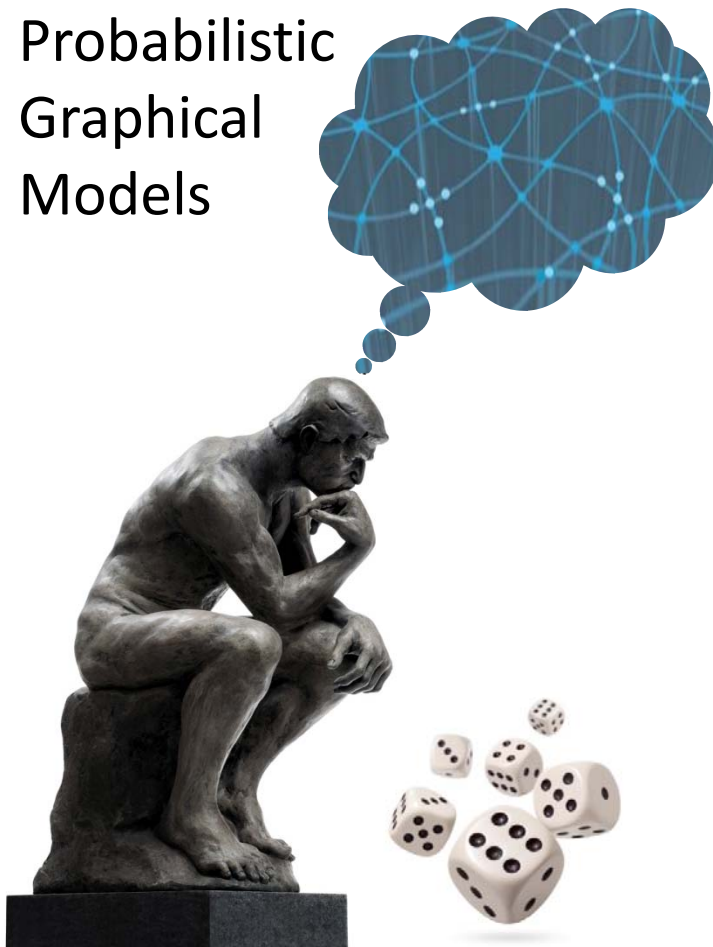Probabilistic Graphical Models
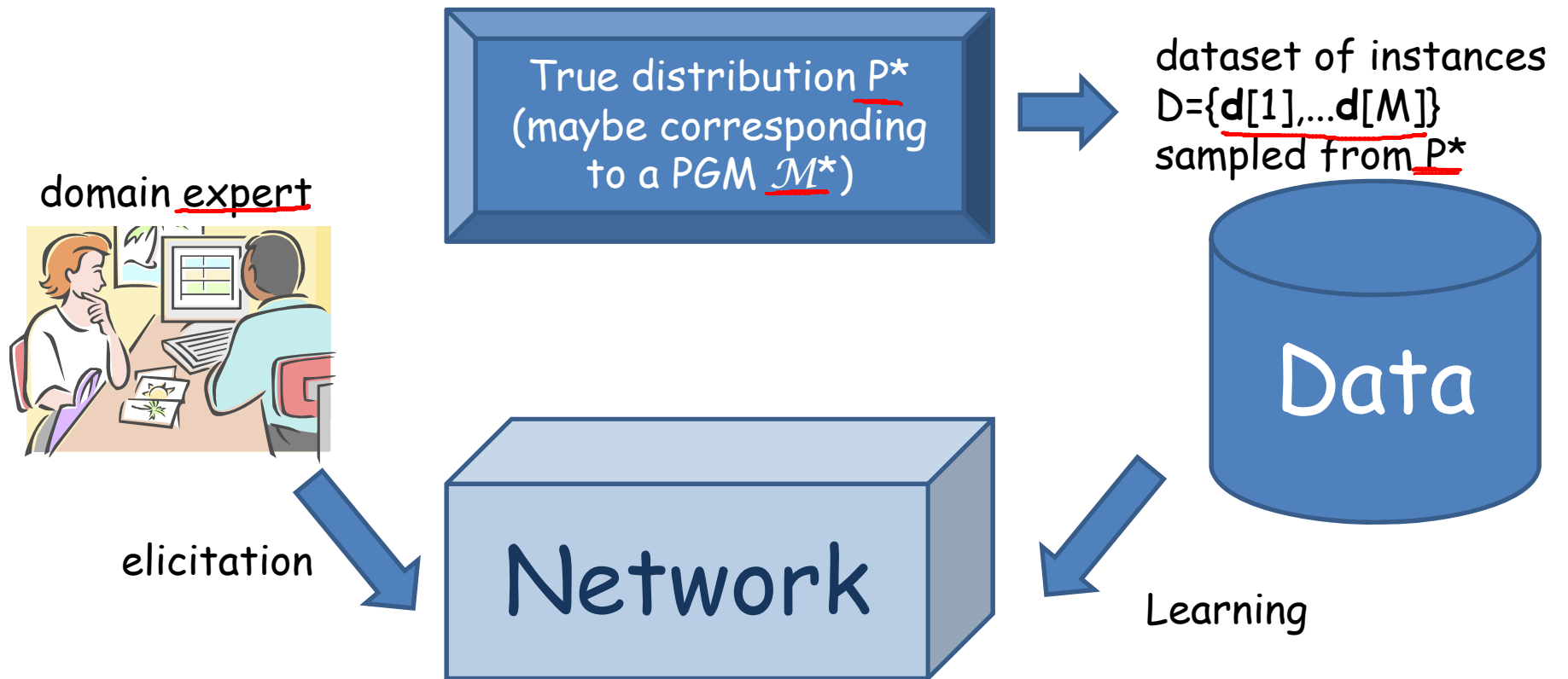
Learning
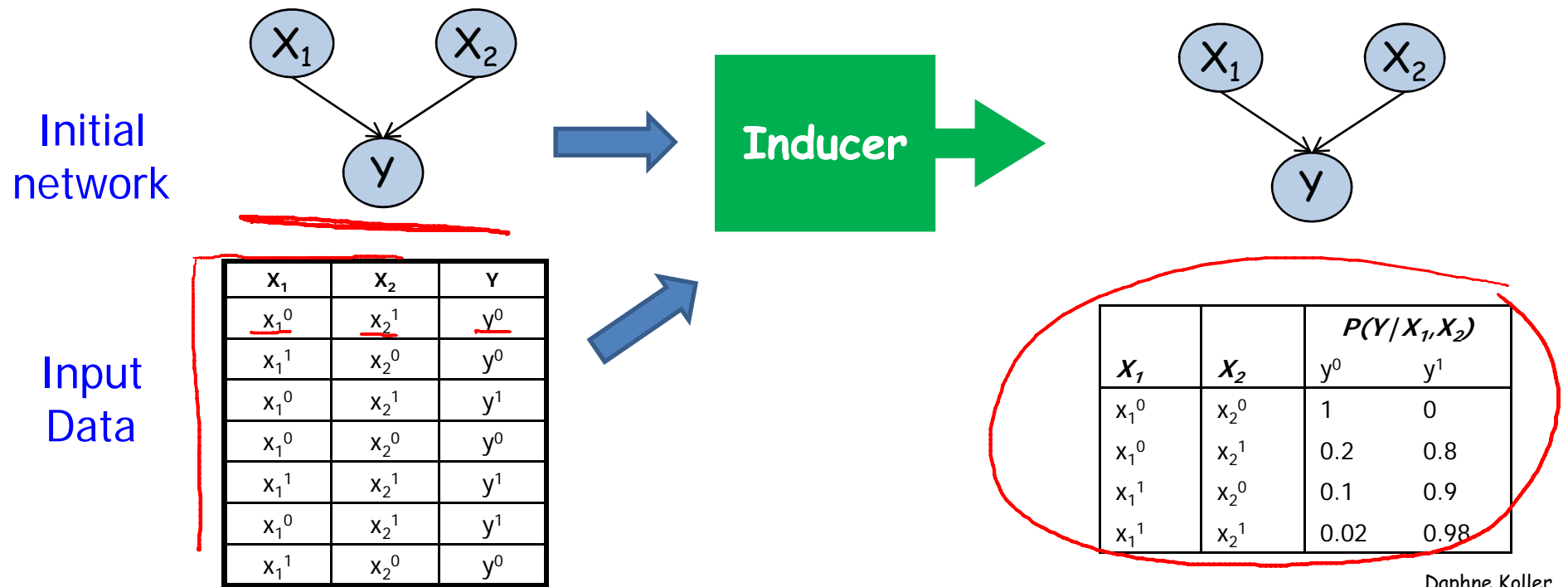
Overview

