



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

PROBABILISTIC GRAPHICAL MODEL SESSION # 13: LEARNING

SEETHA PARAMESWARAN
seetha.p@pilani.bits-pilani.ac.in

TABLE OF CONTENTS

1 LEARNING

2 PARAMETER ESTIMATION IN BAYESIAN NETWORKS

3 MLE FOR BAYESIAN NETWORKS

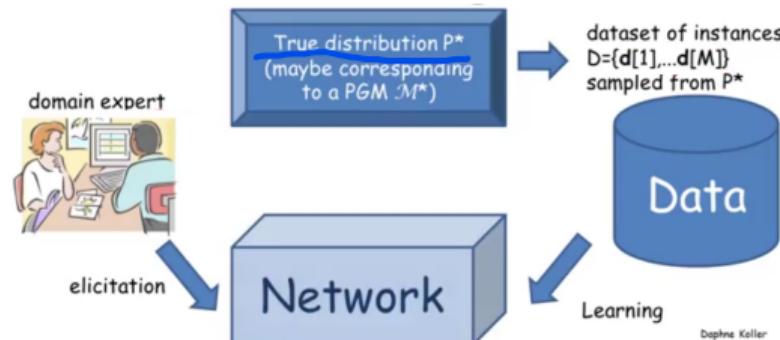
4 BAYESIAN PARAMETER ESTIMATION

MOTIVATION

- For Inference or predictions, the starting point was the PGM. The structure and parameters were part of the input.
- How to acquire a model?
 - ▶ Construct the network by hand, with the help of an expert – “manual” network construction
 - ▶ Learn a model using a set of examples generated from the distribution we wish to model.
- Predictions of structured objects; sequences, graphs, trees
- Incorporate prior knowledge into model.
- Learning a single model for multiple tasks.
- Framework for knowledge discovery.

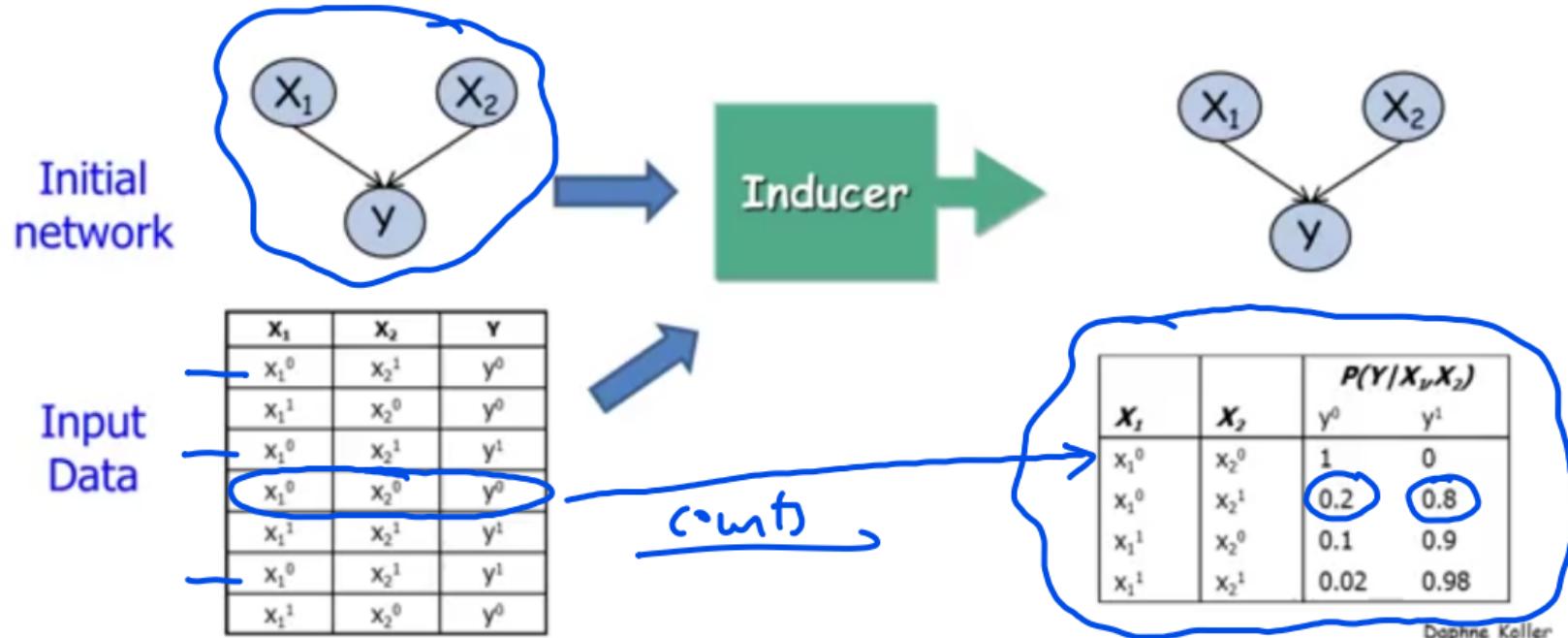
MODEL LEARNING

- The task of constructing a model from a set of instances is generally called **model learning**.
- Goal: Learn a model \tilde{M} from a family of models that defines a distribution $P_{\tilde{M}}$.

