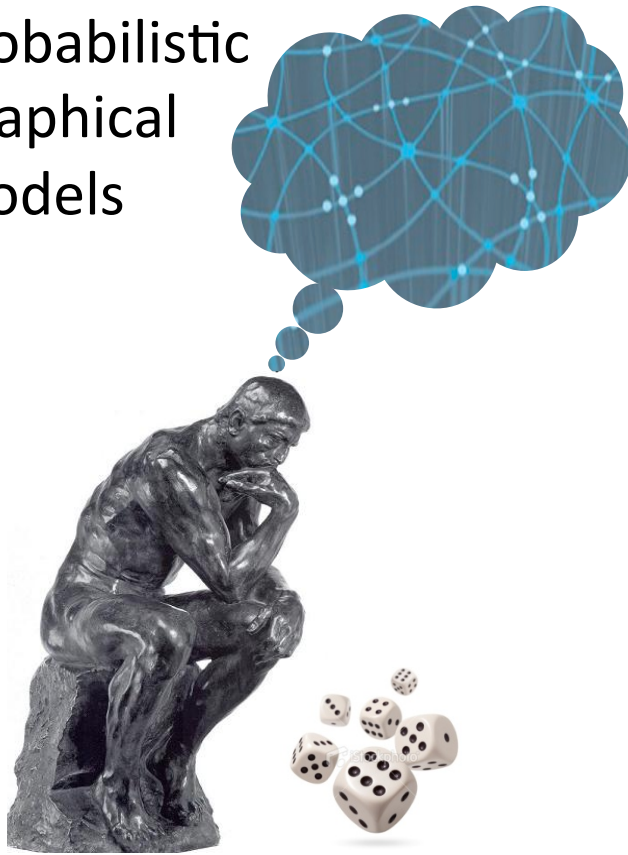
- Multinomial MLE: $\hat{\theta}^i = \frac{M_i}{\sum_{i=1}^m M_i}$ *fraction of x^i in data*

- 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

Probabilistic
Graphical
Models



Learning

Parameter Estimation

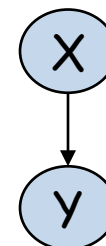
Max Likelihood
for BNs

MLE for Bayesian Networks

- Parameters: $\rightarrow \theta_{x^0}, \theta_{x^1}$
 $\theta_{y^0|x^0}, \theta_{y^1|x^0}, \theta_{y^0|x^1}, \theta_{y^1|x^1}$
- Data instances: $\langle x[m], y[m] \rangle$

x	
x^0	x^1
0.7	0.3

$P(x)$

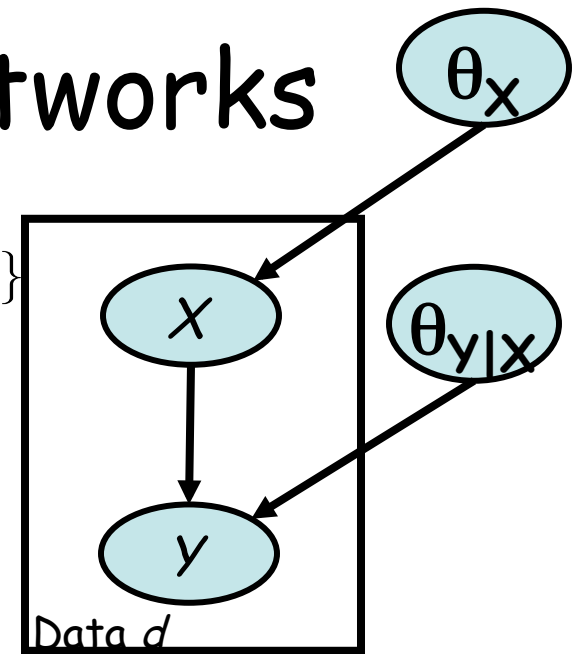


x	y	
	y^0	y^1
x^0	0.95	0.05
x^1	0.2	0.8

$P(y|x)$

MLE for Bayesian Networks

- Parameters: $\{\theta_x : x \in \text{Val}(X)\}$
 $\{\theta_{y|x} : x \in \text{Val}(X), y \in \text{Val}(Y)\}$



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

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

// chain rule for BNs

$$= \left(\prod_{m=1}^M P(x[m] : \theta) \right) \left(\prod_{m=1}^M P(y[m] | x[m] : \theta) \right)$$

$$= \left(\prod_{m=1}^M P(x[m] : \theta_x) \right) \left(\prod_{m=1}^M P(y[m] | x[m] : \theta_{y|x}) \right)$$

product of two
local likelihood

MLE for Bayesian Networks

- Likelihood for Bayesian network

$$\begin{aligned}
 L(\Theta : D) &= \prod_m P(x[m] : \Theta) \\
 &= \prod_m \prod_i P(x_i[m] | \underline{U_i[m]} : \Theta_i) \\
 &= \prod_i \prod_m P(x_i[m] | \underline{U_i[m]} : \Theta_i) \\
 &\stackrel{\text{local likelihood}}{\rightarrow} \prod_i L_i(D : \Theta_i)
 \end{aligned}$$

parents of X_i (points to U_i)
chain rule (points to the first two steps)
 $L_i(\Theta_i : D)$ (points to the final term)

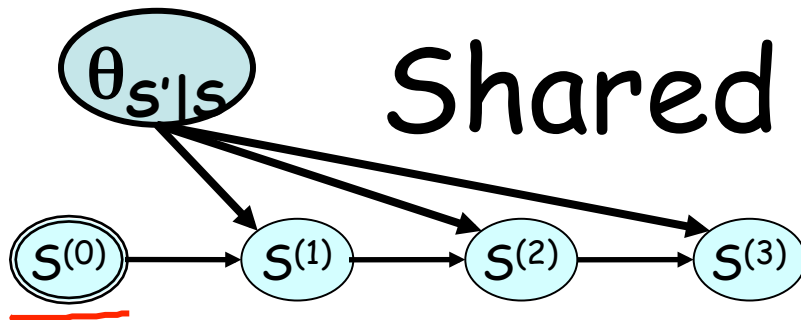
\Rightarrow if $\theta_{X_i|U_i}$ are disjoint, then MLE can be computed by maximizing each local likelihood separately

MLE for Table CPDs

$$\begin{aligned}
 \prod_{m=1}^M P(x[m] \mid \mathbf{u}[m] : \theta) &= \prod_{m=1}^M P(x[m] \mid \mathbf{u}[m] : \theta_{X|U}) \\
 &= \prod_{x, \mathbf{u}} \left(\prod_{m: x[m]=x, \mathbf{u}[m]=\mathbf{u}} P(x[m] \mid \mathbf{u}[m] : \theta_{X|U}) \right) \\
 &= \prod_{x, \mathbf{u}} \left(\prod_{m: x[m]=x, \mathbf{u}[m]=\mathbf{u}} \theta_{x|\mathbf{u}} \right) \\
 &= \prod_{x, \mathbf{u}} \frac{\theta_{x|\mathbf{u}}^{M[x, \mathbf{u}]}}{p(\mathbf{u})}
 \end{aligned}$$