# PGM Learning Tasks and Metrics

# Learning
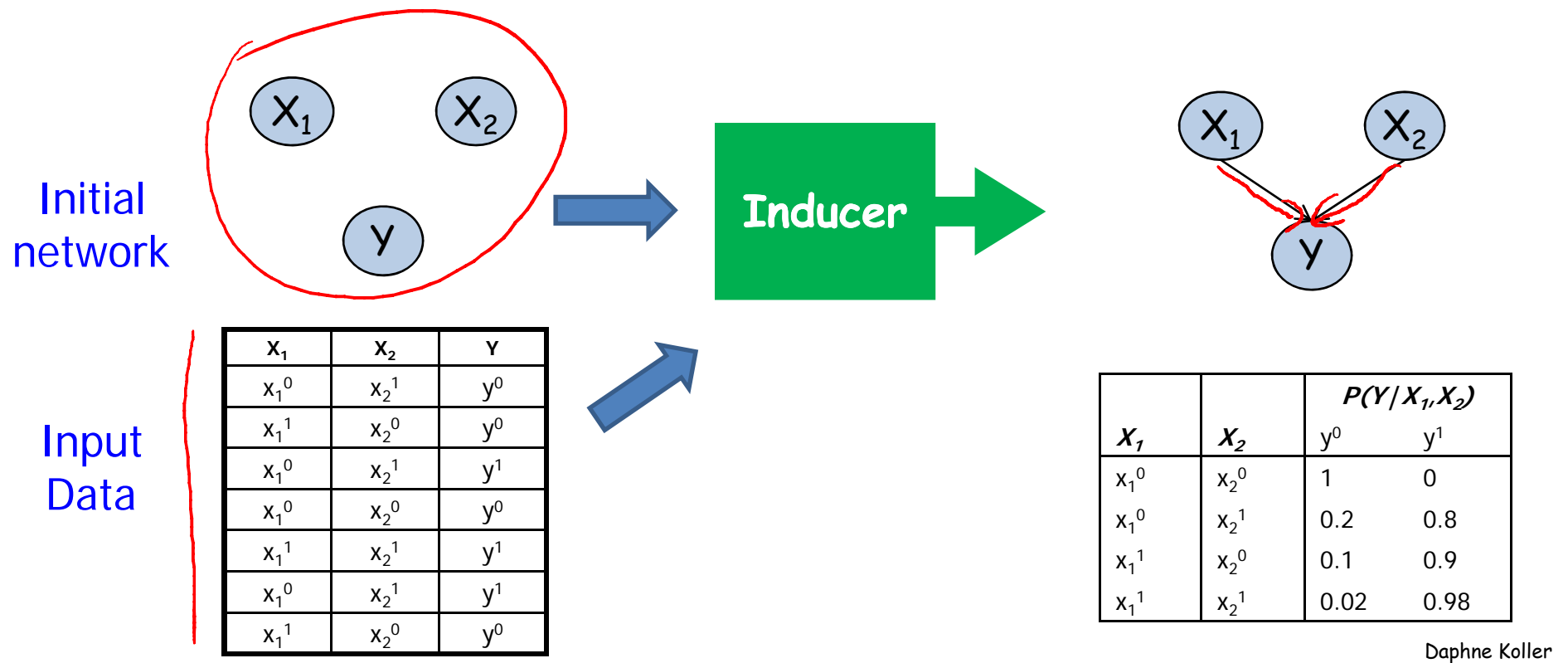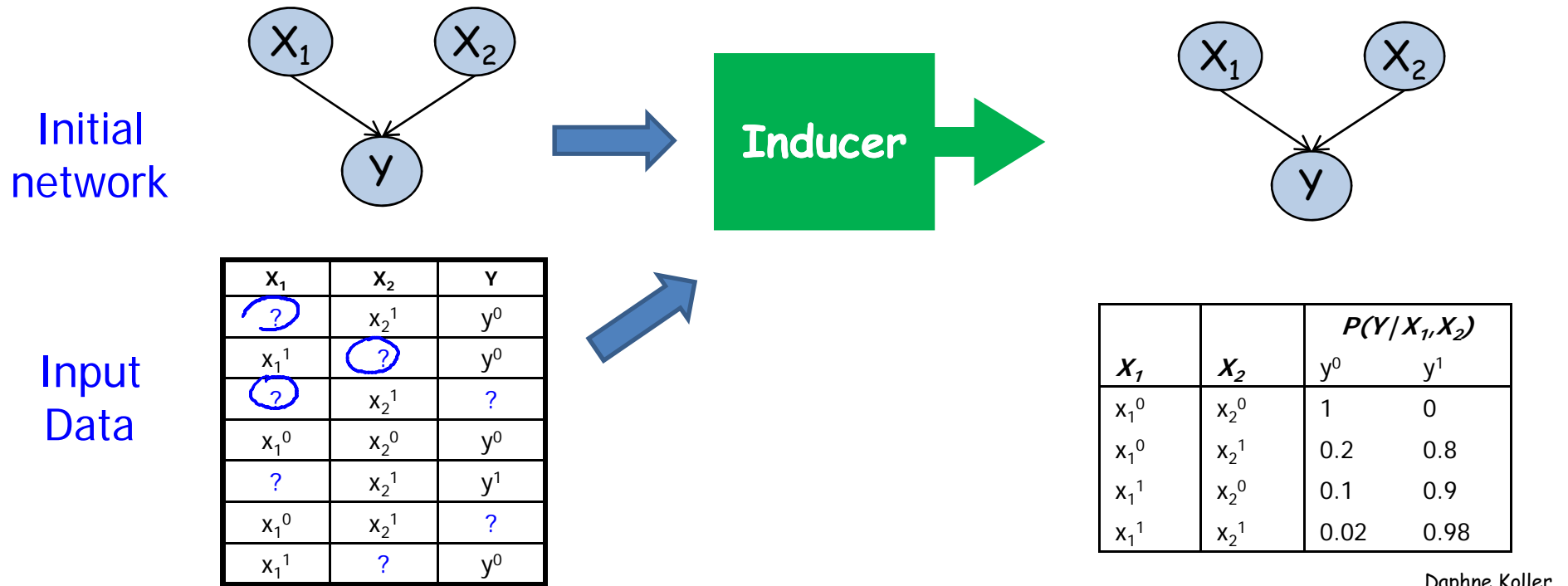


True distribution P*
(maybe corresponding
to a PGM 𝕄*)

dataset of instances
D={**d**[1],...**d**[M]}
sampled from P*

domain expert

Data

elicitation

Network

Learning

Daphne Koller

# Known Structure, Complete Data

Initial network

Input Data

| $X_1$ | $X_2$ | Y |
|---|---|---|
| $x_1^0$ | $x_2^1$ | $y^0$ |
| $x_1^1$ | $x_2^0$ | $y^0$ |
| $x_1^0$ | $x_2^1$ | $y^1$ |
| $x_1^0$ | $x_2^0$ | $y^0$ |
| $x_1^1$ | $x_2^1$ | $y^1$ |
| $x_1^0$ | $x_2^1$ | $y^1$ |
| $x_1^1$ | $x_2^0$ | $y^0$ |

**Inducer**

|  |  | $P(Y \mid X_1, X_2)$ | |
|---|---|---|---|
| $X_1$ | $X_2$ | $y^0$ | $y^1$ |
| $x_1^0$ | $x_2^0$ | 1 | 0 |
| $x_1^0$ | $x_2^1$ | 0.2 | 0.8 |
| $x_1^1$ | $x_2^0$ | 0.1 | 0.9 |
| $x_1^1$ | $x_2^1$ | 0.02 | 0.98 |

Daphne Koller

# Unknown Structure, Complete Data



**Initial network**

**Input Data**

| $X_1$ | $X_2$ | Y |
|-------|-------|---|
| $x_1^0$ | $x_2^1$ | $y^0$ |
| $x_1^1$ | $x_2^0$ | $y^0$ |
| $x_1^0$ | $x_2^1$ | $y^1$ |
| $x_1^0$ | $x_2^0$ | $y^0$ |
| $x_1^1$ | $x_2^1$ | $y^1$ |
| $x_1^0$ | $x_2^1$ | $y^1$ |
| $x_1^1$ | $x_2^0$ | $y^0$ |

**Inducer**

| $X_1$ | $X_2$ | $P(Y\|X_1,X_2)$ | |