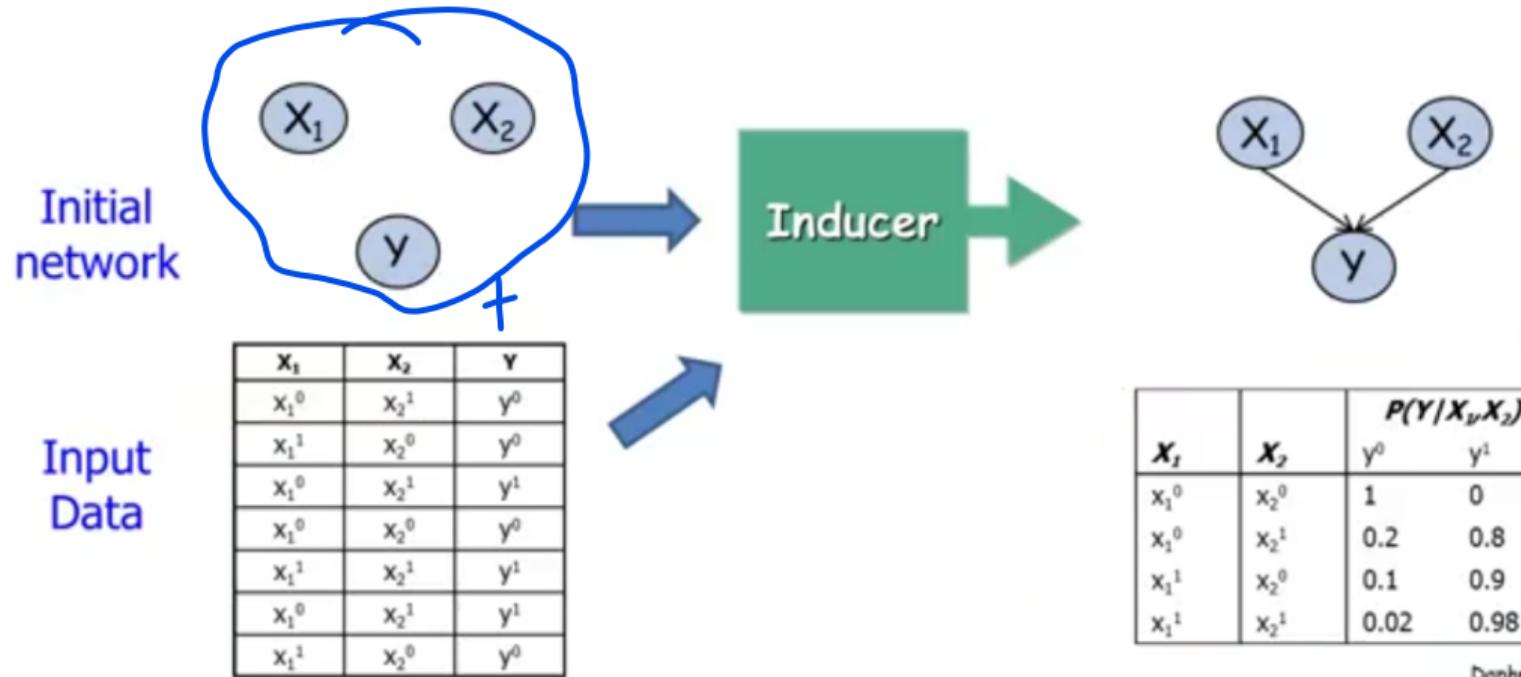


- Amount of data is insufficient, esp for high dimensional distributions. Select \tilde{M} so as to construct the “best” approximation to M^* .

