

**Lecture 11: Bayesian  
Networks – More on  
Inference, Maybe Learning**

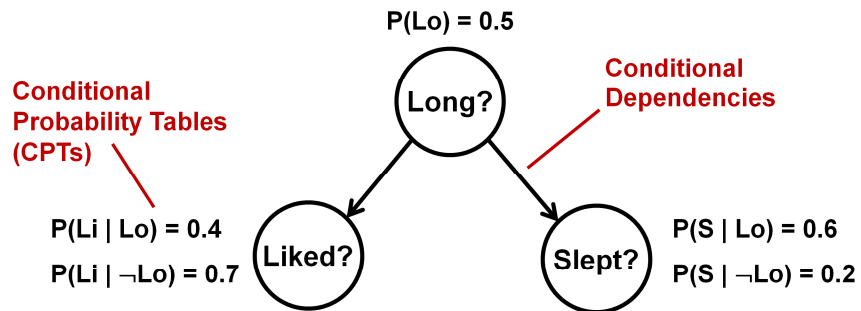
- **Homework 2 due NOW!**
- **Homework 3 out this evening**
- **Due MONDAY, Oct 12<sup>th</sup>**

**(inspect HW3)**

**Lecture 11: Bayesian  
Networks – More on  
Inference, Maybe Learning**

# Bayesian Networks

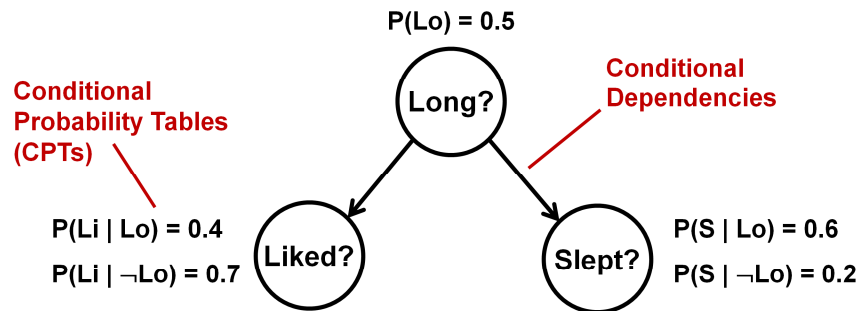
Bayesian networks are directed acyclic graphs with nodes representing random variables and edges representing dependency assumptions



# Bayesian Networks

Joint distribution (factorization):

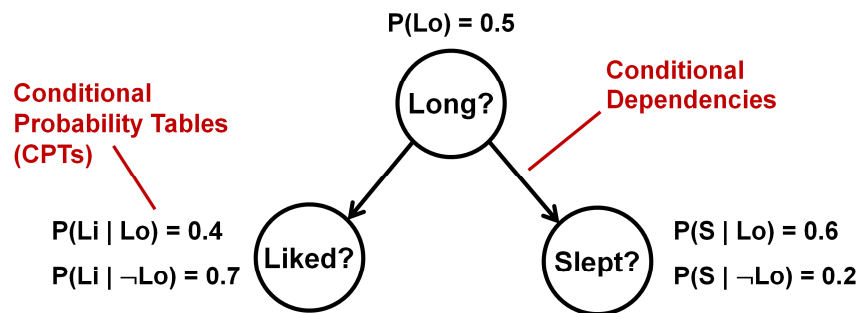
$$P(\text{Lo}, \text{Li}, \text{S}) = P(\text{Lo}) P(\text{Li} \mid \text{Lo}) P(\text{S})$$



# Bayesian Networks

Conditional independence:

$P(Li | Lo, S)?$



# Bayesian Networks

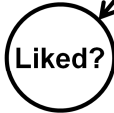
Conditional independence:

$$\begin{aligned} P(\text{Li} \mid \text{Lo}, \text{S}) &= P(\text{Li}, \text{Lo}, \text{S}) / P(\text{Lo}, \text{S}) \\ &= P(\text{Li} \mid \text{Lo}) P(\text{Lo}) P(\text{S}) / P(\text{Lo}) P(\text{S} \mid \text{Lo}) \\ &= P(\text{Li} \mid \text{Lo}) \end{aligned}$$

$$\text{Li} \perp \text{S} \mid \text{Lo}$$

Conditional  
Probability Tables  
(CPTs)

$$\begin{aligned} P(\text{Li} \mid \text{Lo}) &= 0.4 \\ P(\text{Li} \mid \neg \text{Lo}) &= 0.7 \end{aligned}$$



$$P(\text{Lo}) = 0.5$$

Long?

Conditional  
Dependencies



$$\begin{aligned} P(\text{S} \mid \text{Lo}) &= 0.6 \\ P(\text{S} \mid \neg \text{Lo}) &= 0.2 \end{aligned}$$

## **Bayesian networks: Inference**

- **Algorithms for inferring the values of unobserved variables.**
- **Last time: Sampling**

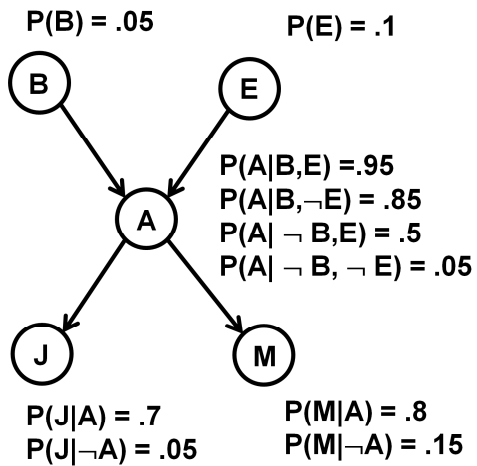


# Stochastic Inference

We can easily sample the joint distribution to obtain possible instances:

1. Sample the free variables
2. For every other variable: If all parents have been sampled, sample based on conditional distribution

**We end up with a new set of assignments for B,E,A,J and M which are a random sample from the joint**



## Weighted Sampling for Computing

**Problem: What if the condition rarely happens?**

- Set  $N_B, N_c = 0$
- Repeat:
  - Sample the joint setting the values for  $J$  and  $M$ , compute the weight,  $w$ , of this sample
  - $N_c = N_c + w$
  - If  $B = 1$ ,  $N_B = N_B + w$
- After many iterations, set:  
$$P(B | J, \neg M) = N_B / N_c$$

**We would need lots and lots of samples, and most would be wasted**

## **Bayesian networks: Inference**

- **Algorithms for inferring the values of unobserved variables.**
- **Last time: Sampling**
  - fast, (often) approximate
- **Last time: Exact inference**

# Inference

We are interested in queries of the form:  
 $P(B \mid J, \neg M)$

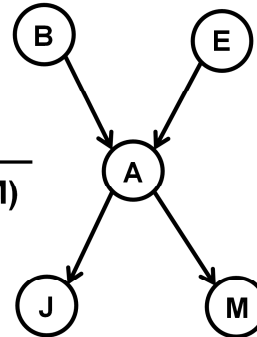
This can also be written as a joint:

$$P(B \mid J, \neg M) = \frac{P(B, J, \neg M)}{P(B, J, \neg M) + P(\neg B, J, \neg M)}$$

How do we compute the new joint?

$$P(B, J, \neg M) = \sum_a \sum_e P(B, J, \neg M, a, e)$$

Sum all probabilities with these settings (B, J,  $\neg M$ ): the sum is over the possible assignments to the other two variables, E and A)