|-------|-------|------|------|
| | | $y^0$ | $y^1$ |
| $x_1^0$ | $x_2^0$ | 1 | 0 |
| $x_1^0$ | $x_2^1$ | 0.2 | 0.8 |
| $x_1^1$ | $x_2^0$ | 0.1 | 0.9 |
| $x_1^1$ | $x_2^1$ | 0.02 | 0.98 |

Daphne Koller

# Known Structure, Incomplete Data

Initial network

Input Data



| $X_1$ | $X_2$ | Y |
|---|---|---|
| ? | $x_2^1$ | $y^0$ |
| $x_1^1$ | ? | $y^0$ |
| ? | $x_2^1$ | ? |
| $x_1^0$ | $x_2^0$ | $y^0$ |
| ? | $x_2^1$ | $y^1$ |
| $x_1^0$ | $x_2^1$ | ? |
| $x_1^1$ | ? | $y^0$ |

| $X_1$ | $X_2$ | $P(Y \mid X_1, X_2)$ | |
|---|---|---|---|
| | | $y^0$ | $y^1$ |
| $x_1^0$ | $x_2^0$ | 1 | 0 |
| $x_1^0$ | $x_2^1$ | 0.2 | 0.8 |
| $x_1^1$ | $x_2^0$ | 0.1 | 0.9 |
| $x_1^1$ | $x_2^1$ | 0.02 | 0.98 |

Daphne Koller

# Unknown Structure, Incomplete Data

**Initial network**



**Input Data**

| X₁ | X₂ | Y |
|---|---|---|
| ? | $x_2^1$ | $y^0$ |
| $x_1^1$ | ? | $y^0$ |
| ? | $x_2^1$ | ? |
| $x_1^0$ | $x_2^0$ | $y^0$ |
| ? | $x_2^1$ | $y^1$ |
| $x_1^0$ | $x_2^1$ | ? |
| $x_1^1$ | ? | $y^0$ |

**Inducer**

| $X_1$ | $X_2$ | $P(Y \mid X_1, X_2)$ | |
|---|---|---|---|
| | | $y^0$ | $y^1$ |
| $x_1^0$ | $x_2^0$ | 1 | 0 |
| $x_1^0$ | $x_2^1$ | 0.2 | 0.8 |
| $x_1^1$ | $x_2^0$ | 0.1 | 0.9 |
| $x_1^1$ | $x_2^1$ | 0.02 | 0.98 |