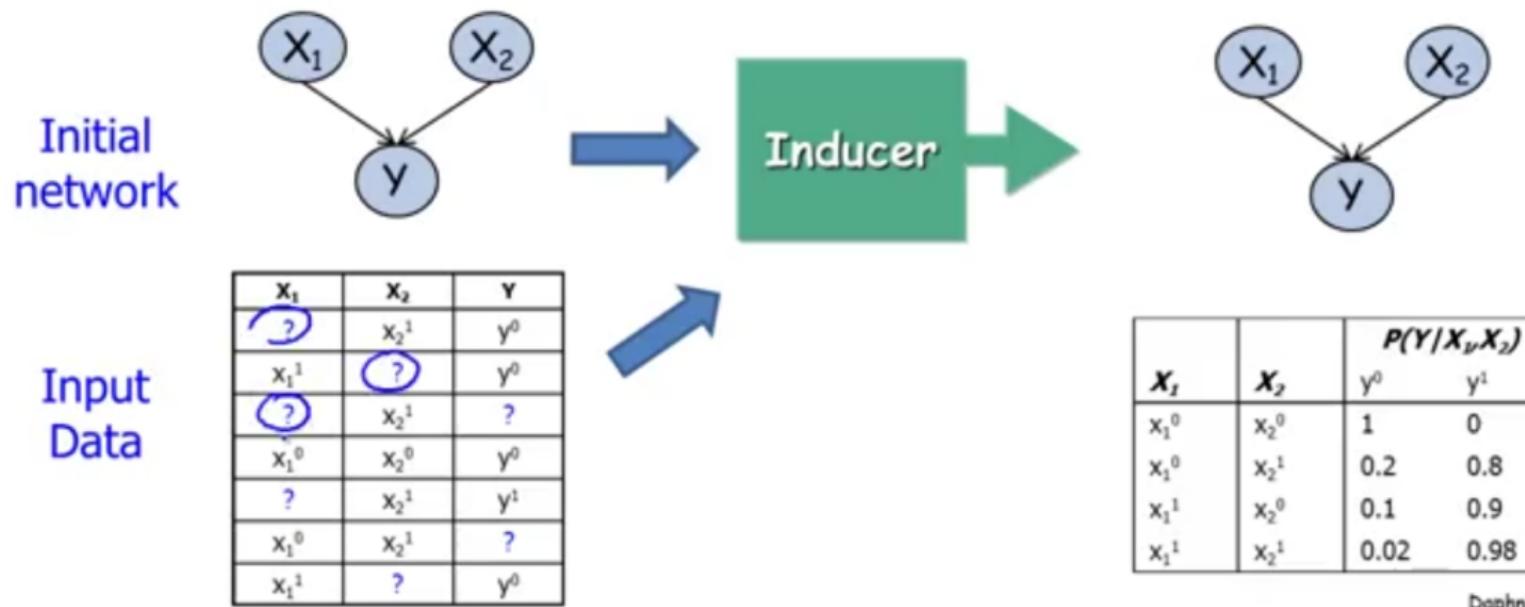
KNOWN STRUCTURE COMPLETE DATA



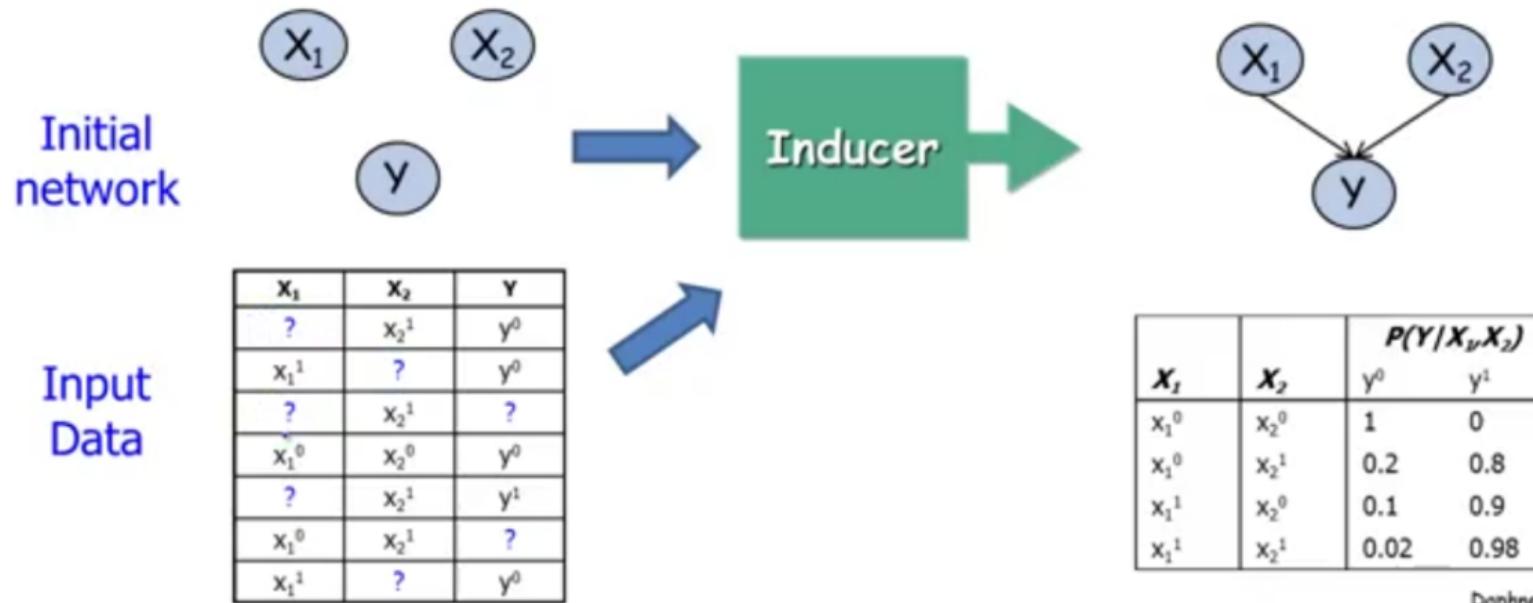
UNKNOWN STRUCTURE COMPLETE DATA



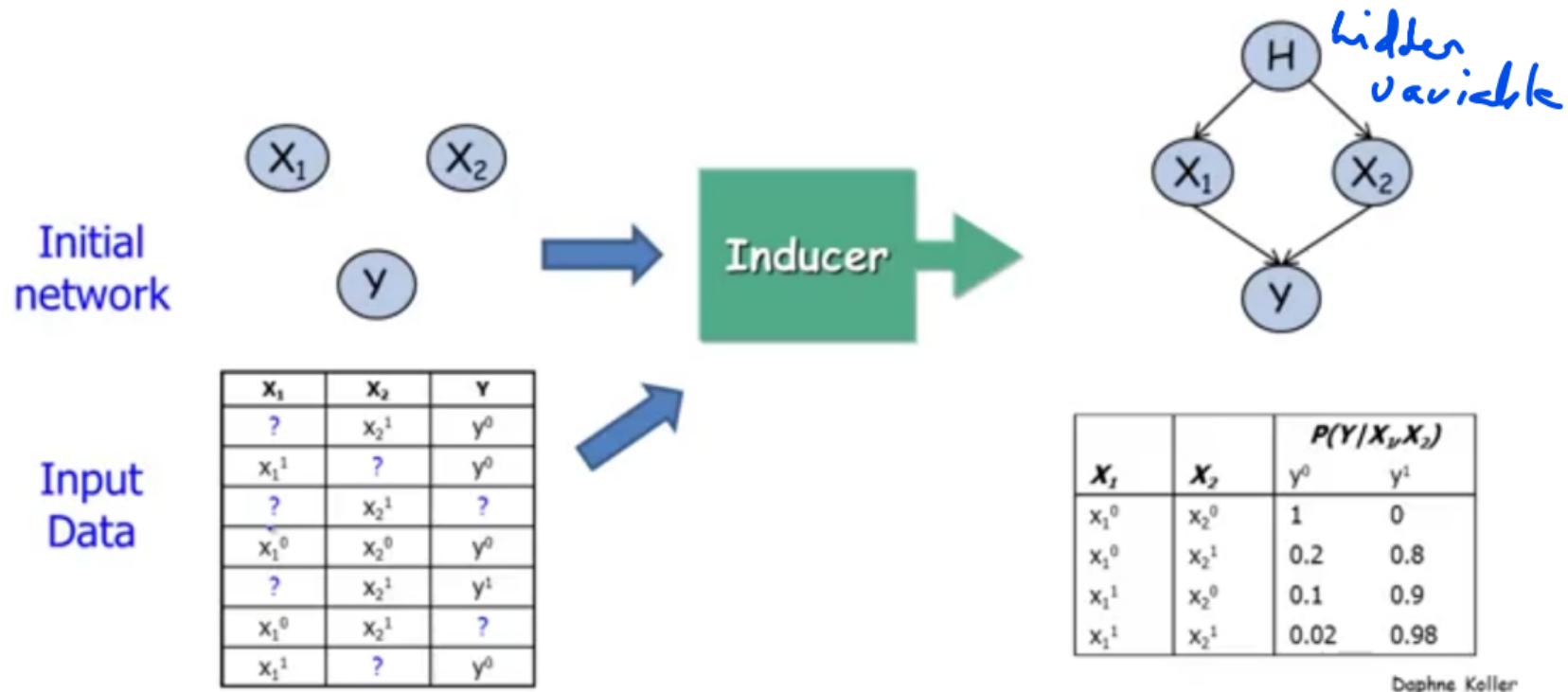
KNOWN STRUCTURE INCOMPLETE DATA



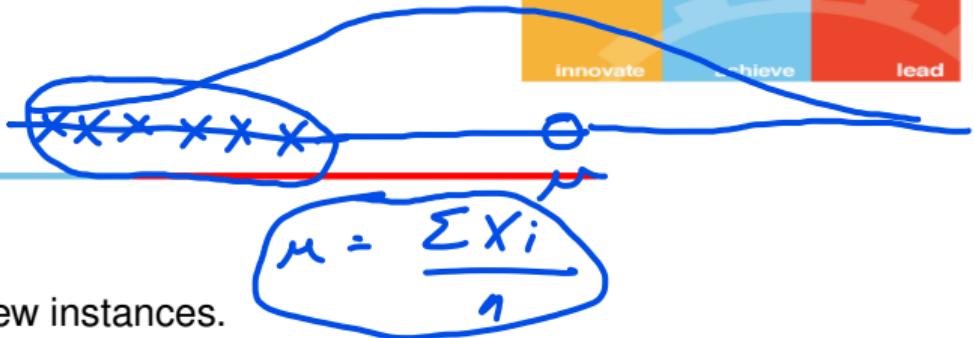
UNKNOWN STRUCTURE INCOMPLETE DATA



LATENT (HIDDEN) VARIABLES INCOMPLETE DATA



GOALS OF LEARNING I



- Density estimation:
 - ▶ Answer general queries about new instances.
 - ▶ Metric: Training set likelihood

$$P(\mathcal{D} : \mathcal{M}) = \prod_m P(d[m] : \mathcal{M})$$

- ▶ Care about new data: generalization performance – Evaluate on test set likelihood $P(\mathcal{D}' : \mathcal{M})$
- ▶ Minimize the expected loss :

$$\mathbb{E}_{\mathcal{D}}[\text{loss}(d : \mathcal{M})] = \frac{1}{|\mathcal{D}|} \sum_{d \in \mathcal{D}} \underline{\text{loss}(d : \mathcal{M})}$$

GOALS OF LEARNING II

- Prediction or Classification Task:

- ▶ Specific prediction task on new instances.
- ▶ Select a MAP assignment to predict a set of variables \mathbf{Y} , given a set of observed variables \mathbf{X} .
$$\arg\max_{\mathbf{Y}} P(\mathbf{Y} | \mathbf{e})$$
- ▶ Eg: Text document classification, Image segmentation, speech recognition
- ▶ Select model to optimize

- ★ likelihood

$$\prod_m P(d[m] : \mathcal{M})$$

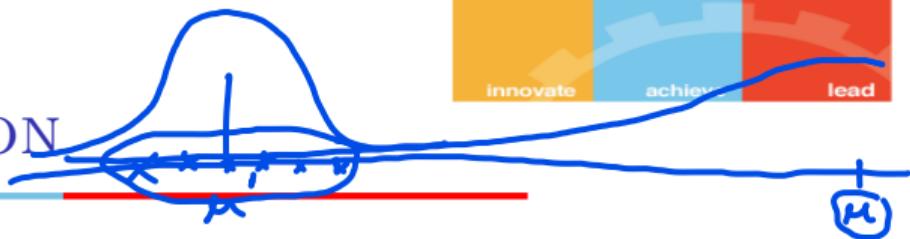
- ★ conditional likelihood

$$\prod_m P(y[m] | x[m] : \mathcal{M})$$

GOALS OF LEARNING III

- Knowledge Discovery of \mathcal{M}^* :
 - ▶ Discover knowledge about P^* .
 - ▶ Distinguish direct and indirect dependencies.
 - ▶ Possibly directionality of edges.
 - ▶ Presence and location of hidden variables.
 - ▶ Reconstruct correct model \mathcal{M}^* .
 - ▶ Measure the success in terms of the model = differences between \mathcal{M}^* and $\tilde{\mathcal{M}}$.

