

Conditional Independence

Sargur Srihari

srihari@cedar.buffalo.edu

Conditional Independence Topics

1. What is Conditional Independence?
Factorization of probability distribution into marginals
2. Why is it important in Machine Learning?
3. Conditional independence from graphical models
4. Concept of “Explaining away”
5. “D-separation” property in directed graphs
6. Examples
 1. Independent identically distributed samples in
 1. Univariate parameter estimation
 2. Bayesian polynomial regression
 2. Naïve Bayes classifier
7. Directed graph as filter
8. Markov blanket of a node

1. Conditional Independence

- Consider three variables a, b, c
- Conditional distribution of a given b and c , is $p(a|b, c)$
- If $p(a|b, c)$ does not depend on value of b
 - we can write $p(a|b, c) = p(a|c)$
- We say that
 - a is conditionally independent of b given c

Factorizing into Marginals

- If conditional distribution of a given b and c does not depend on b
 - we can write $p(a|b,c) = p(a|c)$
- Can be expressed in slightly different way
 - written in terms of joint distribution of a and b conditioned on c
$$p(a,b|c) = p(a|b,c)p(b|c) \quad \text{using product rule}$$
$$= p(a|c)p(b|c) \quad \text{using earlier statement}$$
- That is joint distribution factorizes into product of marginals
 - Says that variables a and b are statistically independent given c
- Shorthand notation for conditional independence

$$a \perp b | c$$

An example with Three Binary Variables

$a \in \{A, \sim A\}$, where A is Red

$b \in \{B, \sim B\}$ where B is Blue

$c \in \{C, \sim C\}$ where C is Green

There are 90 different probabilities

In this problem!

Marginal Probabilities (6)

$p(a)$: $P(A)=16/49$ $P(\sim A)=33/49$

$p(b)$: $P(B)=18/49$ $P(\sim B)=31/49$

$p(c)$: $P(C)=12/49$ $P(\sim C)=37/49$

Joint (12 with 2 variables)

$p(a,b)$: $P(A,B)=6/49$ $P(A,\sim B)=10/49$

$P(\sim A,B)=6/49$ $P(\sim A,\sim B)=21/49$

$p(b,c)$: $P(B,C)=6/49$ $P(B,\sim C)=12/49$

$P(\sim B,C)=6/49$ $P(\sim B,\sim C)=25/49$

$p(a,c)$: $P(A,C)=4/49$ $P(A,\sim C)=12/49$

$P(\sim A,C)=8/49$ $P(\sim A,\sim C)=25/49$

Joint (8 with 3 variables)

$p(a,b,c)$: $P(A,B,C)=2/49$

$P(A,B,\sim C)=4/49$

$P(A,\sim B,C)=2/49$ $P(A,\sim B,\sim C)=8/49$

$P(\sim A,B,C)=4/49$ $P(\sim A,B,\sim C)=8/49$

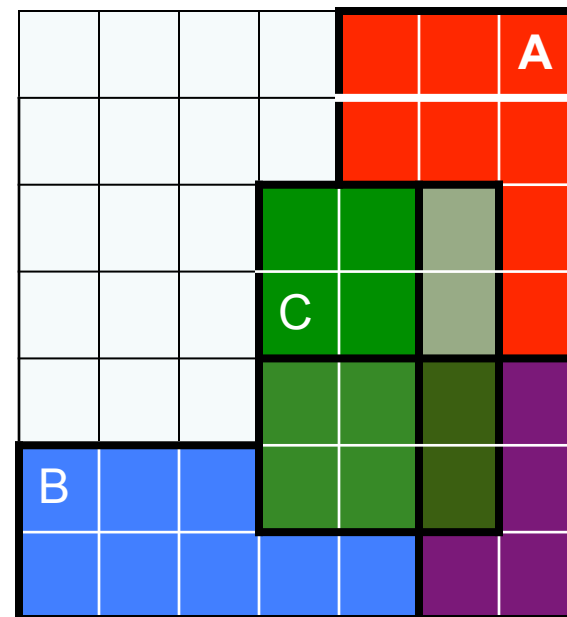
$P(\sim A,\sim B,C)=4/49$ $P(\sim A,\sim B,\sim C)=17/49$

$$p(a,b,c)=p(a)p(b,c/a)=p(a)p(b/a,c)p(c/a)$$

Are there any conditional independences?