Daphne Koller

# Latent Variables, Incomplete Data



Daphne Koller

# PGM Learning Tasks I

- Goal: Answer general probabilistic queries about new instances

- Simple metric: Training set likelihood

  *data*

  $- P(D : \mathcal{M}) = \Pi_m P(\mathbf{d}[m] : \mathcal{M})$   *(IID)*

- But we really care about new data

  $-$ Evaluate on test set likelihood $- P(D' : \mathcal{M})$

  *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 $\Pi_m P(\mathbf{d}[m] : \mathcal{M})$ or
  - conditional likelihood $\Pi_m P(\mathbf{y}[m] \mid \mathbf{x}[m] : \mathcal{M})$
- Model evaluated on "true" objective over test data

Daphne Koller

# PGM Learning Tasks III

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

Daphne Koller

# Avoiding Overfitting

- Selecting $M$ 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

Daphne Koller

# Selecting <u>Hyperparameters</u>

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

Daphne Koller

# 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

Daphne Koller

Probabilistic Graphical Models

Learning

Parameter Estimation

# Maximum Likelihood Estimation

# Biased Coin Example

P is a Bernoulli distribution:
$$P(X=1) = \theta, \; P(X=0) = 1-\theta$$

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

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

# IID as a PGM



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

Daphne Koller

# Maximum Likelihood Estimation

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

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

$$L(\theta : \langle H, T, T, H, H \rangle)$$

$\underset{\theta}{P(H \mid \theta)} \cdot \underset{(1-\theta)}{P(T \mid \theta)} \cdot \underset{(1-\theta)}{P(T \mid \theta)} \cdot \underset{\theta}{P(H \mid \theta)} \cdot \underset{\theta}{P(H \mid \theta)} = \theta^5 (1-\theta)^2$

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



*(plot with vertical axis labeled $L(D:\theta)$ and horizontal axis labeled $\theta$ with marks at 0, 0.2, 0.4, 0.6, 0.8, 1)*

Daphne Koller

# Maximum Likelihood Estimator

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

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

Daphne Koller

# Sufficient Statistics

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

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

- → $M_H$ and $M_T$ are sufficient statistics

# Sufficient Statistics

- A function s(D) is a <u>sufficient statistic</u> from instances to a vector in $\Re^k$ if for any two datasets D and D' and any $\theta \in \Theta$ we have

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

$s(D) \qquad s(D')$



Datasets

Statistics

# Sufficient Statistic for Multinomial

- For a dataset D over variable X with k values, the sufficient statistics are counts $\langle M_1,...,M_k \rangle$ where $M_i$ is the # of times that $X[m]=x^i$ in D

- Sufficient statistic s(x) is a tuple of dimension k
  - $s(x^i)=(0,...0,1,0,...,0)$       $\sum_n s(x[m]) = (M_1, M_2, ..., M_k)$

  i

  $$L(\theta : D) = \prod_{i=1}^{k} \theta_i^{M_i}$$

  $\theta_i$ param for $X=x^i$

# Sufficient Statistic for Gaussian

- Gaussian distribution:

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

- 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)=<1,x,x²>   $s(D) = \left( \sum_m x[m]^2, \sum_m x[m], m \right)$

# Maximum Likelihood Estimation

- MLE Principle: Choose θ to maximize L(D:Θ)

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

- Gaussian MLE: $\hat{\mu} = \dfrac{1}{M} \sum_m x[m]$   empirical mean

  $\hat{\sigma} = \sqrt{\dfrac{1}{M} \sum_m (x[m] - \hat{\mu})^2}$   empirical st dev

Daphne Koller

# 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: $\theta_{x^0}, \theta_{x^1}$
  $$\theta_{y^0|x^0}, \theta_{y^1|x^0}, \theta_{y^0|x^1}, \theta_{y^1|x^1}$$

- Data instances: <x[m],y[m]>

| X | |
|---|---|
| x$^0$ | x$^1$ |
| 0.7 | 0.3 |

$P(x)$

X

↓

Y

| 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

$\theta_X$

- Parameters: $\{\theta_x : x \in Val(X)\}$

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

$X$

$\theta_{Y|X}$

$Y$

Data $d$

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

chain rule for BNs

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

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

product of two local likelihood

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

Daphne Koller

# MLE for Bayesian Networks

- Likelihood for Bayesian network

$$L(\Theta : D) \quad = \prod_m P(x[m] : \Theta)$$

*parents of $X_i$*

*chain rule*

$$= \prod_m \prod_i P(x_i[m] \mid U_i[m] : \Theta_i)$$

$$= \prod_i \prod_m P(x_i[m] \mid U_i[m] : \Theta_i)$$

*local likelihood* $\Longrightarrow$

$$= \prod_i L_i(\Theta_i : D) \quad L_i(\Theta_{x_i} : D)$$

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

Daphne Koller

# MLE for Table CPDs

$$\prod_{m=1}^{M} P(x[m] \mid \boldsymbol{u}[m] : \theta) = \prod_{m=1}^{M} P(x[m] \mid \boldsymbol{u}[m] : \theta_{X \mid U})$$

$$= \prod_{x, \boldsymbol{u}} \left( \prod_{m : x[m] = x, \boldsymbol{u}[m] = \boldsymbol{u}} P(x[m] \mid \boldsymbol{u}[m] : \theta_{X \mid U}) \right)$$

$$P(x[m] = x \mid u[m] = u : \theta_{x \mid u}) = \theta_{x \mid u}$$

$$= \prod_{x, \boldsymbol{u}} \left( \prod_{m : x[m] = x, \boldsymbol{u}[m] = \boldsymbol{u}} \theta_{x \mid \boldsymbol{u}} \right)$$

fraction of $X = x$ among cases where $\bar{u} = \bar{u}$

$$= \prod_{x, \boldsymbol{u}} \theta_{x \mid \boldsymbol{u}}^{M[x, \boldsymbol{u}]}$$

$P(x|u)$

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

Daphne Koller

# Shared Parameters



$$\theta_{S'|S}$$

$S^{(0)} \rightarrow S^{(1)} \rightarrow S^{(2)} \rightarrow S^{(3)}$

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

$$= \prod_{i,j} \prod_{t:S^{(t)}=s^i, S^{(t+1)}=s^j} P(S^{(t+1)} \mid S^{(t)} : \theta_{S'|S})$$

$$s^i \rightarrow s^j$$

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

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

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

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

Daphne Koller

# Shared Parameters

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

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



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

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

Daphne Koller

# Summary

- For BN with <u>disjoint sets of parameters in CPDs</u>, likelihood decomposes as product of <u>local likelihood functions</u>, one per variable

- For table CPDs, local likelihood further <u>decomposes</u> as product of <u>likelihood for multinomials</u>, one for each parent combination