$P(x[m]=x \mid \mathbf{u}[m]=\mathbf{u} ; \theta_{X|U}) = \theta_{x|\mathbf{u}}$
 fraction of $X=x$ among cases where $\mathbf{u}=\mathbf{u}$
 $\theta_{x|\mathbf{u}} = \frac{M[x, \mathbf{u}]}{\sum_{x'} M[x', \mathbf{u}]} = \frac{M[x, \mathbf{u}]}{M[\mathbf{u}]}$

Shared Parameters



$$L(\theta : S^{(0:T)}) = \prod_{t=1}^T P(S^{(t)} \mid S^{(t-1)} : \theta)$$

$$= \prod_{i,j} \prod_{t: S^{(t)}=s^i, S^{(t+1)}=s^j} P(S^{(t+1)} \mid S^{(t)} : \theta_{s^i \rightarrow s^j})$$

$$= \prod_{i,j} \prod_{t: S^{(t)}=s^i, S^{(t+1)}=s^j} \theta_{s^i \rightarrow s^j}$$

$$= \prod_{i,j} \theta_{s^i \rightarrow s^j}^{M[s^i \rightarrow s^j]}$$

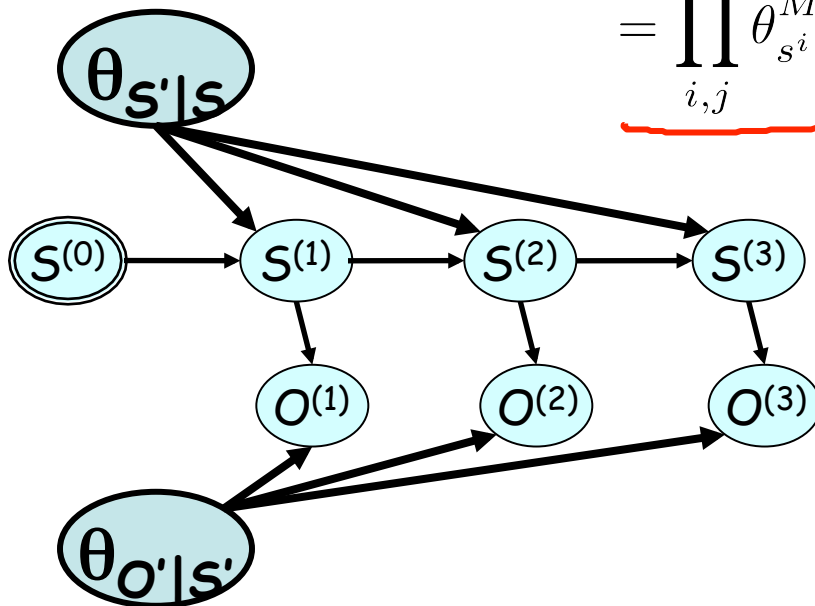
$$\hat{\theta}_{s^i \rightarrow s^j} = \frac{M[s^i \rightarrow s^j]}{M[s^i]}$$

$$M[s^i \rightarrow s^j] = |\{t : S^{(t)} = s^i, S^{(t+1)} = s^j\}|$$

Shared Parameters

$$L(\Theta : S^{(0:T)}, O^{(0:T)}) = \prod_{t=1}^T P(S^{(t)} | S^{(t-1)} : \theta_{S'|S}) \prod_{t=1}^T P(O^{(t)} | S^{(t)} : \theta_{O'|S'})$$

$$= \prod_{i,j} \theta_{S^i \rightarrow S^j}^{M[s^i \rightarrow s^j]} \prod_{i,k} \theta_{O^k | S^i}^{M[o^k, s^i]}$$



$$M[s^i \rightarrow s^j] = |\{t : S^{(t)} = s^i, S^{(t+1)} = s^j\}|$$

$$M[o^k, s^i] = |\{t : S^{(t)} = s^i, O^{(t)} = o^k\}|$$

Summary

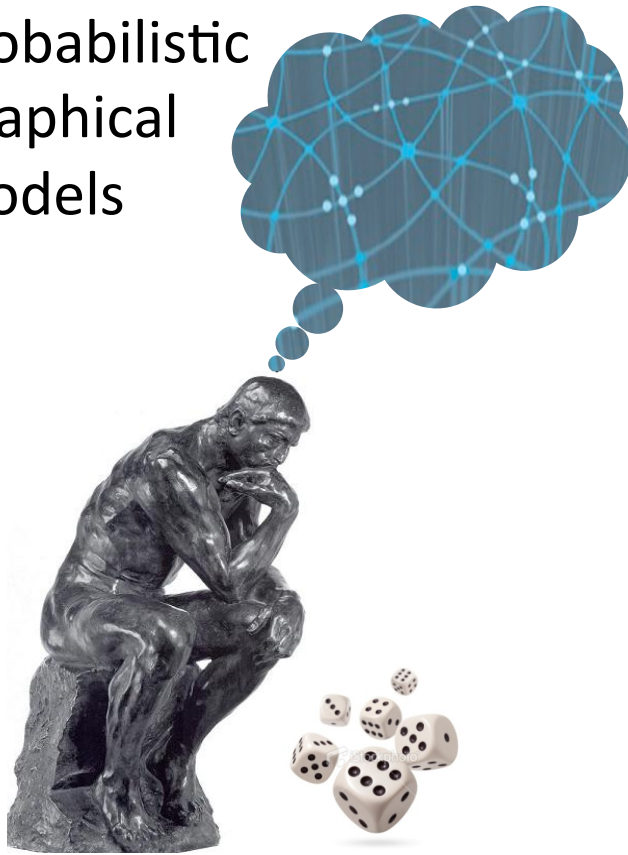
- For BN with disjoint sets of parameters in CPDs, likelihood decomposes as product of local likelihood functions, one per variable
- For table CPDs, local likelihood further decomposes as product of likelihood for multinomials, one for each parent combination
- For networks with shared CPDs, sufficient statistics accumulate over all uses of CPD

Fragmentation & Overfitting

$$\theta_{x|u} = \frac{M[x, u]}{\sum_{x'} M[x', u]} = \frac{M[x, u]}{M[u]}$$

- # of "buckets" increases exponentially with $|U|$
- For large $|U|$, most "buckets" will have very few instances
 \Rightarrow very poor parameter estimates \Leftarrow
- With limited data, we often get better generalization with simpler structures
even when wrong