LEARNING AS OPTIMIZATION



- **Hypothesis space** – a set of candidate models.
- **Objective function** – a criterion for quantifying our preference for different models.
- **Learning Task** – find a high-scoring model within the hypothesis space.
- Use data \mathcal{D} to define an **empirical distribution** $\hat{P}_{\mathcal{D}}$

$$\hat{P}_{\mathcal{D}}(A) = \frac{1}{M} \sum_m \mathbf{1}\{d[m] \in A\}$$

→ indicator Function
count

- The probability of the event A is simply the fraction of training examples that satisfy A .
- Use of the empirical log-loss (or log-likelihood) as the objective.
- This type of objective tends to over-fit the learned model to the training data . Use regularization, cross-validation techniques.

GENERATIVE TRAINING

- Perform a particular task such as predicting \mathbf{Y} from \mathbf{X} .
- Goal: get \tilde{M} close to overall joint distribution $P^*(\mathbf{Y}, \mathbf{X})$.
- Model trained to generate all the variables.
- Naive Bayes model
- Higher bias
- Encode independence assumption about feature variables \mathbf{X} .
- Defines $\tilde{P}(\mathbf{Y}, \mathbf{X})$ and induces $\tilde{P}(\mathbf{Y} | \mathbf{X})$ and $\tilde{P}(\mathbf{X})$ using the same overall model for both.
- Works better when learning from limited amounts of data.

*high bias
(lack of expansion)*

$$P(x_1, x_2, \dots, x_n | y) = P(x_1 | y) P(x_2 | y) \dots$$

assumption

DISCRIMINATIVE TRAINING

- Goal: get $\tilde{P}(\mathbf{Y} \mid \mathbf{X})$ to be close to $P^*(\mathbf{Y} \mid \mathbf{X})$
- Undirected model
- Train a conditional random field (CRF)
- Model directly encodes a conditional distribution $P(\mathbf{Y} \mid \mathbf{X})$.
- Encode independence assumptions about \mathbf{Y} and their dependence on \mathbf{X} .
- Find a good fit only to $P^*(\mathbf{Y} \mid \mathbf{X})$ without containing the same model to provide a good fit to $P^*(\mathbf{X})$.

LEARNING TASKS

Input to learning tasks

- Prior knowledge about $\tilde{\mathcal{M}}$
- Set of \mathcal{D} of data instances $\{d[1], \dots, d[M]\}$ which are sample IID from P^* .

Output of learning tasks

- Model $\tilde{\mathcal{M}}$ with structure and parameters.

3 Axes

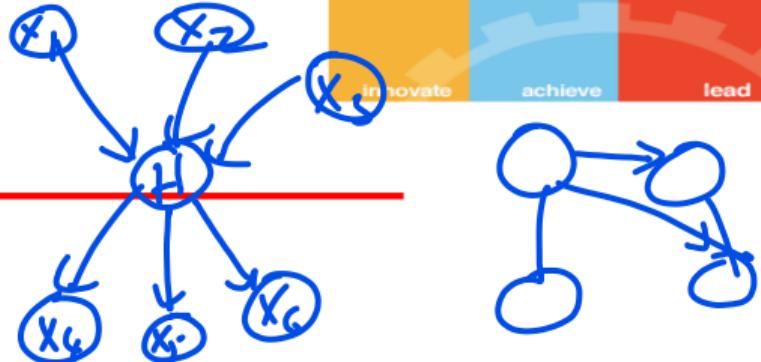
- ① The type of graphical model we are trying to learn – a Bayesian network or a Markov network.
- ② Hypothesis space
- ③ Data observability

LEARNING TASKS

Hypothesis space

- Given a graph structure, and learn only (some of) the parameters.
- We may not know the structure, and we have to learn both parameters and structure from the data.
- We may not even know the complete set of variables over which the distribution P^* is defined. We may only observe some subset of the variables in the domain and possibly be unaware of others.

LEARNING TASKS



Data Observability

- The data are **complete**, or **fully observed**, so that each of the training instances $d[m]$ is a full instantiation to all of the variables in \mathcal{X}^* .
- The data are **incomplete**, or **partially observed**, so that, in each training instance, some variables are not observed.
- The data contain **hidden variables** whose value is never observed in any training instance. The inclusion of a hidden variable in the network can greatly simplify the structure, reducing the complexity of the network that needs to be learned.

TABLE OF CONTENTS

1 LEARNING

2 PARAMETER ESTIMATION IN BAYESIAN NETWORKS

3 MLE FOR BAYESIAN NETWORKS

4 BAYESIAN PARAMETER ESTIMATION

PARAMETER ESTIMATION IN BAYESIAN NETWORKS

- Network structure is fixed.
- Data-set \mathcal{D} consists of fully observed data instances.

$$\mathcal{D} = \{d[1], \dots, d[M]\}$$

- Two approaches
 - ① Maximum likelihood estimation
 - ② Bayesian estimation

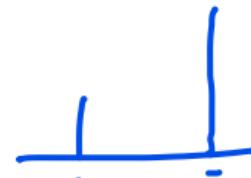
BIASED COIN EXAMPLE

- Head and tails outcome are controlled by a parameter θ .
- θ = frequency of heads in the coin tosses.

$$\max_{\theta} P(D|\theta)$$

$$P(X = H) = \theta$$

$$P(X = T) = 1 - \theta$$



- The distribution P is Bernoulli distribution.
- The data-set D is sampled IID from P .

$$\mathcal{D} = \{x[1], \dots, x[M]\}$$

- ▶ Tosses are independent of each other.
- ▶ Tosses are sampled from the same distribution.
- Goal: Take \mathcal{D} and reconstruct θ .

BIASED COIN EXAMPLE

$x[1]$ and $x[2]$ are independent
given θ

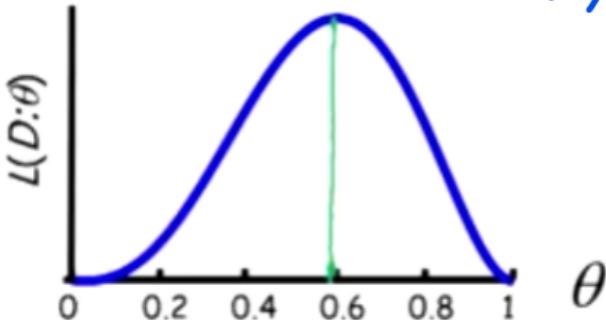
$$P(x[1], x[2] | \theta) = P(x[1] | \theta) P(x[2] | \theta)$$

- Suppose we observe: H, T, T, H, H

$$P(X[1] = H) = \theta$$

$$P(X[2] = T) = (1 - \theta) \quad \text{IID samples}$$

$$\begin{aligned} P(< H, T, T, H, H >: \theta) &= \underline{\theta(1 - \theta)(1 - \theta)\theta\theta} \\ &= \underline{\theta^3(1 - \theta)^2} \end{aligned}$$



$$\max P(\theta_1, \theta_2, \dots, \theta_N | \theta) = \max P(\theta_1 / \theta) I(\theta_2 / \theta) \dots I(\theta_N / \theta)$$

BIASED COIN EXAMPLE

- The probability of the data changes as a function of θ .
- Define the **likelihood function**:

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

- Use the likelihood function as the measure of quality for different parameter values.
- Select the parameter value that maximizes the likelihood; this value is called the **maximum likelihood estimator (MLE)**.
- Observations: M_H heads and M_T tails .
- M_H and M_T are sufficient statistics.