# Computing: $P(B, J, \neg M)$

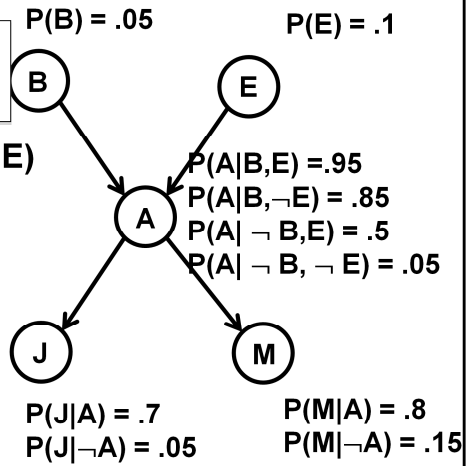
$$P(B, J, \neg M) = \sum_a \sum_e P(B, J, \neg M, a, e)$$

$$= P(B, J, \neg M, A, E) + P(B, J, \neg M, \neg A, E) \\ + P(B, J, \neg M, A, \neg E) + \\ P(B, J, \neg M, \neg A, \neg E)$$

$$= 0.0007 + 0.00001 + 0.005 + 0.0003$$

$$= 0.00601$$

**How can we reuse computations?**



# Computing: $P(B, J, \neg M)$

$$P(B, J, \neg M) = \sum_a \sum_e P(B, J, \neg M, a, e)$$

$P(B) = .05$

$P(E) = .1$

```

graph TD
    B((B)) --> A((A))
    E((E)) --> A
    A --> M((M))
    A --> J((J))
    style J stroke-dasharray: 5 5
    
```

$P(A|B, E) = .95$   
 $P(A|B, \neg E) = .85$   
 $P(A|\neg B, E) = .5$   
 $P(A|\neg B, \neg E) = .05$

$P(J|A) = .7$   
 $P(J|\neg A) = .05$

$P(M|A) = .8$   
 $P(M|\neg A) = .15$

$$= P(B, J, \neg M, A, E) + P(B, J, \neg M, \neg A, E)$$

$$+ P(B, J, \neg M, A, \neg E) + P(B, J, \neg M, \neg A, \neg E)$$

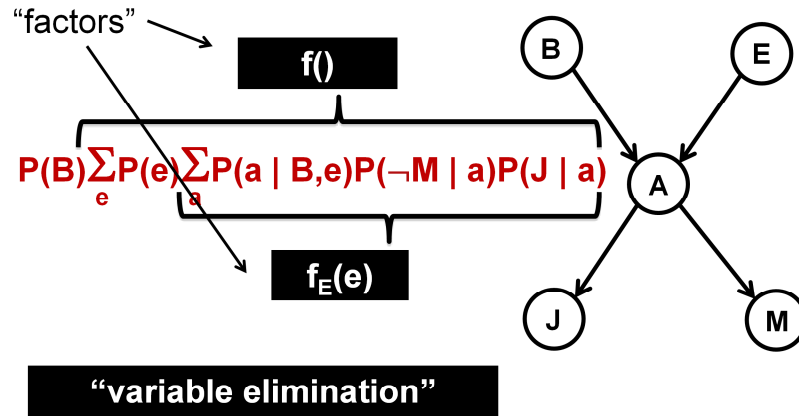
$$= \sum_a \sum_e P(B)P(e)P(a | B, e)P(\neg M | a)P(J | a)$$

$$= P(B) \sum_e P(e) \sum_a P(a | B, e)P(\neg M | a)P(J | a)$$

“sum out” A, and store as a function of e and use whenever necessary (no need to recompute for each value of e)

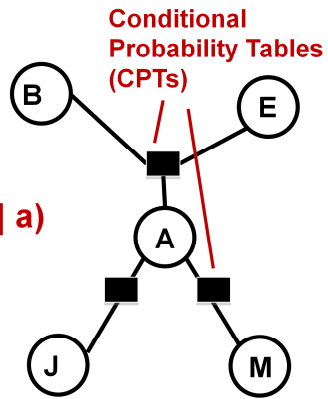
Instead of computing the value for every value of e