- For <u>networks with shared CPDs</u>, <u>sufficient statistics</u> 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
  - ⇨ **very poor parameter estimates** ⇦
- **With <u>limited data</u>, we often get better generalization with <u>simpler structures</u>**
  - *even when wrong*

Daphne Koller

# 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

Daphne Koller

# Parameter Estimation as a PGM



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

Daphne Koller

# Bayesian Inference

$P(\theta)$    **θ**    PGM

- Joint probabilistic model    X[1]   . . .   X[M]

$$P(x[1],...,\ x[M],\theta) = P(x[1],...,\ x[M]\,|\,\theta)P(\theta)$$

$$= P(\theta)\prod_{i=1}^{M}P(x[i]\,|\,\theta) \quad \leftarrow \text{likelihood function}$$

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

likelihood        prior

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

data D      constant relative to θ

Daphne Koller

# Dirichlet Distribution

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

  hyperparameters

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

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

# Dirichlet Distributions



Beta = Dirichlet($\alpha_1, \alpha_2$)

mix $\alpha_H, \alpha_T$ determines position of peak

$\alpha = \alpha_H + \alpha_T$ determines how sharp it is

Legend:
- Dirichlet(1,1)
- Dirichlet(2,2)
- Dirichlet(0.5 0.5)
- Dirichlet(5,5)

# Dirichlet Priors & Posteriors

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

$$\theta_i^{M_i + \alpha_i - 1}$$

$$M_i = \text{\# instances with } x_i \qquad P(D \mid \theta) = \prod_{i=1}^{k} \underbrace{\theta_i^{M_i}}_{\text{multinomial } \theta} \qquad P(\theta) \propto \prod_{i=1}^{k} \underbrace{\theta_i^{\alpha_i - 1}}$$

- If P(θ) is Dirichlet and the likelihood is multinomial, then the posterior is also Dirichlet
  – Prior is Dir($\alpha_1$,...,$\alpha_k$)
  – Data counts are $M_1$,...,$M_k$
  – Posterior is Dir($\alpha_1 + M_1$,...$\alpha_k + M_k$)
- 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

Daphne Koller

Probabilistic
Graphical
Models

Learning

Parameter Estimation

# Bayesian Prediction

# Bayesian Prediction



$\theta$ ~ Dirichlet($\alpha_1,...,\alpha_k$)

$$P(X) = \int_\theta P(X \mid \theta) P(\theta) d\theta$$

*marginalizing over $\theta$*

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

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

*fraction of instances we've seen where $x^i$*

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

Daphne Koller

# Bayesian Prediction

$\theta$ ~ Dirichlet($\alpha_1, ..., \alpha_k$)

X[1] . . . X[M] X[M+1]

$P(x[M+1] \mid x[1], ..., x[M])$

$$= \int_\theta P(x[M+1] \mid x[1], ..., x[M], \theta) P(\theta \mid x[1], ..., x[M]) \, d\theta$$

~Dirichlet($\alpha_1 + M_1, ..., \alpha_k + M_k$)

posterior over $\theta$ given D

$$= \int_\theta P(x[M+1] \mid \theta) P(\theta \mid x[1], ..., x[M]) \, d\theta$$

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

$\alpha = \sum \alpha_i$

$M = \sum M_i$

- <u>Equivalent sample size</u> $\alpha = \alpha_1 + ... + \alpha_K$
  - Larger $\alpha \Rightarrow$ more confidence in our prior

Daphne Koller

# Example: Binomial Data

- Prior: uniform for $\theta$ in [0,1]

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

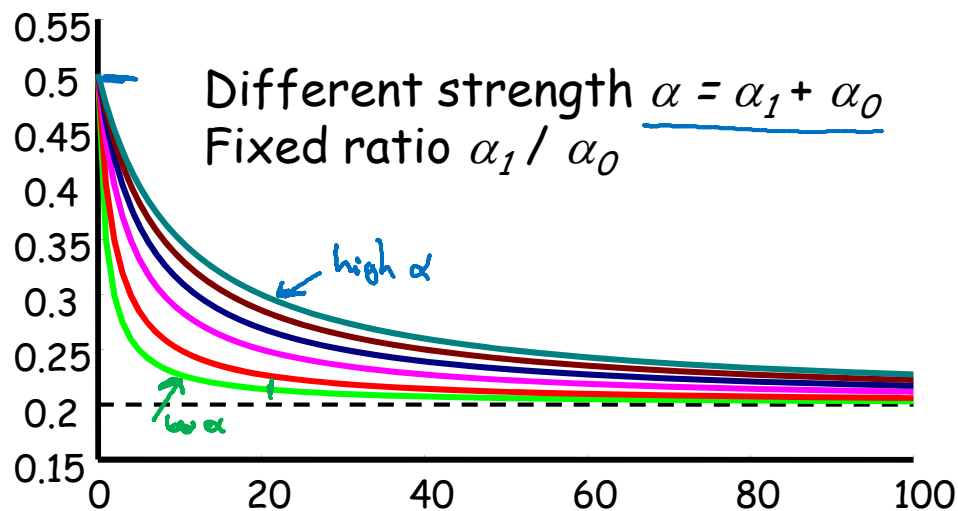Dirichlet (1,1)