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

$$\hat{\theta} = \frac{M_H}{M_H + M_T} = \frac{3}{5} = 0.6$$

$\max_{\theta} \log \text{likelihood}$
 $\log(\theta^3(1 - \theta)^2)$
 $3\log\theta + 2\log(1 - \theta)$
 $\frac{3}{\theta} + \frac{2}{1 - \theta}(-1) = 0$

BERNOULLI DISTRIBUTION AS PGM

$$P(x[m+1] | D) \neq P(x[m+1] | D, \theta)$$

Distribution

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

$D = x[1], x[2], \dots, x[m]$

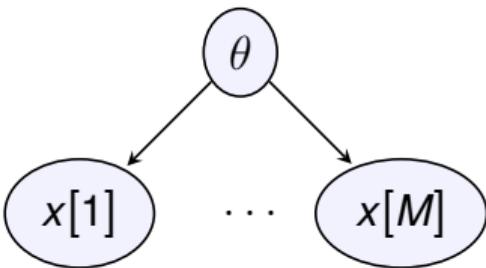
- Hypothesis space

Diagram:

```

graph TD
    theta((θ)) --> X((X))
    subgraph Data_D [Data D]
        X
    end

```



$$P(\theta | D) \neq P(\theta)$$

$\Theta = [0, 1] : \sum_i \theta_i = 1$

- Likelihood function

$$P(x[1], x[2], \dots, x[m] | \theta) = \prod_{i=1}^m P(x[i] | \theta)$$

$$\mathcal{L}(\theta : D) = P(D : \theta) = \prod_{m=1}^M P(x[m] | \theta) = \theta^{M_1} (1 - \theta)^{M_0}$$

BERNOULLI DISTRIBUTION AS PGM

- Observations: M_1 and M_0
- The counts M_1 and M_0 give a compact distribution of the likelihood. These are called **sufficient statistics**.
- M_1 and M_0 are sufficient statistics.
- Find θ that maximizes likelihood:

$$\underline{\mathcal{L}(\theta : \mathcal{D})} = P(\mathcal{D} : \theta) = \underline{\theta^{M_1}} (1 - \theta)^{\underline{M_0}}$$

- Equivalently maximize log-likelihood:

$$\underline{\ell(\theta : \mathcal{D})} = M_1 \log \theta + M_0 \log(1 - \theta)$$

- Differentiate log-likelihood and solve for θ as

$$\hat{\theta} = \frac{M_1}{M_1 + M_0}$$

SUFFICIENT STATISTICS

- Sufficient statistic is a function of the data that summarizes the relevant information for computing the likelihood.

DEFINITION (SUFFICIENT STATISTICS)

A function $\tau(\xi)$ from instances \mathcal{X} to \mathbb{R}^ℓ is a sufficient statistic, if for any two data sets \mathcal{D} and \mathcal{D}' and any $\theta \in \Theta$

If $\sum_{\xi[m] \in \mathcal{D}} \tau(\xi[m]) = \sum_{\xi'[m] \in \mathcal{D}'} \tau(\xi'[m])$
 then $\mathcal{L}(\theta : \mathcal{D}) = \mathcal{L}(\theta : \mathcal{D}')$

\downarrow $\left[\begin{array}{l} \mathcal{D} \in \mathbb{R}^\ell \\ \mathcal{D}' \in \mathbb{R}^\ell \end{array} \right]$
 \downarrow $\left[\begin{array}{l} \tau \in \mathbb{R}^\ell \\ \tau' \in \mathbb{R}^\ell \end{array} \right]$

- The tuple $\sum_{\xi[m] \in \mathcal{D}} \tau(\xi[m])$ is referred to as sufficient statistics of the data-set \mathcal{D} .

MULTINOMIAL DISTRIBUTION AS PGM

- Multinomial variable X takes the values x^1, \dots, x^K .
- Distribution :

$$P(X : \theta) = \theta_k \quad \text{if } x = x^k$$

- Hypothesis space:

$$\Theta = \left\{ \theta \in [0, 1]^K : \sum_i \theta_i = 1 \right\}$$

- Sufficient Statistics:

$$\underline{\tau(x^k)} = \underbrace{0, \dots, 0}_{k-1}, 1, \underbrace{0, \dots, 0}_{n-k}$$

k dimensional vector.

MULTINOMIAL DISTRIBUTION AS PGM

- Likelihood function:

$$\frac{\partial \mathcal{L}}{\partial \theta_1} = 0 \quad \frac{\partial \mathcal{L}}{\partial \theta_2} = 0$$

$$\mathcal{L}(\theta : \mathcal{D}) = \prod_k \theta_k^{M[k]}$$

- Maximum Likelihood Estimation:

$$\theta_1^{M[1]} \theta_2^{M[2]}$$

$$\theta_1^{M[1]} (1 - \theta_1)^{M[2]}$$

$$\frac{\partial \mathcal{L}}{\partial \theta_1} = 0$$

$$\mathcal{L}(\theta : \mathcal{D}) = \max_{\theta \in \Theta} \mathcal{L}(\theta : \mathcal{D})$$

$$\hat{\theta} = \frac{M[k]}{M}$$

MLE SUMMARY

- MLE is a simple principle for estimating or parameter selection given a data-set \mathcal{D} .
- Likelihood function uniquely determined by sufficient statistic that summarize \mathcal{D} .
- MLE has a closed form solution for many parametric distributions.

TABLE OF CONTENTS

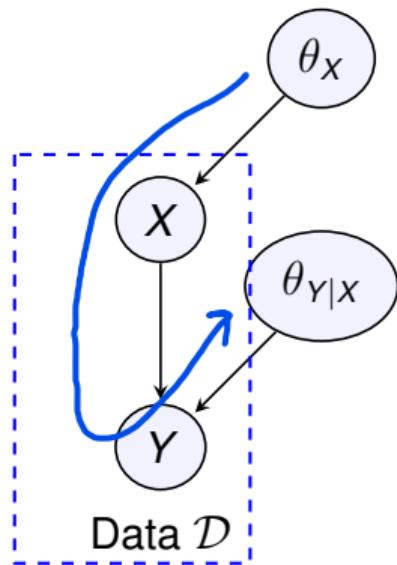
1 LEARNING

2 PARAMETER ESTIMATION IN BAYESIAN NETWORKS

3 MLE FOR BAYESIAN NETWORKS

4 BAYESIAN PARAMETER ESTIMATION

MLE FOR BAYESIAN NETWORKS



- A network consisting of two binary variables.
- Each assignment or training instance is given by $< x[m], y[m] >$.
- The network is parametrized by θ , which defines the set of parameters for all the CPDs in the network.
- Parameters for X : θ_{x^1} and θ_{x^0}
- Parameters for Y : $\theta_{Y|x} = \theta_{Y|x^1} \cup \theta_{Y|x^0}$

$$\theta_{Y|x^1} = \{\theta_{y^1|x^1}, \theta_{y^0|x^1}\}$$

$$\theta_{Y|x^0} = \{\theta_{y^1|x^0}, \theta_{y^0|x^0}\}$$

MLE FOR BAYESIAN NETWORKS

- Goal: Maximize the likelihood function.