# Computing: $P(B, J, \neg M)$



# Computing: $P(B, J, \neg M)$

$$P(B) \sum_e P(e) \sum_a P(a | B, e) P(\neg M | a) P(J | a)$$



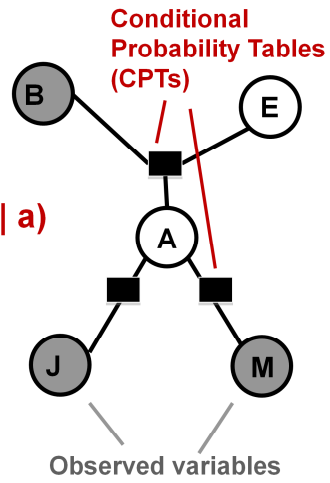
**“variable elimination”**



# Computing: $P(B, J, \neg M)$

$$P(B) \sum_e P(e) \sum_a P(a | B, e) P(\neg M | a) P(J | a)$$

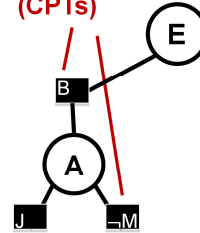
**“variable elimination”**



# Computing: $P(B, J, \neg M)$

$$P(B) \sum_e P(e) \sum_a P(a | B, e) P(\neg M | a) P(J | a)$$

Conditional  
Probability Tables  
(CPTs)



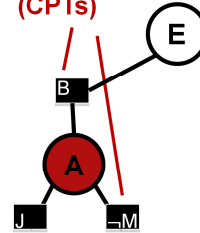
Observed variables  
don't need to be  
eliminated (summed  
out)

**“variable elimination”**

# Computing: $P(B, J, \neg M)$

$$P(B) \sum_e P(e) \underbrace{\sum_a P(a | B, e) P(\neg M | a) P(J | a)}_{f_E(e)}$$

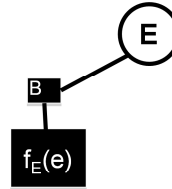
Conditional  
Probability Tables  
(CPTs)



**“variable elimination”**

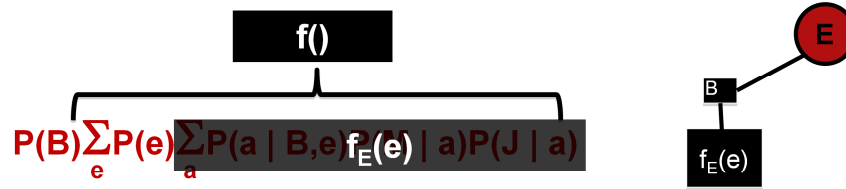
## Computing: $P(B, J, \neg M)$

$$P(B) \sum_e P(e) \sum_a P(a | B, e) f_E(e | a) P(J | a)$$



“variable elimination”

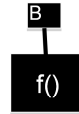
# Computing: $P(B, J, \neg M)$



“variable elimination”

## Computing: $P(B, J, \neg M)$

$$P(B) \sum_e P(e) \sum_a P(a | B, f()) P(M | a) P(J | a)$$



**“variable elimination”**

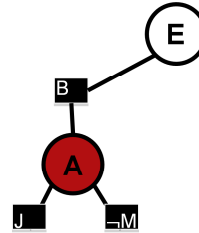
# Actually Computing $f_E(e)$

$$f_A(a) = \begin{pmatrix} P(\neg M|A) P(J|A) \\ P(\neg M|\neg A) P(J|\neg A) \end{pmatrix}$$

$$P(B) \sum_e P(e) \sum_a P(a | B, e) P(\neg M | a) P(J | a)$$

$f_E(e)$

$$f_E(e) = \begin{pmatrix} P(A|B, E) f_A(A) \\ + P(\neg A|B, E) f_A(\neg A) \\ \hline P(A|B, \neg E) f_A(A) \\ + P(\neg A|B, \neg E) f_A(A) \end{pmatrix}$$



Reuse,  
don't recompute!

btw, we computed  $P(B, J, \neg M)$ ,  
but wanted  $P(B|J, \neg M)$

$$P(B | J, \neg M) = \frac{P(B, J, \neg M)}{P(B, J, \neg M) + P(\neg B, J, \neg M)}$$

“normalization”

Also need to compute,  
but can reuse some  
computation again!