Probabilistic
Graphical
Models



Learning

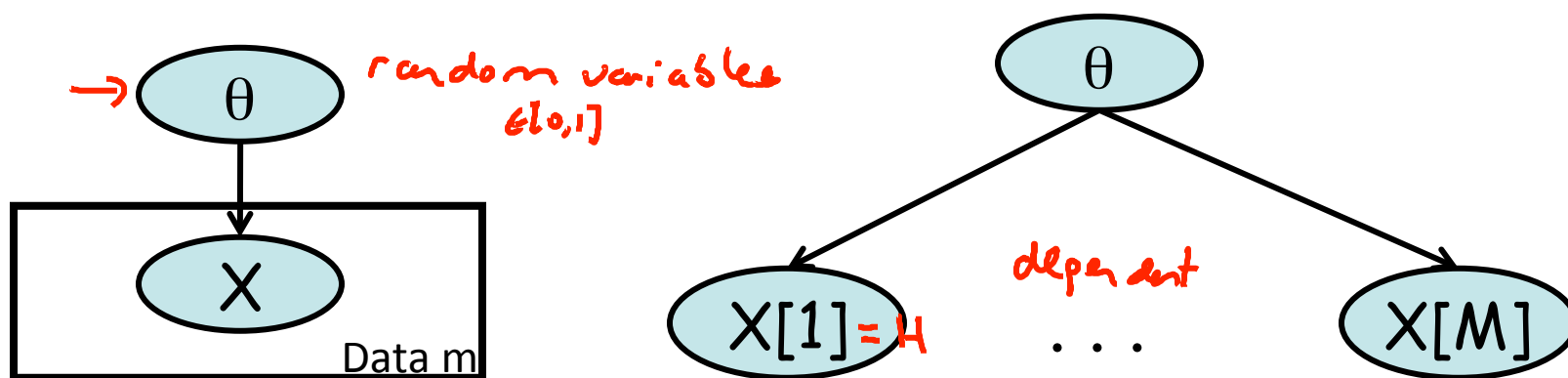
Parameter Estimation

Bayesian
Estimation

Limitations of MLE

- Two teams play 10 times, and the first wins 7 of the 10 matches
 - ⇒ Probability of first team winning = 0.7
- A coin is tossed 10 times, and comes out 'heads' 7 of the 10 tosses
 - ⇒ Probability of heads = 0.7
- A coin is tossed 10000 times, and comes out 'heads' 7000 of the 10000 tosses
 - ⇒ Probability of heads = 0.7

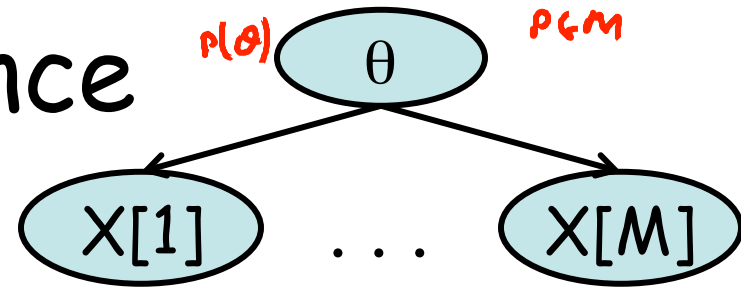
Parameter Estimation as a PGM



- Given a fixed θ , tosses are independent
- If θ is unknown, tosses are not marginally independent
 - each toss tells us something about θ

Bayesian Inference

- Joint probabilistic model



$$\underline{P(x[1], \dots, x[M], \theta)} = \underline{P(x[1], \dots, x[M] | \theta)} \underline{P(\theta)}$$

$$= P(\theta) \prod_{i=1}^M P(x[i] | \theta)$$

$$= \underline{P(\theta) \theta^{M_H} (1 - \theta)^{M_T}}$$

likelihood function

$$\underline{P(\theta | x[1], \dots, x[M])} = \frac{\overset{\text{likelihood}}{\underline{P(x[1], \dots, x[M] | \theta)}} \overset{\text{prior}}{\underline{P(\theta)}}}{\underset{\substack{\text{data } D \\ \text{constant relative to } \theta}}{P(x[1], \dots, x[M])}}$$

Dirichlet Distribution

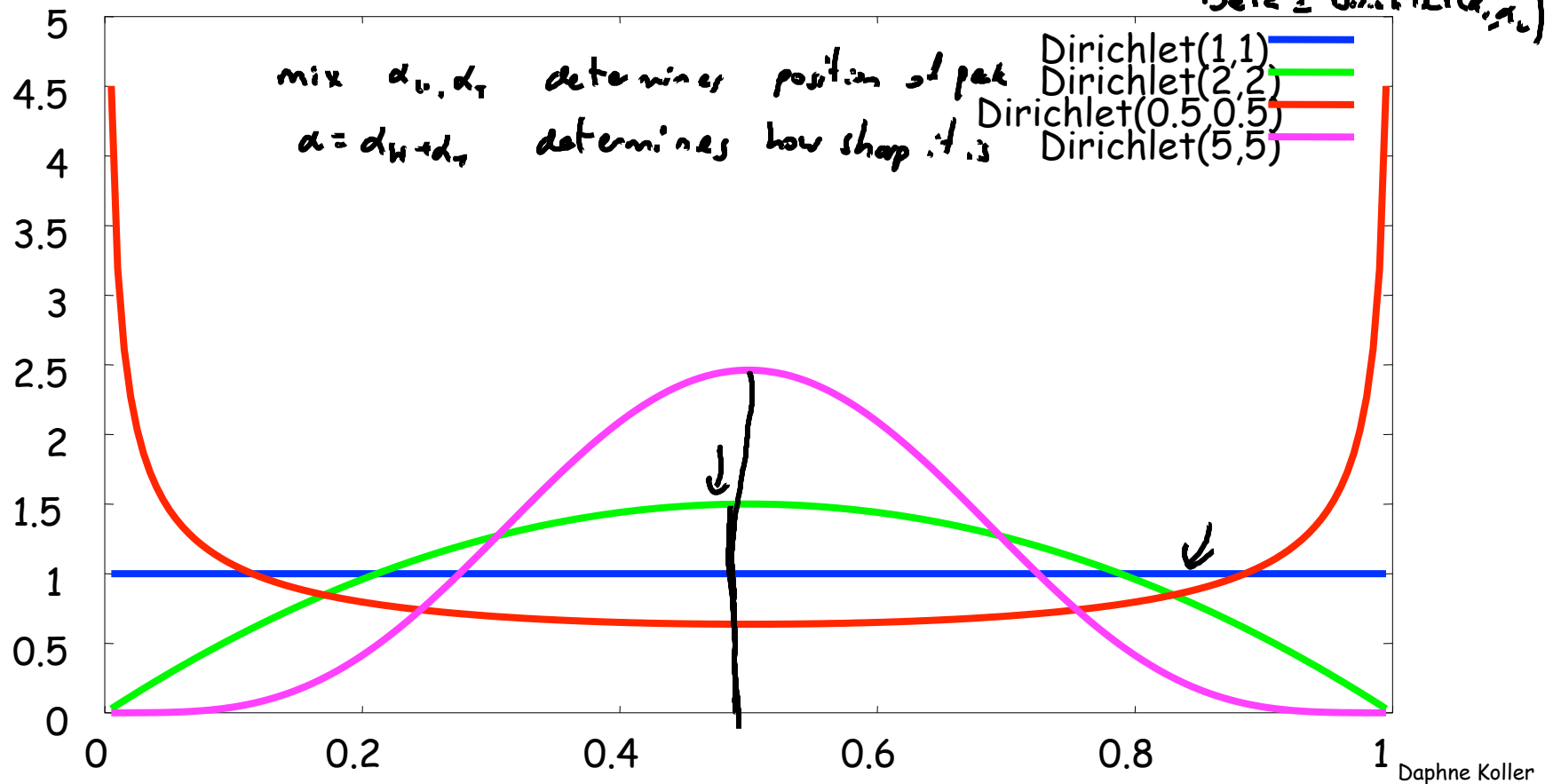
- θ is a multinomial distribution over k values
- Dirichlet distribution $\theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$

– where $P(\theta) = \frac{1}{Z} \prod_{i=1}^k \theta_i^{\alpha_i - 1}$ and $Z = \frac{\prod_{i=1}^k \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^k \alpha_i)}$ $\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$

hyperparameters

- Intuitively, hyperparameters α_i correspond to the number of samples we have seen