$$\begin{aligned}
 \mathcal{L}(\theta : \mathcal{D}) &= \prod_{m=1}^M P(x[m], y[m] : \underline{\theta}) \\
 &= \prod_m \underbrace{P(x[m] : \theta)}_{\text{no dependence on } \theta_{Y|X}} \underbrace{P(y[m] | x[m] : \theta)}_{\theta_{Y|X}} \\
 &= \left[\prod_m \underbrace{P(x[m] : \theta_X)}_{\theta_X} \right] \left[\prod_m \underbrace{P(y[m] | x[m] : \theta_{Y|X})}_{\theta_{Y|X}} \right] \\
 &= \left[\prod_m P(x[m] : \theta_X) \right] \left[\prod_{m: x[m] = x^0} P(y[m] | x[m] : \theta_{Y|x^0}) \right] \left[\prod_{m: x[m] = x^1} P(y[m] | x[m] : \theta_{Y|x^1}) \right]
 \end{aligned}$$

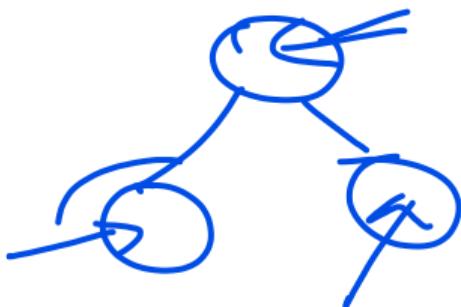
- The likelihood function decomposes into local likelihood terms, each one depends only on the parameters for that variable's CPD.

$$\max_{\mathbf{x}, \mathbf{y}} F(\mathbf{x}, \mathbf{y}) = \max_{\mathbf{x}, \mathbf{y}} g(\mathbf{x})h(\mathbf{y}) = \left[\max_{\mathbf{x}} g(\mathbf{x}) \right] \left[\max_{\mathbf{y}} h(\mathbf{y}) \right]$$

innovate achieve lead

MLE FOR BAYESIAN NETWORKS

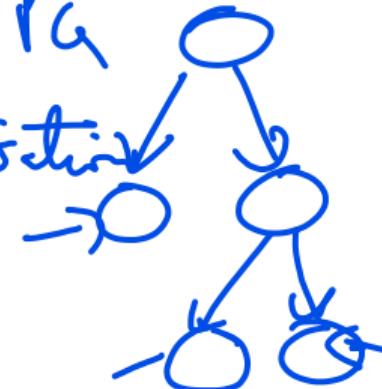
- Goal: Maximize the likelihood function.



$$\begin{aligned}
 \mathcal{L}(\theta : \mathcal{D}) &= \prod_m P_{\mathcal{G}}(x[m], y[m] : \theta) \\
 &= \prod_m \prod_i P(x_i[m] | pa_{X_i}[m] : \theta) \\
 &= \prod_i \left[\prod_m P(x_i[m] | pa_{X_i}[m] : \theta) \right] \\
 &= \prod_i L_i(\theta_{X_i|Pa_{X_i}} : \mathcal{D})
 \end{aligned}$$

all samples

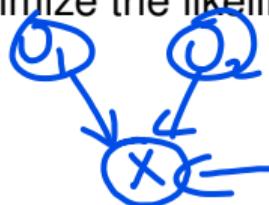
factorization



- When the parameter sets $\theta_{X_i|Pa_{X_i}}$ are disjoint, then MLE can be computed by maximizing each local likelihood separately. The likelihood decomposes as a product of independent terms, one for each CPD in the network.

MLE FOR BAYESIAN NETWORKS FOR TABLE CPDs

- Goal: Maximize the likelihood function.



$$\mathcal{L}_x(\theta_{x|u} : \mathcal{D}) = \prod_m \theta_{x[m]|u[m]}$$

$\dots - \underline{u_1} - \underline{u_2} \times$
 o o o

- MLE parameters

optimal

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

$M[u, x]$ = frequency of

$$\Pr(x=x, u=u) = \Pr(x=x \text{ and } u=u) = \Pr(u=u)$$

$$M[u] = \sum_x M[u, x]$$

MLE FOR BAYESIAN NETWORK SUMMARY

- For Bayesian Network with disjoint sets of parameters in CPDs, likelihood decomposes as a product of local likelihood functions, one per variable.
 - For table CPDs, local likelihood functions further decomposes as a product of likelihood for multinomials, one for each parent combination.
-

TABLE OF CONTENTS

- 1 LEARNING
- 2 PARAMETER ESTIMATION IN BAYESIAN NETWORKS
- 3 MLE FOR BAYESIAN NETWORKS
- 4 BAYESIAN PARAMETER ESTIMATION

DRAWBACK OF MLE

MLE does not distinguish between

- ~~a biased coin and unbiased coin.~~ thumbtack and coin
- 10 tosses and 1,000,000 tosses of the coin.

Another approach – Bayesian Estimation.

JOINT PROBABILISTIC MODEL

X
H H H H . H H H ?
 $P(X = H | D) \neq P(X = H)$

- If θ is unknown, tosses are not marginally independent.
- Each toss tells us something about θ and about the probability of next toss.
- Tosses are conditionally independent given θ .
- Treat θ as a random variable.

$$P(X = H | D, \theta) = P_{\theta}(X = H | \theta)$$

PARAMETER ESTIMATION AS PGM



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

Prior distribution $P(\theta) \in [0, 1]$ Uniform Prior

PARAMETER ESTIMATION AS PGM

- Joint Distribution or Likelihood

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

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

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

$$= P(\theta) \theta^{M[1]} (1 - \theta)^{M[0]}$$

$P(\theta | \theta)$ = likelihood

- Posterior Distribution

$$\underline{P(\theta)} \quad \underline{P(\theta | \theta)} \quad P(\theta | x[1], \dots, x[M]) = \frac{\underline{P(x[1], \dots, x[M] | \theta)} P(\theta)}{\underline{P(x[1], \dots, x[M])}} \rightarrow \text{prior}$$

Posterior \propto likelihood \times prior

handwritten notes:
any

$$P(x[M+1] = x^i | x[1], \dots, x[M]) = \frac{M[1] + 1}{M + 2} \quad \text{vs}$$

$$\frac{M[1]}{M}$$

PARAMETER ESTIMATION AS PGM

- To predict the next value $X[M + 1]$; integrate the posterior over θ .