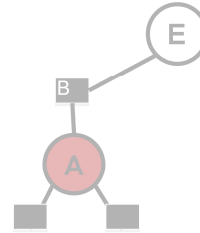
# Actually Computing $f_E(e)$

$$f_A(a) = \begin{pmatrix} P(\neg M|A) P(J|A) \\ P(\neg M|\neg A) P(J|\neg A) \end{pmatrix}$$

$$P(B) \sum_e P(e) \sum_a P(a | B, e) P(\neg M | a) P(J | a)$$

$f_E(e)$

$$f_E(e) = \begin{pmatrix} P(A|B, E) f_A(A) \\ + P(\neg A|B, E) f_A(\neg A) \\ \hline P(A|B, \neg E) f_A(A) \\ + P(\neg A|B, \neg E) f_A(\neg A) \end{pmatrix}$$



Reuse,  
don't recompute!

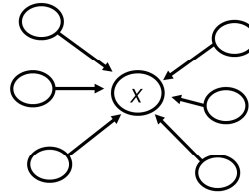
# Algorithm

- **e** - evidence (the variables that are observed)
- **vars** - the conditional probabilities derived from the network in reverse order (bottom up)
- For each *var* in *vars*
  - *factors*  $\leftarrow$  *make\_factor* (*var*,*e*)
  - if *var* is a hidden variable then create a new factor by summing out *var*
- Compute the product of all factors
- Normalize

# Computational Complexity

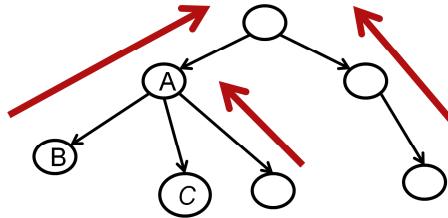
- We can reuse computations to reduce the running time
- However, there are still cases in which this algorithm will lead to exponential running time.
  - Exact Bayesian Inference is NP-Hard
- Consider the case of  $f_x(y_1 \dots y_n)$ . When factoring  $x$  out we would need to account for all possible values of the  $y$ 's.
- e.g. binary:

$$f_x(y_1 \dots y_n) = \begin{pmatrix} f_x(0, \dots 0) \\ \dots \\ f_x(1, \dots 1) \end{pmatrix} \leftarrow 2^n \text{ values}$$



# Computational Complexity

- We can reuse computations to reduce the running time
- However, there are still cases in which this algorithm will lead to exponential running time.
  - Exact Bayesian Inference is NP-Hard
- Easy on trees:



$\sum_B P(B|A) \rightarrow f_1(A)$

$\sum_C P(C|A) \rightarrow f_2(A)$

→ never get functions (factors) with more than 1 argument (size 2)

## **Bayesian networks: Inference**

- **Algorithms for inferring the values of unobserved variables.**
- **Last time: Sampling**
  - fast, (often) approximate
- **Last time: Exact inference**
  - variable elimination
- **Also: “belief propagation”, “variational inference”**