Dirichlet Distributions



Dirichlet Priors & Posteriors

$$\overbrace{P(\theta | D)}^{\text{posterior}} \propto \overbrace{P(D | \theta)}^{\text{likelihood}} \overbrace{P(\theta)}^{\text{prior}}$$

$$d_{x_i} = \# \text{ instances with } x_i \quad \overbrace{P(D | \theta)}^{\text{multinomial } \theta} = \prod_{i=1}^k \underbrace{\theta_i^{M_i}}_{\theta_i^{n_i + \alpha_i - 1}} \quad P(\theta) \propto \prod_{i=1}^k \underbrace{\theta_i^{\alpha_i - 1}}$$

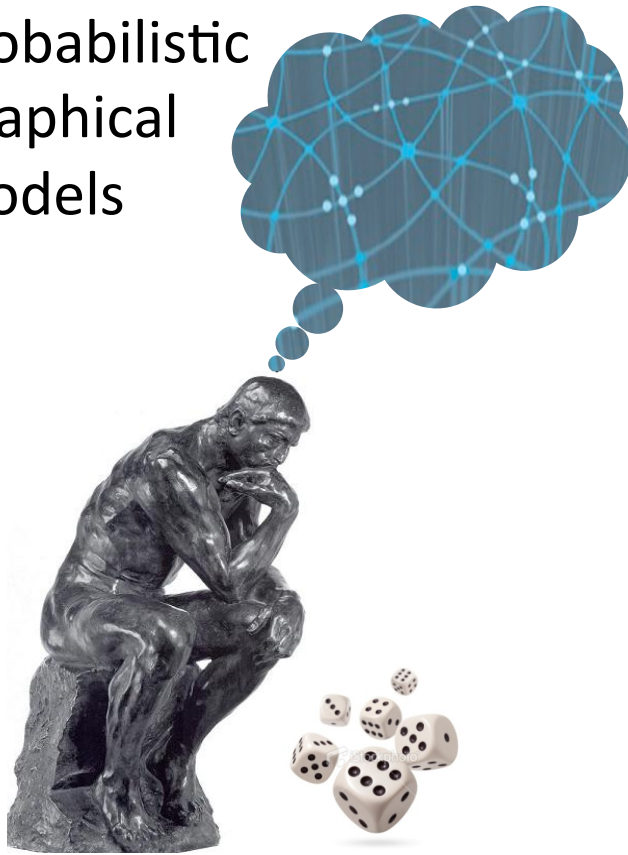
- If $P(\theta)$ is Dirichlet and the likelihood is multinomial, then the posterior is also Dirichlet
 - Prior is $\text{Dir}(\alpha_1, \dots, \alpha_k)$
 - Data counts are M_1, \dots, M_k
 - Posterior is $\text{Dir}(\alpha_1 + M_1, \dots, \alpha_k + M_k)$
- Dirichlet is a conjugate prior for the multinomial

prior, posterior have the same form

Summary

- Bayesian learning treats parameters as random variables
 - Learning is then a special case of inference
- Dirichlet distribution is conjugate to multinomial
 - Posterior has same form as prior
 - Can be updated in closed form using sufficient statistics from data

Probabilistic
Graphical
Models

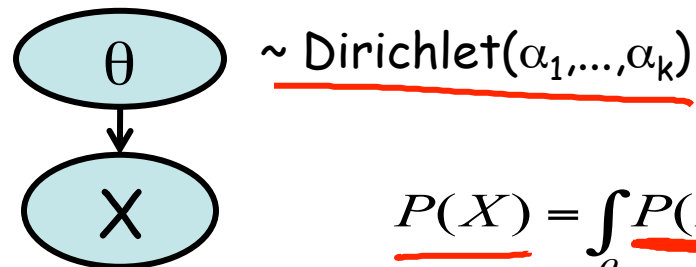


Learning

Parameter Estimation

Bayesian
Prediction

Bayesian Prediction



$$\underline{P(X)} = \int_{\theta} \underline{P(X | \theta)} \underline{P(\theta)} d\theta$$

← marginalizing over θ

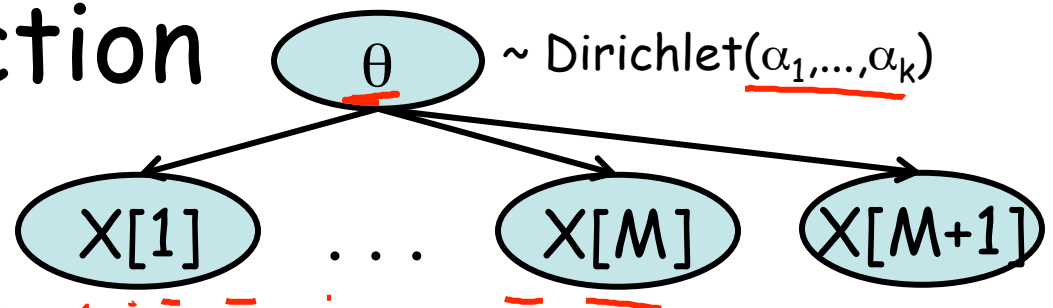
$$\underline{P(X = \underline{x^i} | \theta)} = \frac{1}{Z} \int_{\theta} \theta^{\alpha_i} \cdot \underbrace{\prod_j \theta^{\alpha_j - 1}}_{\text{prior}} d\theta$$

$$= \frac{\alpha_i}{\underline{\sum_j \alpha_j} + \alpha}$$

fraction of instances we've seen where x^i

- Dirichlet hyperparameters correspond to the number of samples we have seen

Bayesian Prediction



$$P(\underline{x[M+1]} | \underline{x[1]}, \dots, \underline{x[M]})$$

$$= \int_{\theta} P(\underline{x[M+1]} | \underline{x[1]}, \dots, \underline{x[M]}, \theta) P(\theta | \underline{x[1]}, \dots, \underline{x[M]}) d\theta$$

$$= \int_{\theta} \underline{P(x[M+1] | \theta)} \underline{P(\theta | x[1], \dots, x[M])} d\theta$$

~ Dirichlet($\alpha_1 + M_1, \dots, \alpha_k + M_k$)
Posterior θ given 0

$$P(X[M+1] = \underline{x^i} | \theta, x[1], \dots, x[M]) = \frac{\alpha_i + M_i}{\underline{\alpha + M}}$$

$\alpha = \sum \alpha_i$
 $M = \sum M_i$

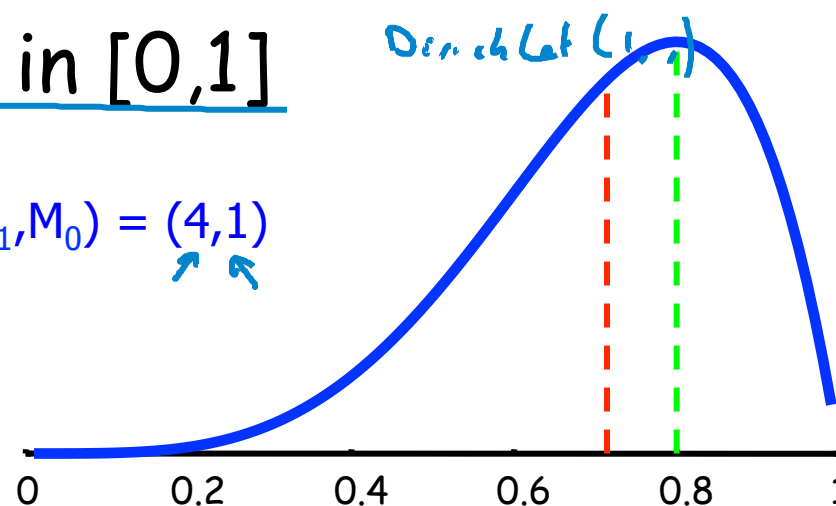
- Equivalent sample size $\alpha = \alpha_1 + \dots + \alpha_k$
 - Larger $\alpha \Rightarrow$ more confidence in our prior