→ how did we get this?

$$\begin{aligned}
 P(x[M+1] | x[1], \dots, x[M]) &= \int P(x[M+1] | \theta) P(\theta | x[1], \dots, x[M]) d\theta \\
 P(X[M+1] = x^1 | x[1], \dots, x[M]) &= \frac{1}{P(x[1], \dots, x[M])} \int \theta^{x[1]} (1 - \theta)^{M[0]} d\theta \\
 P(x[M+1] = x^i | x[1], \dots, x[M]) &= \frac{M[1] + 1}{M + 2} \quad \text{vs} \quad \frac{M[1]}{M}
 \end{aligned}$$

- As the number of samples grows, the Bayesian estimator and the MLE estimator converge to the same value.

More Detailed Derivation

$$P(x[M+1], D) = \int P(x[M+1], D, \theta) d\theta$$

where $D = x[1], x[2], \dots, x[M]$

$$P(x[M+1], D) = \int P(x[M+1]/D, \theta) p(\theta, \phi) d\theta$$

$$p(\theta) P(x[M+1]/D) = \int P(x[M+1]/D, \phi) p(\phi/D) p(D) d\phi$$

$$P(x \in [M+1] / D) = \int P(x \in [M+1] / \emptyset) P(\emptyset / D) d\emptyset$$

Since $P(x \in [M+1] / D, \emptyset) = P(x \in [M+1] / \emptyset)$

PRIORS AND POSTERIORS

- Observe a training set $\mathcal{D} = x[1] \dots x[M]$
- M IID samples of random variable \mathcal{X} from an unknown distribution $P^*(\mathcal{X})$.
- Assume a parametric model $P(\xi | \theta)$ with a parameter space Θ .
- Treat θ as a random variable.
- Joint distribution

$$P(\mathcal{D}, \theta) = \underline{P(\mathcal{D} | \theta)} \underline{P(\theta)}$$

- The first term is the **likelihood function**. Compactly described by using sufficient statistics.
- The second term is the **prior distribution** over the possible values in Θ . Captures initial uncertainty about the parameters.

PRIORS AND POSTERIORS

- Derive the **posterior probability** over parameters

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

- Marginal likelihood**

$$P(\mathcal{D}) = \int_{\Theta} P(\mathcal{D} | \theta)P(\theta)d\theta$$

PRIORS AND POSTERIORS

- For a multinomial distribution, parameter space Θ is a space of

$$\theta = \langle \theta_1, \dots, \theta_K \rangle$$

$$\sum_k \theta_k = 1$$

- Likelihood function

$$\mathcal{L}(\theta : \mathcal{D}) = \prod_k \theta_k^{M[k]}$$

- Prior is assumed to follow Dirichlet distribution.

DIRICHLET DISTRIBUTION

$$\langle \theta_1, \theta_2, \dots, \theta_K \rangle$$

DEFINITION (DIRICHLET DISTRIBUTION)

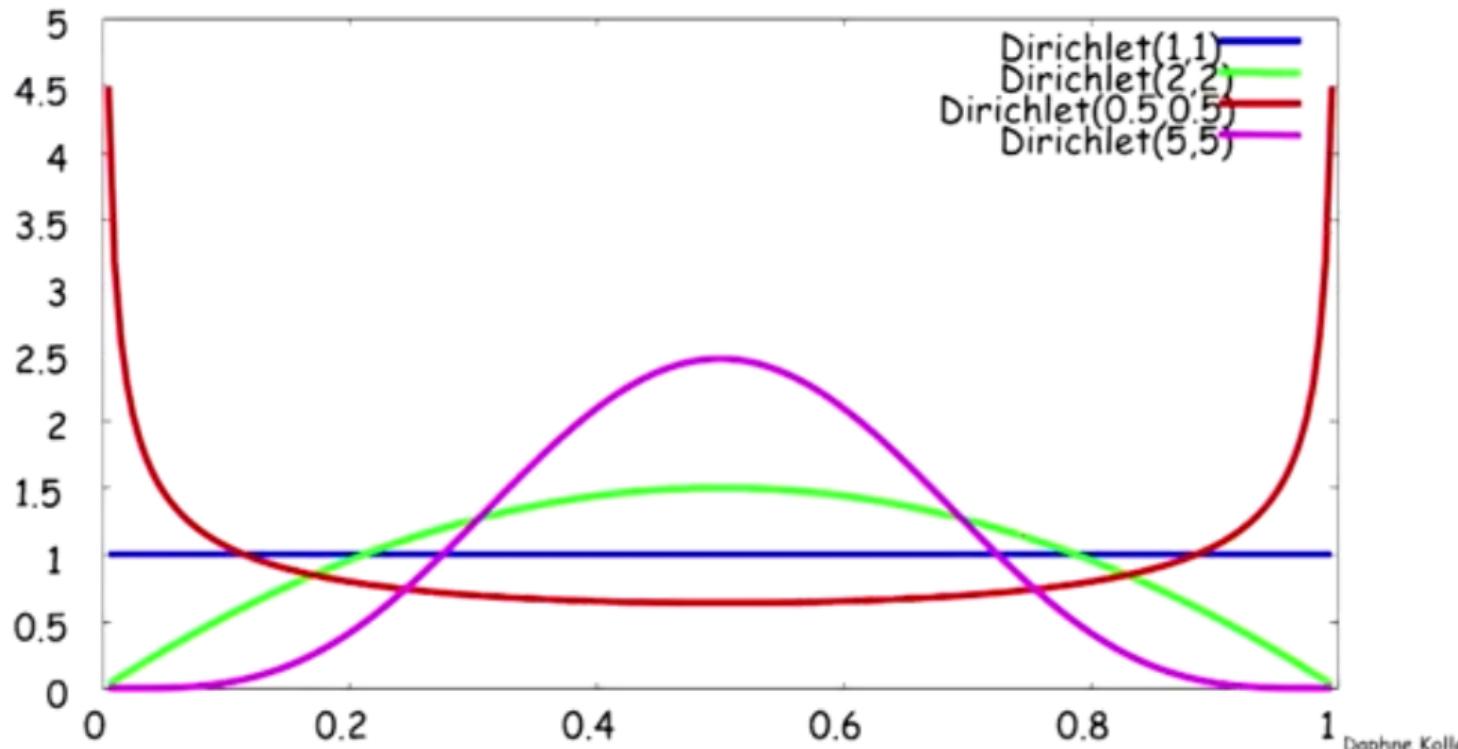
A Dirichlet distribution is specified by a set of hyper-parameters $\langle \underline{\alpha}_1, \dots, \underline{\alpha}_K \rangle$

$$\underline{\theta} \sim Dirichlet(\underline{\alpha}_1, \dots, \underline{\alpha}_K) \quad \text{if} \quad P(\theta) \propto \prod_k \theta_k^{\underline{\alpha}_k - 1}$$

- If $P(\theta)$ is Dirichlet distribution with hyper-parameters $\langle \underline{\alpha}_1, \dots, \underline{\alpha}_K \rangle$ and $\underline{\alpha} = \sum_j \underline{\alpha}_j$, then

$$\mathbb{E}[\theta_k] = \frac{\alpha_k}{\alpha}$$

DIRICHLET DISTRIBUTION



DIRICHLET PRIORS AND POSTERIORS

- Posterior distribution

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

Diagram annotations: The term $P(\theta | \mathcal{D})$ is circled and labeled "Dirichlet". The term $P(\mathcal{D} | \theta) P(\theta)$ is circled and labeled "multinomial". The entire equation is circled and labeled "Dirichlet".