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



- MLE for P(X[6]=1)=4/5
- Bayesian prediction is 5/7
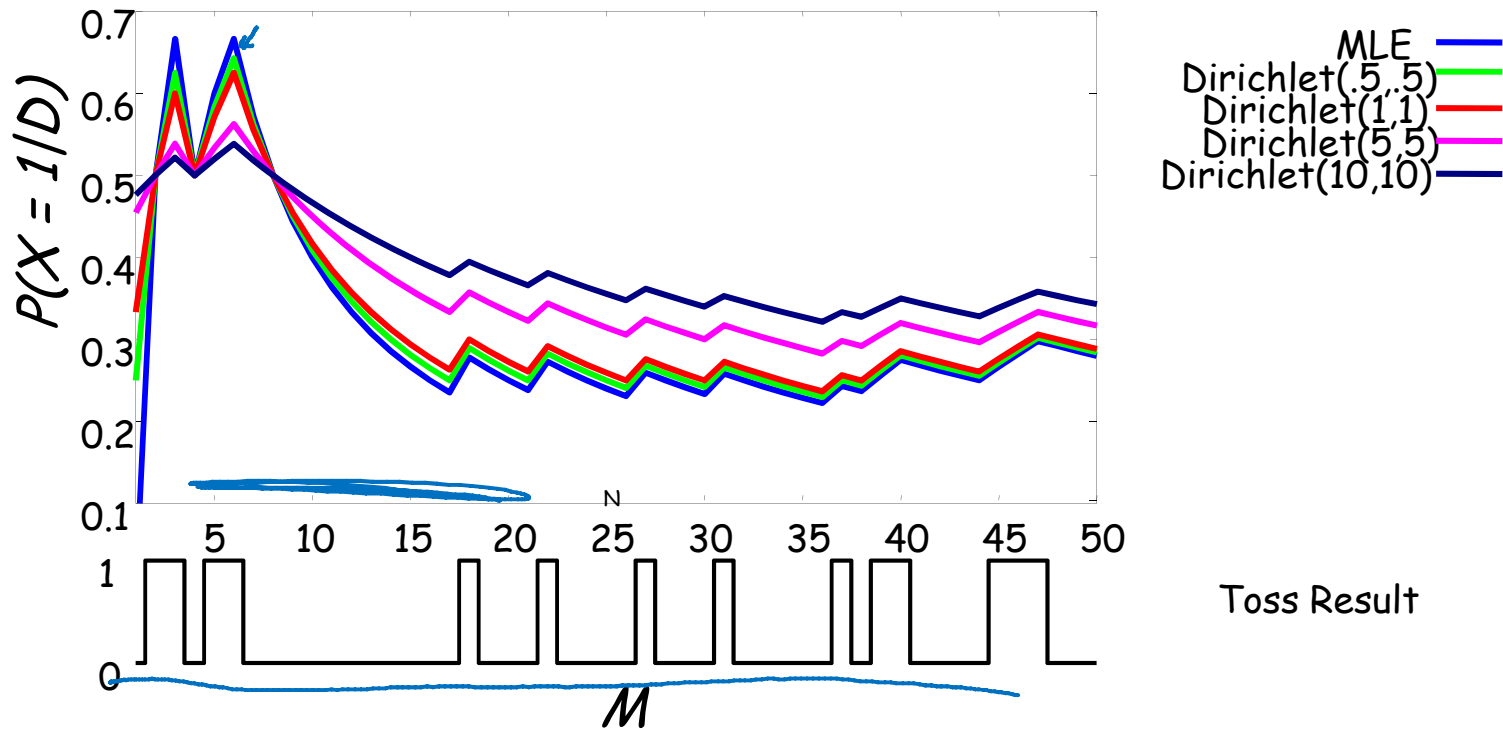
$$\frac{\alpha_1 + M_1}{\alpha + M} = \frac{1+4}{2+5}$$

Daphne Koller

# Effect of Priors

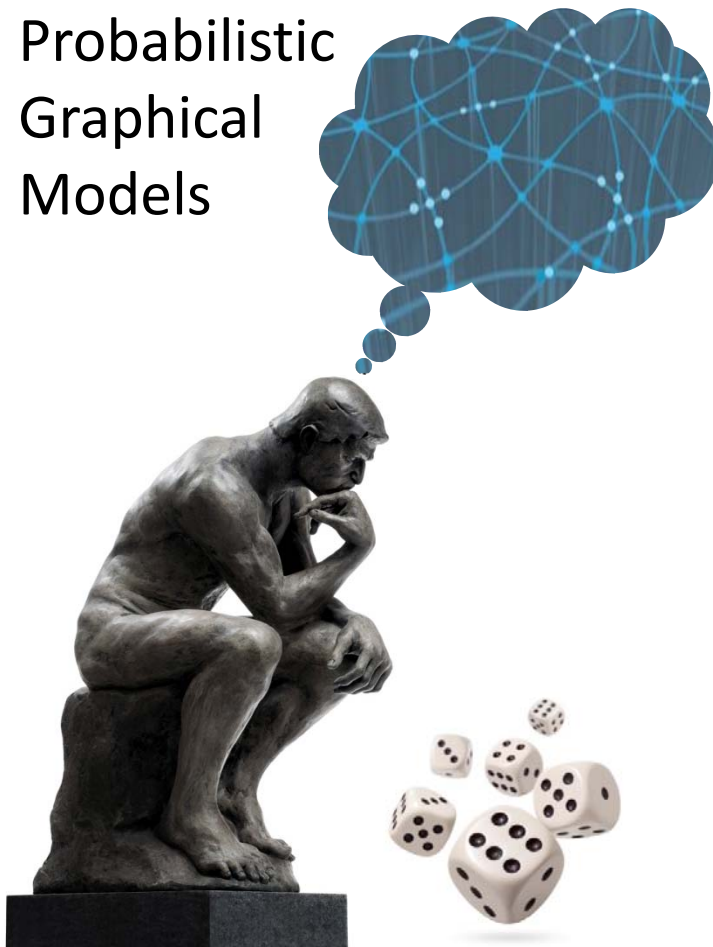- Prediction of P(X=1) after seeing data with $M_1 = \frac{1}{4} M_0$ as a function of sample size M

Different strength $\alpha = \alpha_1 + \alpha_0$
Fixed ratio $\alpha_1 / \alpha_0$

high $\alpha$

$\omega \alpha$

Fixed strength $\alpha = \alpha_1 + \alpha_0$
Different ratio $\alpha_1 / \alpha_0$

Daphne Koller

# Effect of Priors

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



MLE
Dirichlet(.5,.5)
Dirichlet(1,1)
Dirichlet(5,5)
Dirichlet(10,10)

$P(X = 1/D)$

Toss Result

Daphne Koller

# Summary

- Bayesian prediction combines sufficient statistics from imaginary Dirichlet samples and real data samples

- Asymptotically the same as MLE

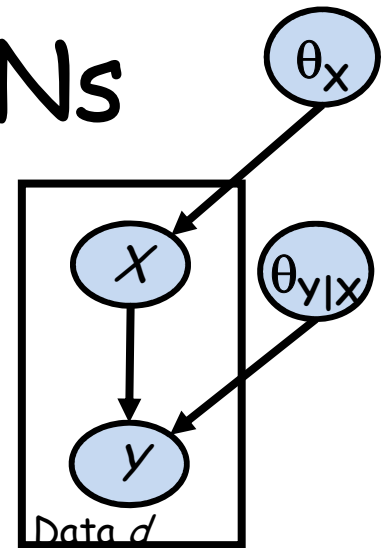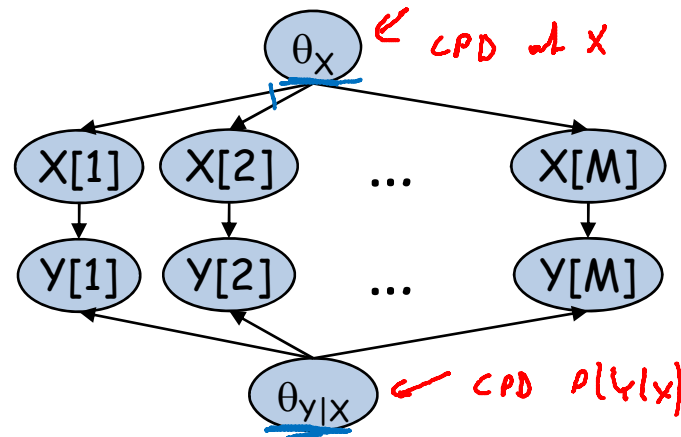- But Dirichlet hyperparameters determine both the prior beliefs and their strength