Example: Binomial Data

- Prior: uniform for θ in $[0,1]$

$$P(\theta) = \frac{1}{Z} \prod_k \theta^{\alpha_k - 1}$$

$$(M_1, M_0) = (4, 1)$$

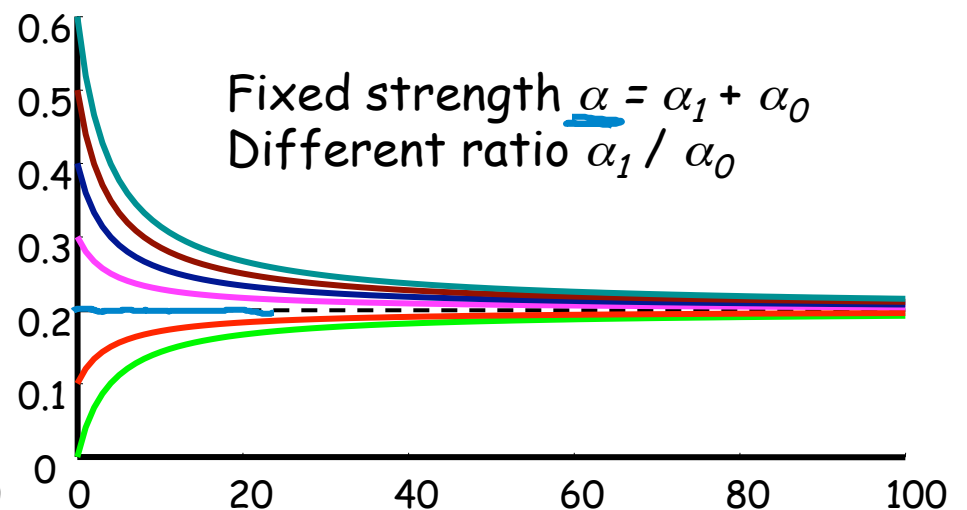
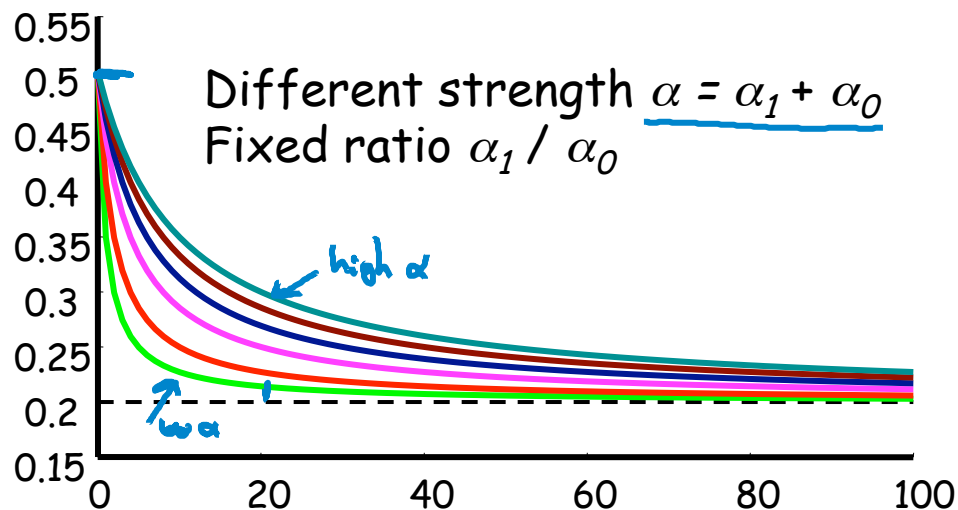


- MLE for $P(X[6]=1)=4/5$
- Bayesian prediction is $5/7$

$$\frac{\alpha_1 + M_1}{\alpha + M} = \frac{1 + 4}{2 + 5}$$

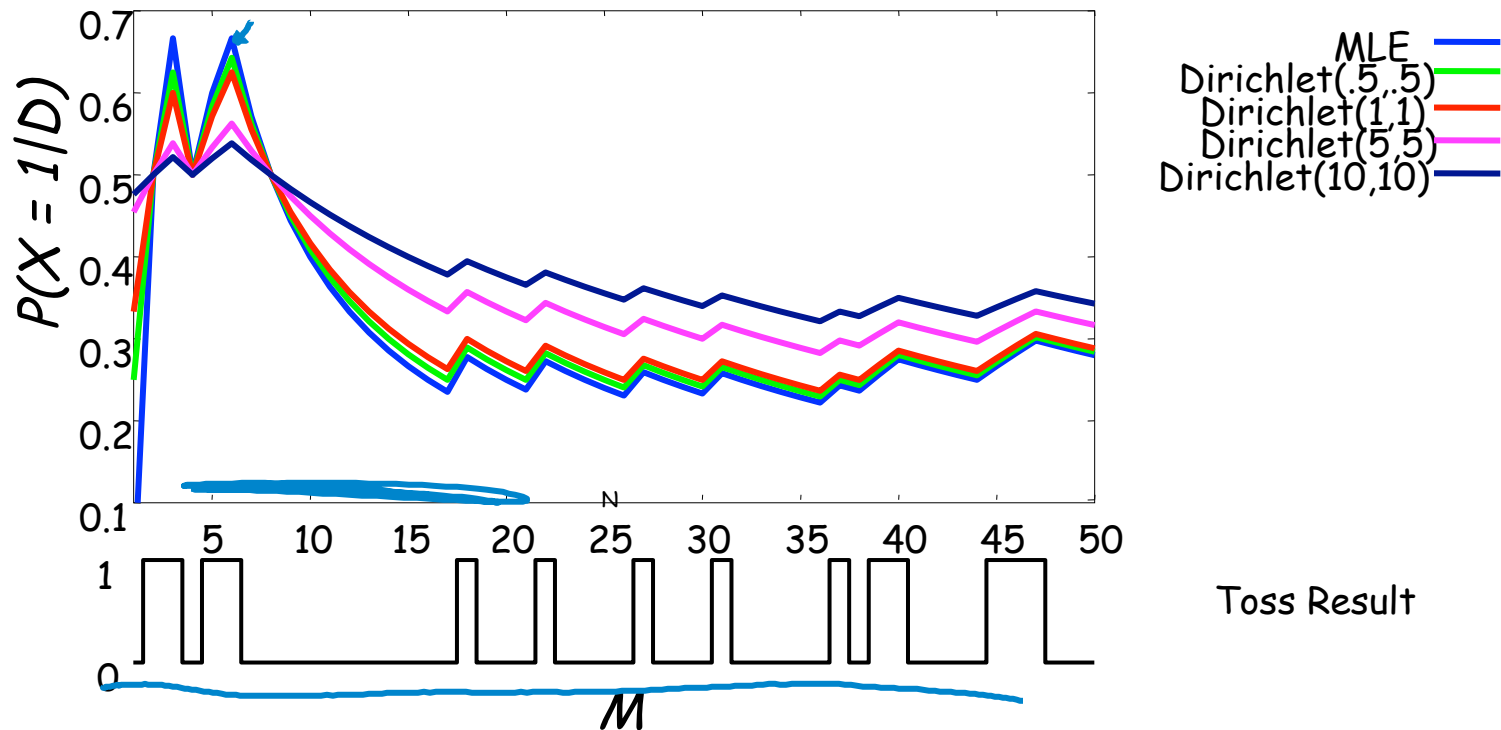
Effect of Priors

- Prediction of $P(X=1)$ after seeing data with $M_1 = \frac{1}{4}M_0$ as a function of sample size M