- Likelihood function – multinomial

$$\underline{P(\mathcal{D} | \theta)} = \mathcal{L}(\theta : \mathcal{D}) = \prod_k \theta_k^{M[k]}$$

- Prior distribution – Dirichlet

$$P(\theta) \propto \prod_k \theta_k^{\alpha_k - 1}$$

- Posterior distribution – Dirichlet

DIRICHLET PRIORS AND POSTERIORS

- If Prior $P(\theta)$ is Dirichlet and the Likelihood $P(\mathcal{D} | \theta)$ is multinomial, then the Posterior $P(\theta|\mathcal{D})$ is also Dirichlet.

Prior = Dirichlet($\underline{\alpha_1}, \dots, \underline{\alpha_K}$)

Data counts = $\underline{M_1}, \dots, \underline{M_K}$

Posterior = Dirichlet($\underline{\alpha_1 + M_1}, \dots, \underline{\alpha_K + M_K}$)

$$\int \frac{P(\theta|\mathcal{D})}{(P(\theta|\mathcal{D}))} \left(\prod_{k=1}^K \theta_k^{M_k} \right) (\underline{\theta_1}, \underline{\theta_2}, \dots, \underline{\theta_K})$$

- Prior and Posterior have the same form of distribution.
- Dirichlet is a conjugate pair for multinomial.

BAYESIAN PREDICTION

- To predict for a new sample

$$P(x[M+1]) = \frac{\int g(x) f(x) dx}{\int g(x) dx}$$

$$\mathbb{E}(g(x)) = \int g(x) f(x) dx$$

$$P(x[M+1] | \mathcal{D}) = \int P(x[M+1] | \mathcal{D}, \theta) P(\theta | \mathcal{D}) d\theta$$

$$= \int P(x[M+1] | \theta) P(\theta | \mathcal{D}) d\theta$$

$$= \mathbb{E}_{P(\theta | \mathcal{D})} [P(x[M+1] | \theta)]$$

$$P(x[M+1] = x^k | \mathcal{D}) = \frac{M[k] + \alpha_k}{M + \alpha}$$

Θ_k

$$\mathbb{E}[\Theta_k] = \frac{\alpha_k}{\alpha}$$

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

$$\mathbb{E}[\Theta_k] =$$

EXAMPLE

- For a given binomial data with uniform distribution for parameter θ , it is observed that

$$(M[1], M[0]) = (5, 2) \quad \text{uniform} = \text{Dirichlet}(1, 1)$$

Predict $P(X[8] = 1)$ using MLE and Bayesian prediction.

- MLE

$$P(X[8] = 1) = \frac{M_1}{M} = \frac{5}{7} = 0.71$$

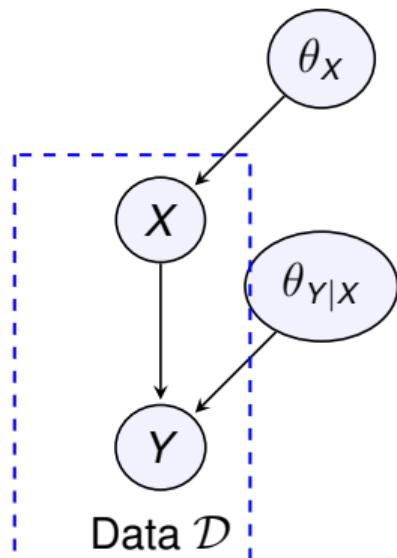
- Bayesian prediction

$$\begin{aligned} P(X[8] = 1) &= \frac{\alpha_1 + M_1}{\alpha + M} \\ &= \frac{1 + 5}{2 + 7} = \frac{6}{9} = 0.66 \end{aligned}$$

BAYESIAN ESTIMATION SUMMARY

- Bayesian Learning treats parameters as random variables.
- Dirichlet distribution as conjugate pair of multinomial distribution.
 - ▶ Posterior has the same form as prior.
 - ▶ Can be updated in closed form using sufficient statistic from data.
- Bayesian prediction combines sufficient statistic from imaginary Dirichlet sample and real data samples.
- Dirichlet hyper-parameters determine both the prior beliefs and their strengths.
- Bayesian Learning is robust in sparse data regime in terms of its generalization ability.

BAYESIAN ESTIMATION ON BAYESIAN NETWORK



- A network consisting of two binary variables.
- Training data consists of M observations given by $\langle x[m], y[m] \rangle$
- Unknown parameters θ_X and $\theta_{Y|X}$
- Instances are independent given the unknown parameters.
- $\langle x[m], y[m] \rangle$ are d-separated from $\langle x[m'], y[m'] \rangle$ once we know the parameter variables.
- Assume that the priors for the individual parameters variables are apriori independent. That is, we believe that knowing the value of one parameter tells us nothing about another.

BAYESIAN ESTIMATION ON BAYESIAN NETWORK

- The Bayesian network have parameters

$$\theta = (\theta_{X_1|Pa_{X_1}}, \dots, \theta_{X_n|Pa_{X_n}})$$

- Prior distribution satisfies **global parameter independence** if

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

- If the complete data $\langle x[m], y[m] \rangle$ are observed for all m , then parameters θ_X and $\theta_{Y|X}$ are d-separated.

$$P(\theta_X, \theta_{Y|X} | \mathcal{D}) = P(\theta_X | \mathcal{D})P(\theta_{Y|X} | \mathcal{D})$$

- Given the data set \mathcal{D} , determine the posterior over θ_X independently of the posterior over $\theta_{Y|X}$.

BAYESIAN ESTIMATION ON BAYESIAN NETWORK

- Let X have parents U . Then the prior $P(\theta_{X|U})$ satisfies **local parameter independence** if

$$P(\theta_{X|U}) = \prod_u P(\theta_{X|u})$$

- If $P(\theta)$ satisfies global and local parameter independence then,

$$P(\theta | \mathcal{D}) = \prod_i \prod_{Pa_{X_i}} P(\theta_{X_i|Pa_{X_i}} | \mathcal{D})$$

BAYESIAN ESTIMATION ON BAYESIAN NETWORK

- Multinomial Parameters $\theta_{X|u}$
- $P(\theta_{X|u})$
 - ▶ Dirichlet prior with hyper-parameters $\alpha_{x^1|u}, \dots, \alpha_{x^K|u}$
- $P(\theta_{X|u} | \mathcal{D})$
 - ▶ Dirichlet posterior with hyper-parameters $\alpha_{x^1|u} + M[u, x^1], \dots, \alpha_{x^K|u} + M[u, x^K]$

SUMMARY

- In Bayesian networks, if parameters are independent apriori, then they are also independent in the posterior.
- For multinomial Bayesian networks, estimation uses sufficient statistic $M[x, u]$ (counts).

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

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

- Bayesian Learning requires a choice of prior.

EXAMPLE

Given the following data and the structure that $X \rightarrow Y$, learn the parameters using MLE.

X	Y
0	1
0	1
1	0
1	1
1	0
1	1
1	0
0	1
0	1

$\hat{\theta}_{X|u} = \frac{M[X, u]}{M[u]}$
 $\theta_X = P(X) = \begin{bmatrix} 4/9 & 5/9 \end{bmatrix}$
 $\theta_{Y|X} = P(Y | X) = \begin{bmatrix} 0/4 & 4/4 \\ 3/5 & 2/5 \end{bmatrix}$