Daphne Koller

Probabilistic
Graphical
Models



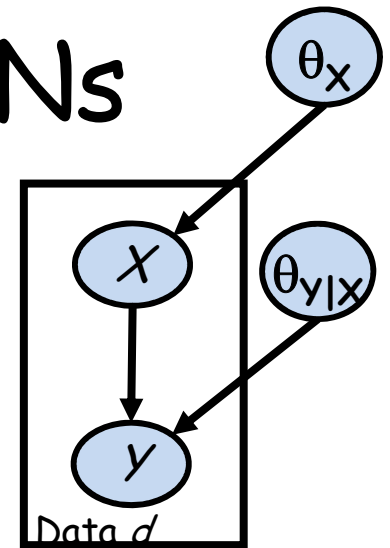Learning

Parameter Estimation
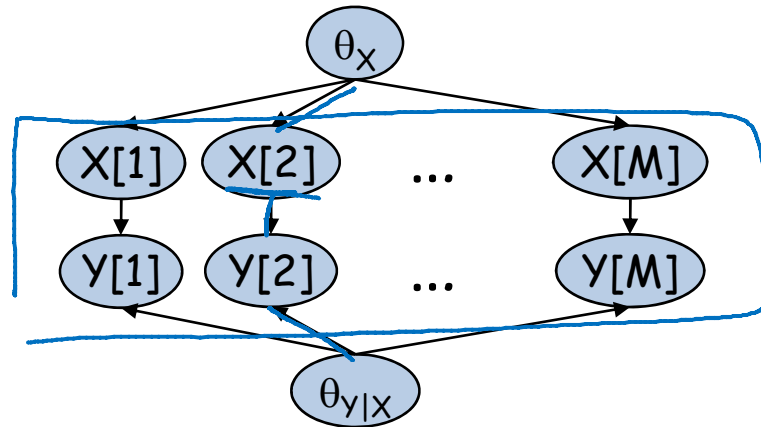
# Bayesian Estimation for BNs

# Bayesian Estimation in BNs



- **Instances are <u>independent</u> given the parameters**
  - (X[m'],Y[m']) are d-separated from (X[m],Y[m]) given θ
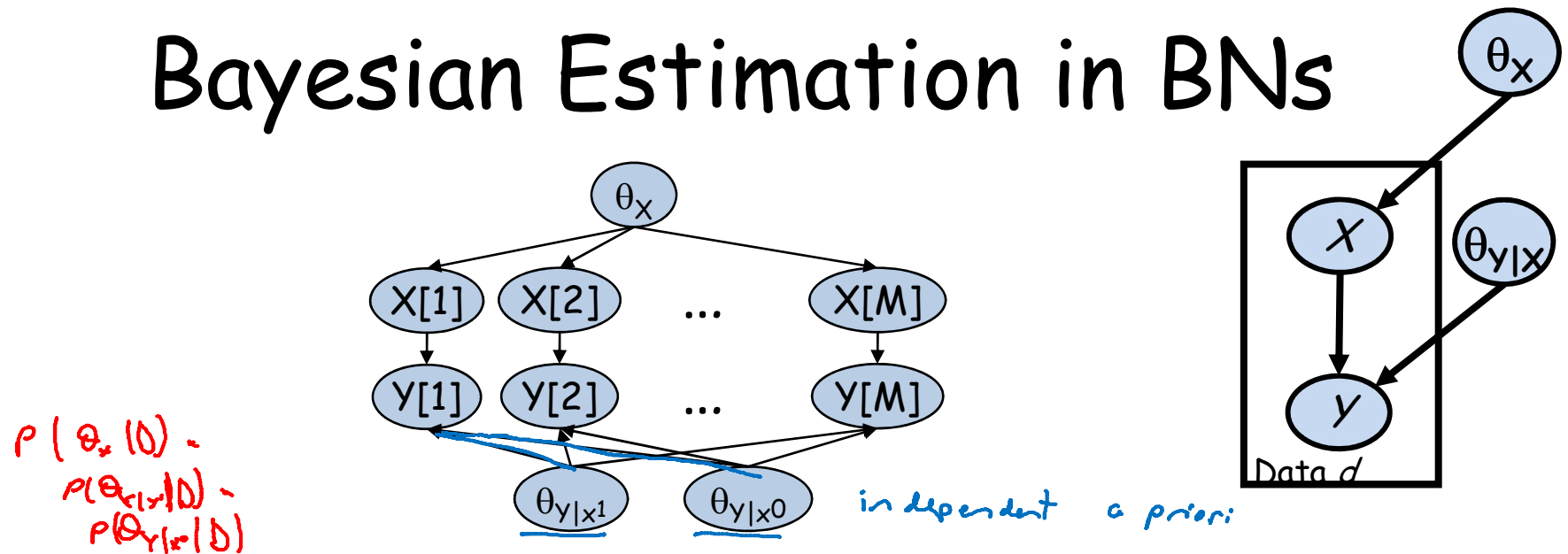- **Parameters for individual variables are independent a priori**
  $$P(\theta) = \prod_i P(\theta_{X_i | Pa(X_i)})$$
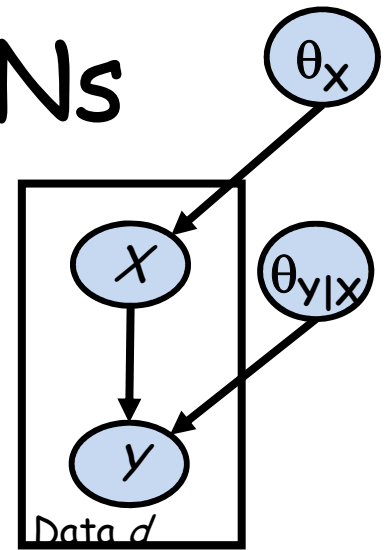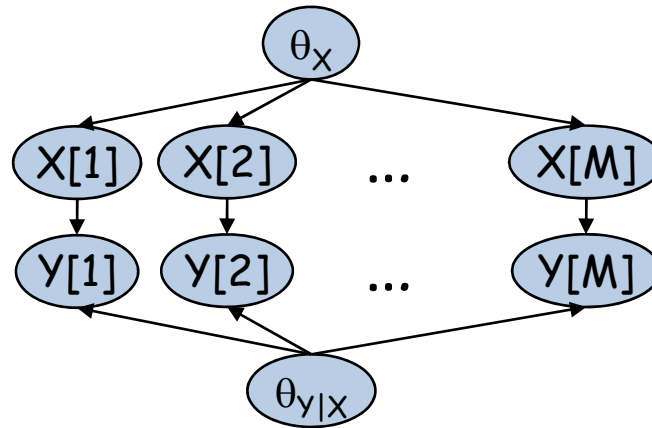
Daphne Koller