Effect of Priors

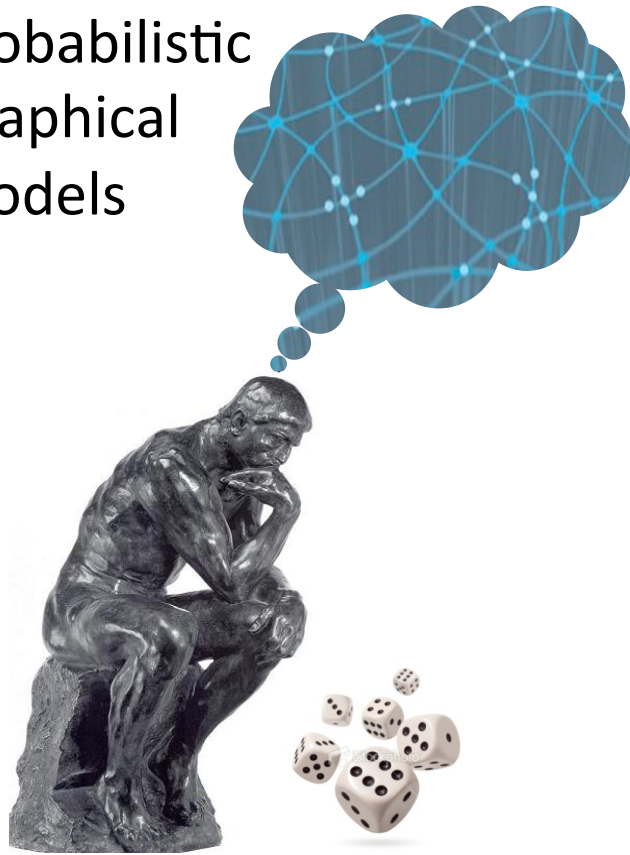
- In real data, Bayesian estimates are less sensitive to noise in the data



Summary

- Bayesian prediction combines sufficient statistics from imaginary Dirichlet samples and real data samples
- Asymptotically the same as MLE
- But Dirichlet hyperparameters determine both the prior beliefs and their strength

Probabilistic
Graphical
Models

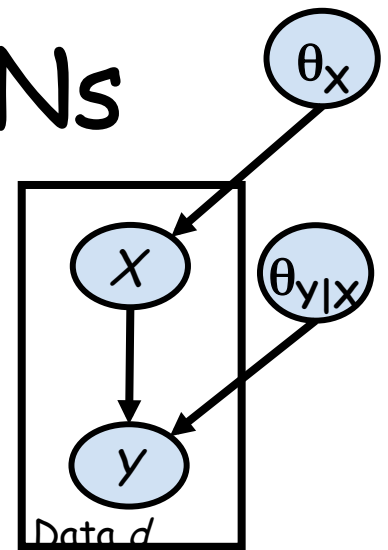
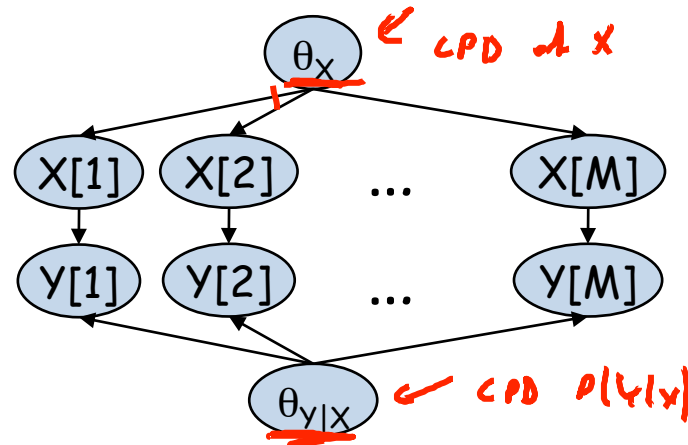


Learning

Parameter Estimation

Bayesian
Estimation
for BNs

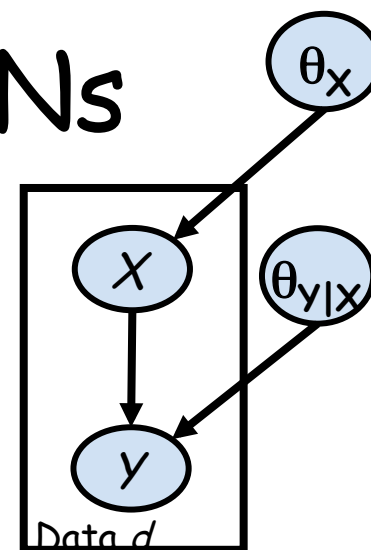
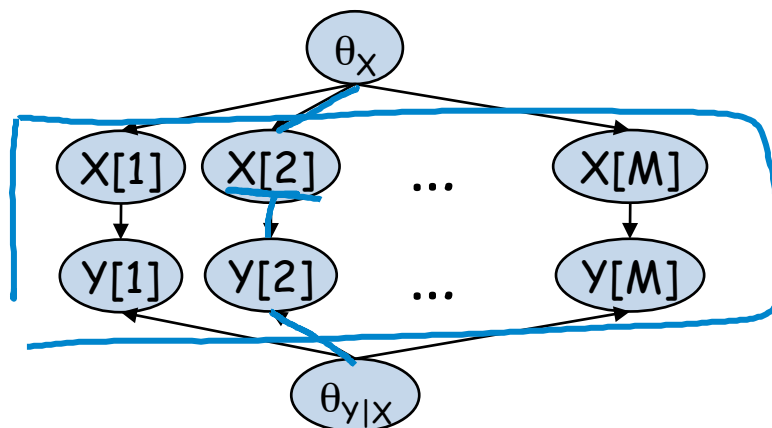
Bayesian Estimation in BNs



- Instances are independent given the parameters
 - $(X[m'], Y[m'])$ are d-separated from $(X[m], Y[m])$ given θ
- Parameters for individual variables are independent a priori