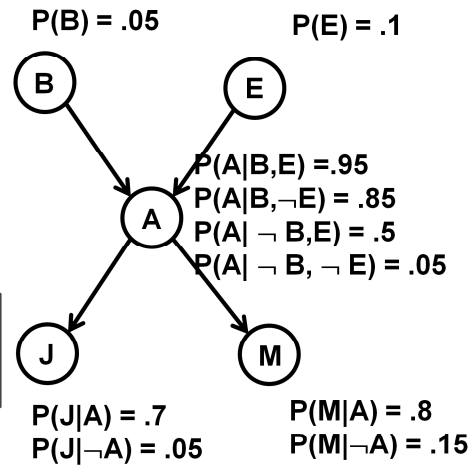
BP on trees = variable elimination

General DAGs need to be

## Inference: compute probabilities from CPTs

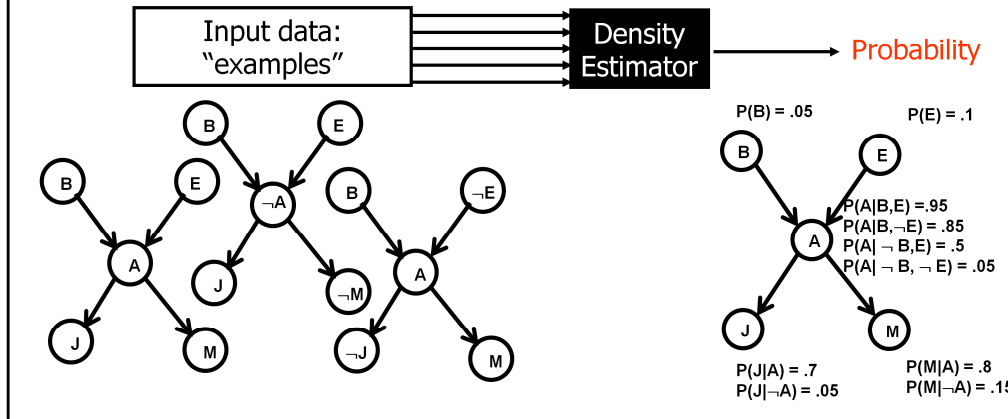
But where do we get them?

Density estimation  
“learning” parameters



# Density Estimation

- A Density Estimator learns a mapping from a set of variables to a Probability, e.g. CPTs



# Density estimation

- **Binary and discrete variables:**

Easy: Just count!

- **Continuous variables:**

Harder (but just a bit): Fit a model



## Learning a density estimator

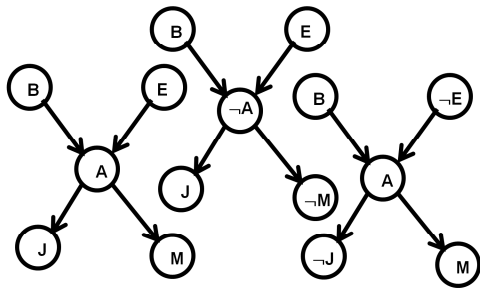
$$\hat{P}(y_i = u) = \frac{\text{\small a variable} \# \text{ examples in which } y_i = u}{\text{total number of examples}}$$

A trivial learning algorithm!

# Learning a density estimator

a variable

$$\hat{P}(y_i = u) = \frac{\text{\# examples in which } y_i = u}{\text{total number of examples}}$$

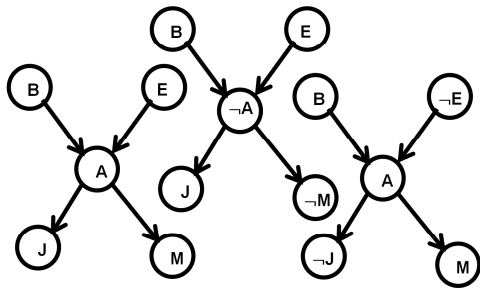


P(B)=  
P(E)=  
P(A|B,E)=  
...

# Learning a density estimator

a variable

$$\hat{P}(y_i = u) = \frac{\text{\# examples in which } y_i = u}{\text{total number of examples}}$$



$P(B)=1$   
 $P(E)=2/3$   
 $P(A|B,E)=1/2$   
...

# Maximum Likelihood Principle

$$\hat{P}(\text{dataset}|M) = \hat{P}(\mathbf{x}_1 \wedge \mathbf{x}_2 \dots \wedge \mathbf{x}_R|M) = \prod_{k=1}^R \hat{P}(\mathbf{x}_k|M)$$

examples

Model: CPTs, net structure, gaussian parameters, ...