# Bayesian Estimation in BNs



- ## Posteriors of $\theta$ 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 $\theta$ are independent given complete data
  - Also holds for parameters within families
  - Note context specific independence between $\theta_{y|x^1}$ and $\theta_{y|x^0}$ when given both X's and Y's

Daphne Koller

# Bayesian Estimation in BNs



- ## Posteriors of $\theta$ can be computed independently
  - For multinomial $\theta_{X|\mathbf{u}}$ if prior is Dirichlet($\alpha_{x^1|\mathbf{u}},..., \alpha_{x^k|\mathbf{u}}$)

    *assignment to x's parents $\bar{u}$*
  - posterior is Dirichlet($\alpha_{x^1|\mathbf{u}}+M[x^1,\mathbf{u}],...,\alpha_{x^k|\mathbf{u}}+M[x^k,\mathbf{u}]$ )

# Assessing Priors for BNs

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

$(X)$   $\Theta_0$ uniform     $(X)$ $\theta_x \sim$ Dirichlet$\left(\frac{\alpha}{2}, \frac{\alpha}{2}\right)$

$(Y)$     $\theta_{Y|x^0} \sim$ Dirichlet$\left(\frac{\alpha}{4}, \frac{\alpha}{4}\right)$
$\theta_{Y|x^1} \sim$ Dirichlet$\left(\frac{\alpha}{4}, \frac{\alpha}{4}\right)$
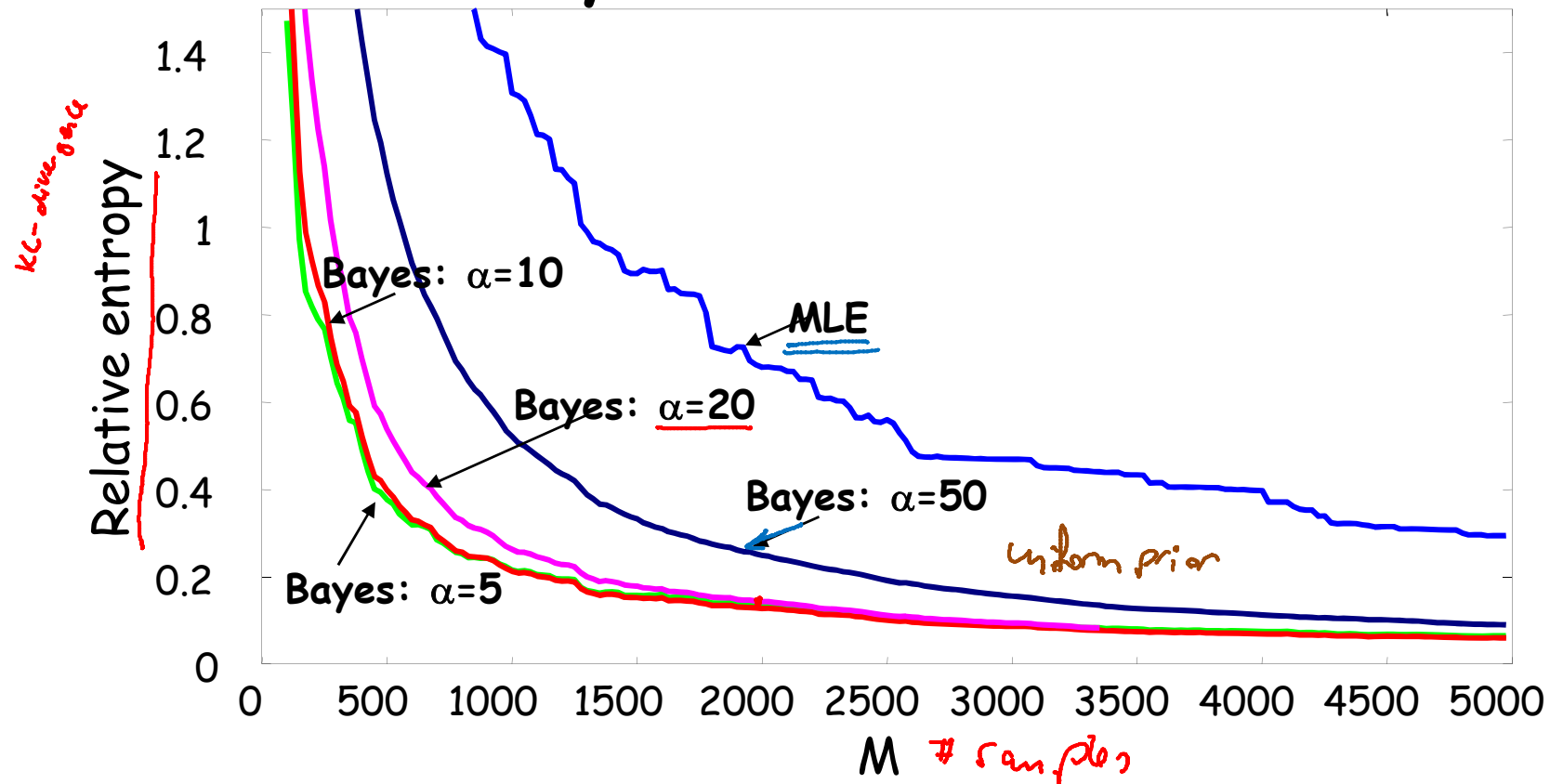
# Case Study

- ICU-Alarm network
  - 37 variables
  - 504 params



- Experiment
  - Sample instances from network
  - Relearn parameters

# Case Study: ICU Alarm Network



Daphne Koller

# 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

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

Bayesian (Dirichlet)

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