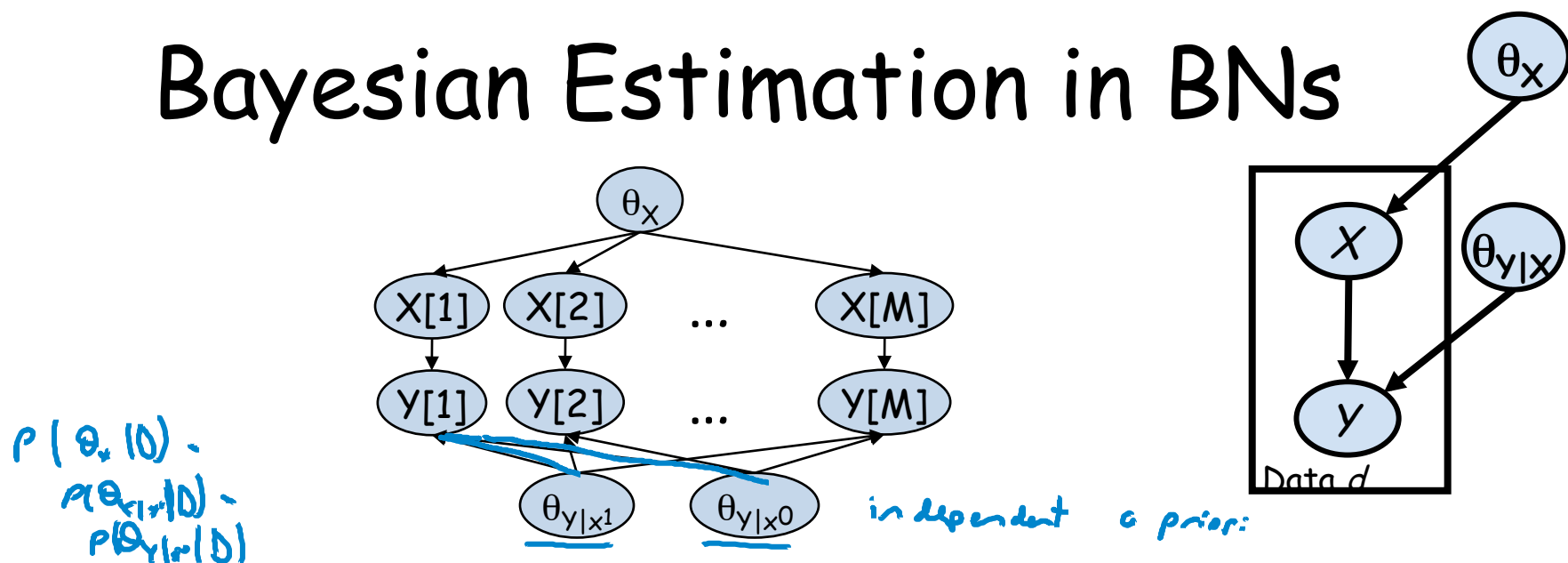
$$P(\theta) = \prod_i P(\theta_{X_i | Pa(X_i)})$$

Bayesian Estimation in BNs



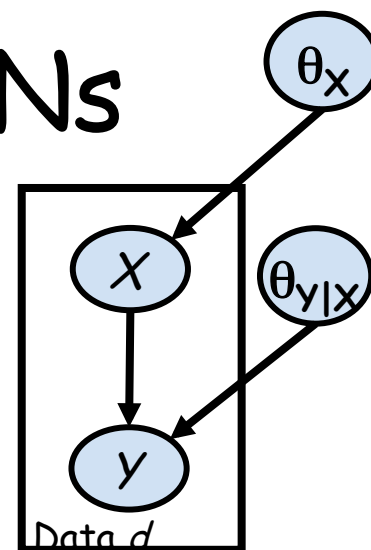
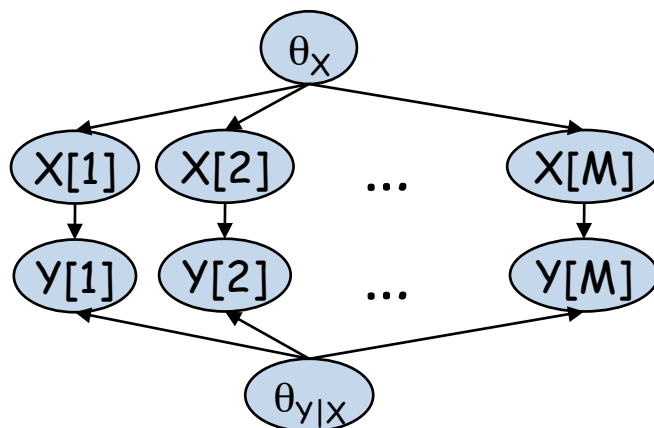
- Posteriors of θ are independent given complete data
 - Complete data d-separates parameters for different CPDs
 - $P(\theta_x, \theta_{y|x} \mid D) = P(\theta_x \mid D)P(\theta_{y|x} \mid D)$
 - As in MLE, we can solve each estimation problem separately

Bayesian Estimation in BNs



- Posteriors of θ are independent given complete data
 - Also holds for parameters within families
 - Note **context specific independence** between $\theta_{y|x1}$ and $\theta_{y|x0}$ when given both X 's and Y 's

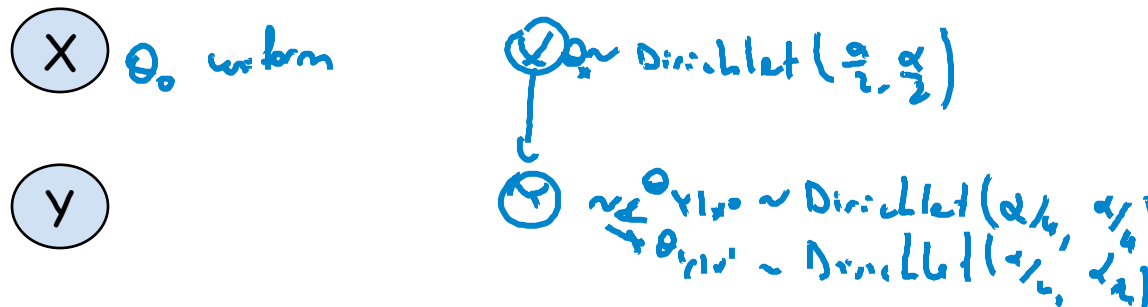
Bayesian Estimation in BNs



- Posteriors of θ can be computed independently
 - For multinomial $\theta_{x|u}$ ^{assigned to x 's parents u} if prior is $\text{Dirichlet}(\alpha_{x^1|u}, \dots, \alpha_{x^k|u})$
 - posterior is $\text{Dirichlet}(\alpha_{x^1|u} + M[x^1, u], \dots, \alpha_{x^k|u} + M[x^k, u])$