- Fit models by maximizing the probability of generating the observed samples:

$$L(x_1, \dots, x_n | \theta) = p(x_1 | \theta) \dots p(x_n | \theta)$$

e.g. "joint probability" from a CPT

- The examples are assumed to be independent
- For a binary random variable A with  $P(A=1)=q$   
 $\text{argmax}_q \text{Likelihood} = \#(A=1)/\#\text{examples}$
- Why?

# Maximum Likelihood Principle

- For a binary random variable A with  $P(A=1)=q$   
 $\operatorname{argmax}_q \text{Likelihood} = \#(A=1)/\#\text{examples}$
- Why?

Data likelihood:  $P(D|q) = q^{n_1} (1-q)^{n_2}$

$n_1$ : #examples w/ A=1  
 $n_2$ : #examples w/ A=0

We would like to find:  $\operatorname{argmax}_q q^{n_1} (1-q)^{n_2}$

**How?**

# Maximum Likelihood Principle

Data likelihood:  $P(D|q) = q^{n_1}(1-q)^{n_2}$

We would like to find:  $\arg \max_q q^{n_1}(1-q)^{n_2}$

$$\frac{\partial}{\partial q} q^{n_1}(1-q)^{n_2} = n_1 q^{n_1-1}(1-q)^{n_2} - q^{n_1} n_2 (1-q)^{n_2-1}$$

$$\frac{\partial}{\partial q} = 0 \Rightarrow$$

$$n_1 q^{n_1-1}(1-q)^{n_2} - q^{n_1} n_2 (1-q)^{n_2-1} = 0 \Rightarrow$$

$$q^{n_1-1}(1-q)^{n_2-1}(n_1(1-q) - q n_2) = 0 \Rightarrow$$

$$n_1(1-q) - q n_2 = 0 \Rightarrow$$

$$n_1 = n_1 q + n_2 q \Rightarrow$$

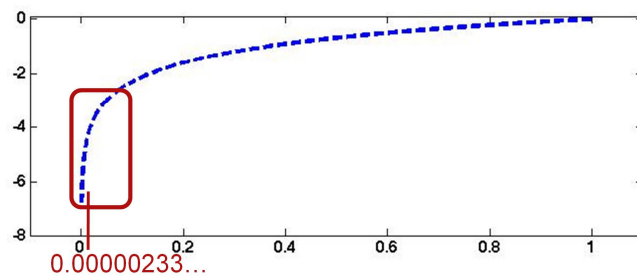
$$q = \frac{n_1}{n_1 + n_2}$$

# Log Probabilities

When working with products, probabilities of entire datasets often get too small. A possible solution is to use the log of probabilities, often termed 'log likelihood'

$$\log \hat{P}(\text{dataset}|M) = \log \prod_{k=1}^R \hat{P}(\mathbf{x}_k|M) = \sum_{k=1}^R \log \hat{P}(\mathbf{x}_k|M)$$

Log values  
between 0 and  
1



Maximize that!

# Density estimation

- **Binary and discrete variables:**

Easy: Just count!

- **Continuous variables:**

Harder (but just a bit): Fit a model

But what if we only have very few samples?



# The danger of joint density estimation

$P(\text{summer} \ \& \ \text{size} \geq 20 \ \& \ \text{evaluation} = 3) = 0$

- No such example in our dataset

Now lets assume we are given a new ("test") dataset. If this dataset contains

Summer	Size	Evaluation
1	30	3

Then the probability we would assign to the *entire* dataset is 0

Summer?	Size	Evaluation
1	19	3
1	17	3
0	49	2
0	33	1
0	55	3
1	20	1

## Naïve Density Estimation

The problem with the Joint Estimator is that it just mirrors the training data.

We need something which generalizes more usefully.

The naïve model generalizes strongly:

Assume that each attribute is distributed independently of any of the other attributes.