Assessing Priors for BNs

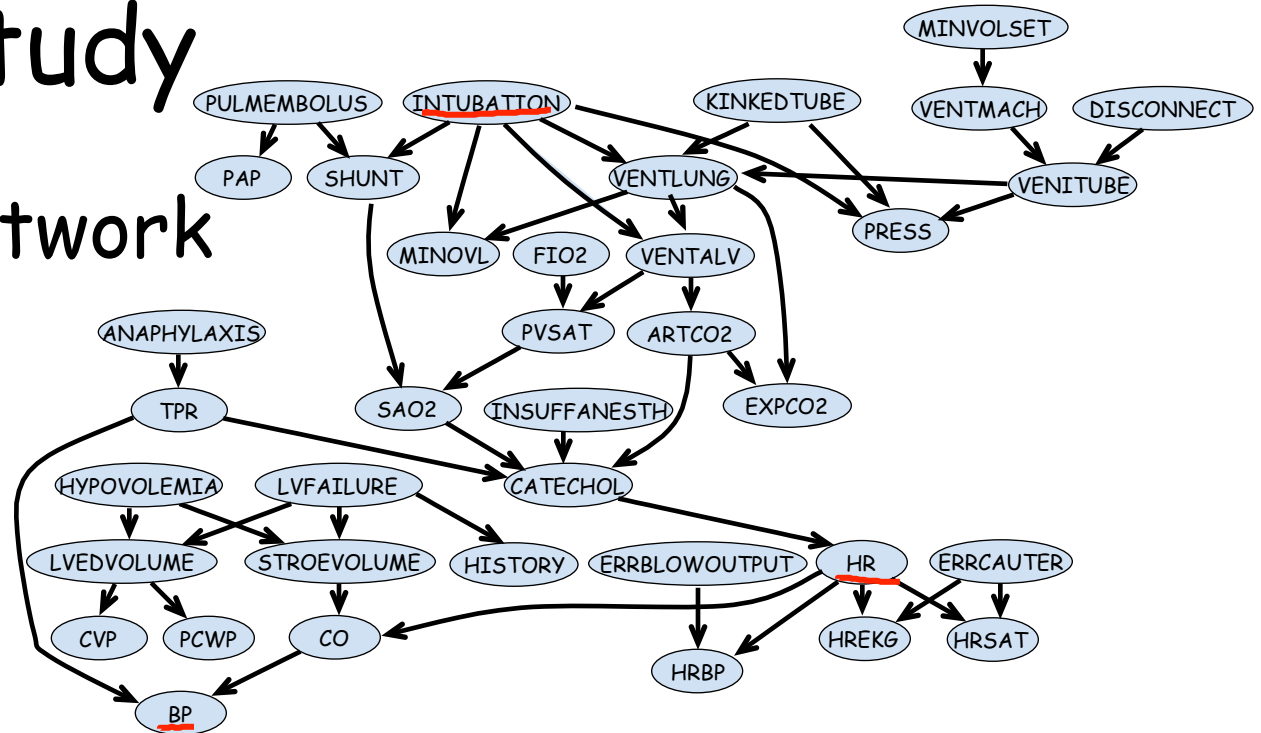
- We need hyperparameter $\alpha_{x|u}$ for each node X , value x , and parent assignment u
 - Prior network with parameters Θ_0
 - Equivalent sample size parameter α
 - $\alpha_{x|u} := \alpha \cdot P(x, u | \Theta_0)$ $X=y, \bar{u}=\bar{u}$



Case Study

- ICU-Alarm network

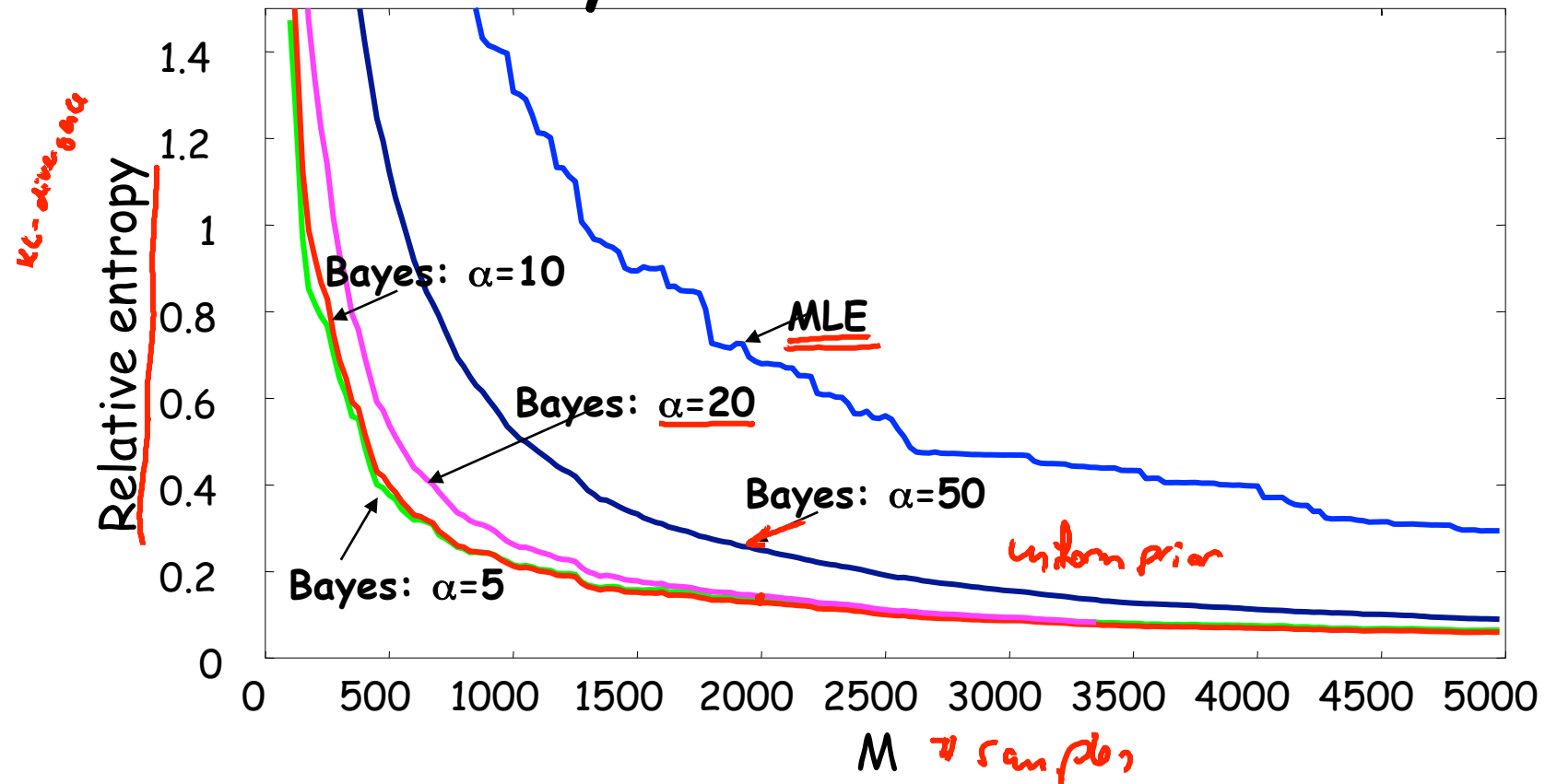
- 37 variables
- 504 params



- Experiment

- Sample instances from network
- Relearn parameters

Case Study: ICU Alarm Network



Summary

- In Bayesian networks, if parameters are independent a priori, then also independent in the posterior
- For multinomial BNs, estimation uses sufficient statistics $M[x, u]$

$$\hat{\theta}_{x|u} = \frac{M[x, u]}{M[u]}$$

MLE

$$P(x | u, D) = \frac{\alpha_{x,u} + M[x, u]}{\alpha_u + M[u]}$$

Bayesian (Dirichlet)

- Bayesian methods require choice of prior
 - can be elicited as prior network and equivalent sample size α