# Joint estimation, revisited

Assuming independence we can compute each probability independently

$$P(\text{Summer}) = \frac{1}{2} = 0.5$$

$$P(\text{Evaluation} = 1) = \frac{1}{3} = 0.33$$

$$P(\text{Size} \geq 20) = \frac{2}{3} = 0.66$$

How do we do on the joint?

$$P(\text{Summer} \& \text{Evaluation} = 1) = \frac{1}{6}$$

$$P(\text{Summer})P(\text{Evaluation} = 1) = \frac{1}{2} * \frac{1}{3} = \frac{1}{6}$$

$$P(\text{size} \geq 20 \& \text{Evaluation} = 1) = \frac{1}{3} = 0.33$$

$$P(\text{size} \geq 20)P(\text{Evaluation} = 1) = \frac{2}{3} * \frac{1}{3} = 0.22$$

Summer?	Size	Evaluation
1	19	3
1	17	3
0	49	2
0	33	1
0	55	3
1	20	1

Okay

# Joint estimation, revisited

Assuming independence we can compute each probability independently

$$P(\text{Summer}) = \frac{1}{2} = 0.5$$

$$P(\text{Evaluation} = 1) = \frac{1}{3} = 0.33$$

$$P(\text{Size} \geq 20) = \frac{2}{3} = 0.66$$

How do we do on the joint?

$$P(\text{Summer} \ \& \ \text{Size} \geq 20) = \frac{1}{6} = 0.16667$$

$$P(\text{Summer})P(\text{Size} \geq 20) = \frac{1}{2} * \frac{2}{3} = \frac{1}{3} = 0.333$$

Summer?	Size	Evaluation
1	19	3
1	17	3
0	49	2
0	33	1
0	55	3
1	20	1

**We must be careful when using the Naïve density estimator**

# Contrast

Joint DE	Naïve DE
Can model anything	Can model only very boring distributions
No problem to model "C is a noisy copy of A"	Outside Naïve's scope
Given 100 records and more than 6 Boolean attributes will screw up badly	Given 100 records and 10,000 multivalued attributes will be fine

# Naïve Density Estimation

The problem with the Joint Estimator is that it just mirrors the training data.

We need something which generalizes more usefully.

Joint estimator:  
 $2^n - 1$  parameters

Naïve estimator:  
 $n$  parameters

The naïve model generalizes strongly:

Assume that each attribute is distributed independently of any of the other attributes.

## another way to deal with small datasets

- We just discussed one possibility: Naïve estimation
- Assume we want to compute the probability of heads in a coin flip (50/50)
  - What if we can only observe 3 flips?

1 1 1 0 0 0  
1 1 0 1 0 1  
0 1 1 1 0 0  
0 1 0 0 0 1



- 25% of the times a maximum likelihood estimator will assign probability of 1 to either the heads or tails

## Pseudo counts

- Use *prior belief* about the 'fairness' of most coins to influence the resulting model.
- We assume that we have "observed" 10 flips with 5 tails and 5 heads
- Thus  $P(\text{heads}) = (\#\text{heads}+5) / (\#\text{flips}+10)$
- Advantages: 1. Never assign a probability of 0 to an event  
2. As more data accumulates we can get very close to the real distribution (the impact of the pseudo counts will diminish rapidly)



# Pseudo counts

- Use *prior belief* about the 'fairness' of most coins to influence the estimate
- We assume 5 heads and 5 tails
- Thus P(5 heads) = P(5 tails)
- Advantages:
  1. As more data is collected, the influence of the prior distribution rapidly diminishes
  2. As more data is collected, the influence of the prior distribution rapidly diminishes

Sometimes you can even justify this by incorporating a *real* distribution into your model!

# Lets go back to Naïve vs.full model

What should I use?

This can be determined based on:

- Training data size
- Cross validation
- Likelihood ratio test

Cross validation is one of the most useful tricks in model fitting

Divide up data set into  $m$  parts, train on  $m-1$ , test on the 1 (do  $m$  times)

Statistically valid!

→ Which model does better?

## **Important points**

- **Showing conditional independence**
- **Inference: sampling & exact (variable elimination)**
  
- **Maximum likelihood estimation (MLE)**
- **Pseudo counts**
- **Cross-validation**