Probabilities are assigned using a Venn diagram: allows us to evaluate every probability by inspection

Shaded areas with respect to total area 7 x 7 square



Three Binary Variables Example

Single variables conditioned on single variables (16)
(obtained from earlier values)

$$p(a/b): P(A/B)=P(A,B)/P(B)=1/3 \quad P(\sim A/B)=2/3$$

$$P(A/\sim B)=P(A,\sim B)/P(\sim B)=10/31 \quad P(\sim A/\sim B)=21/31$$

$$P(A/C)=P(A,C)/P(C)=1/3 \quad P(\sim A/C)=2/3$$

$$P(A/\sim C)=P(A,\sim C)/P(\sim C)=12/37 \quad P(\sim A/\sim C)=25/37$$

$$P(B/C)=P(B,C)/P(C)=1/2 \quad P(\sim B/C)=1/2$$

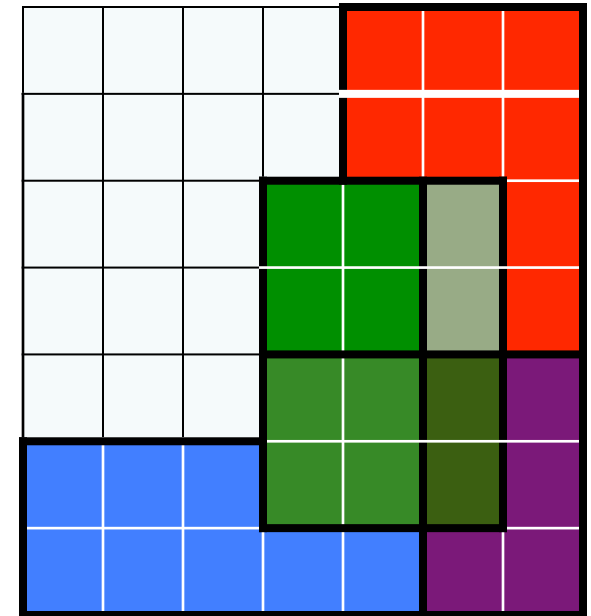
$$P(B/\sim C)=P(B,\sim C)/P(\sim C)=12/37 \quad P(\sim B/\sim C)=25/37$$

$$P(B/A) \quad P(\sim B/A)$$

$$P(B/\sim A) \quad P(\sim B/\sim A)$$

$$P(C/A) \quad P(\sim C/A)$$

$$P(C/\sim A) \quad P(\sim C/\sim A)$$



Three Binary Variables: Conditional Independences

Two variables conditioned on single variable(24)

$$p(a,b/c): P(A,B/C)=P(A,B,C)/P(C)=1/6$$

$$P(A,B/\sim C)=P(A,B,\sim C)/P(\sim C)=4/37$$

$$P(A,\sim B/C)=P(A,\sim B,C)/P(C)=1/6$$

$$P(A,\sim B/\sim C)=P(A,\sim B,\sim C)/P(\sim C)=8/37$$

$$P(\sim A,B/C)=P(\sim A,B,C)/P(C)=1/3$$

$$P(\sim A,B/\sim C)=P(\sim A,B,\sim C)/P(\sim C)=8/37$$

$$P(\sim A,\sim B/C)=P(\sim A,\sim B,C)/P(C)=1/3$$

$$P(\sim A,\sim B/\sim C)=P(\sim A,\sim B,\sim C)/P(\sim C)=17/37$$

$p(a,c/b)$ eight values

$p(b,c/a)$ eight values

Similarly 24 values for one variables conditioned on two: $p(a/b,c)$ etc

There are no Independence Relationships

$$P(A,B) \neq P(A)P(B) \quad P(A,\sim B) \neq P(A)P(\sim B) \quad P(\sim A,B) \neq P(\sim A)P(B)$$

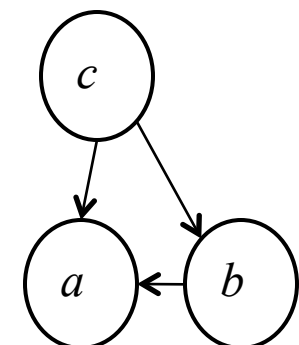
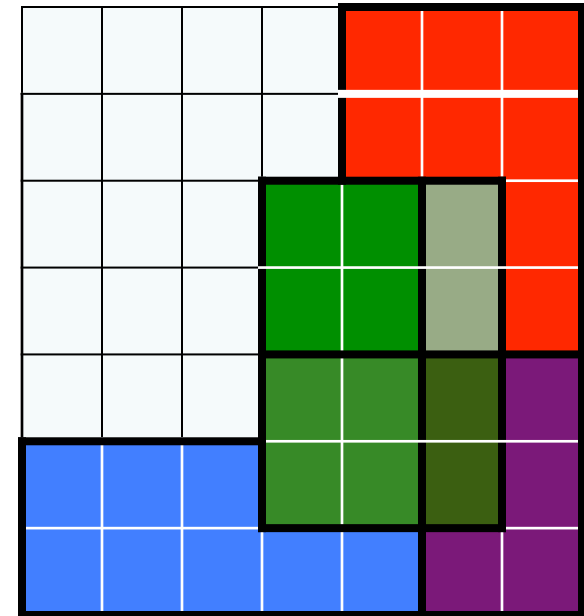
$$P(A,B/C)=1/6=P(A/C)P(B/C) \quad \text{but} \quad P(A,B/\sim C)=4/37 \neq P(A/\sim C)P(B/\sim C)$$

$$p(a,b,c)=p(c)p(a,b/c)=p(c)p(a/b,c)p(b/c)$$

$$\text{If } p(a,b/c)=p(a/c)p(b/c) \text{ then } p(a/b,c)=p(a,b/c)/p(b/c)=p(a/c)$$

Then the arrow from b to a would be eliminated

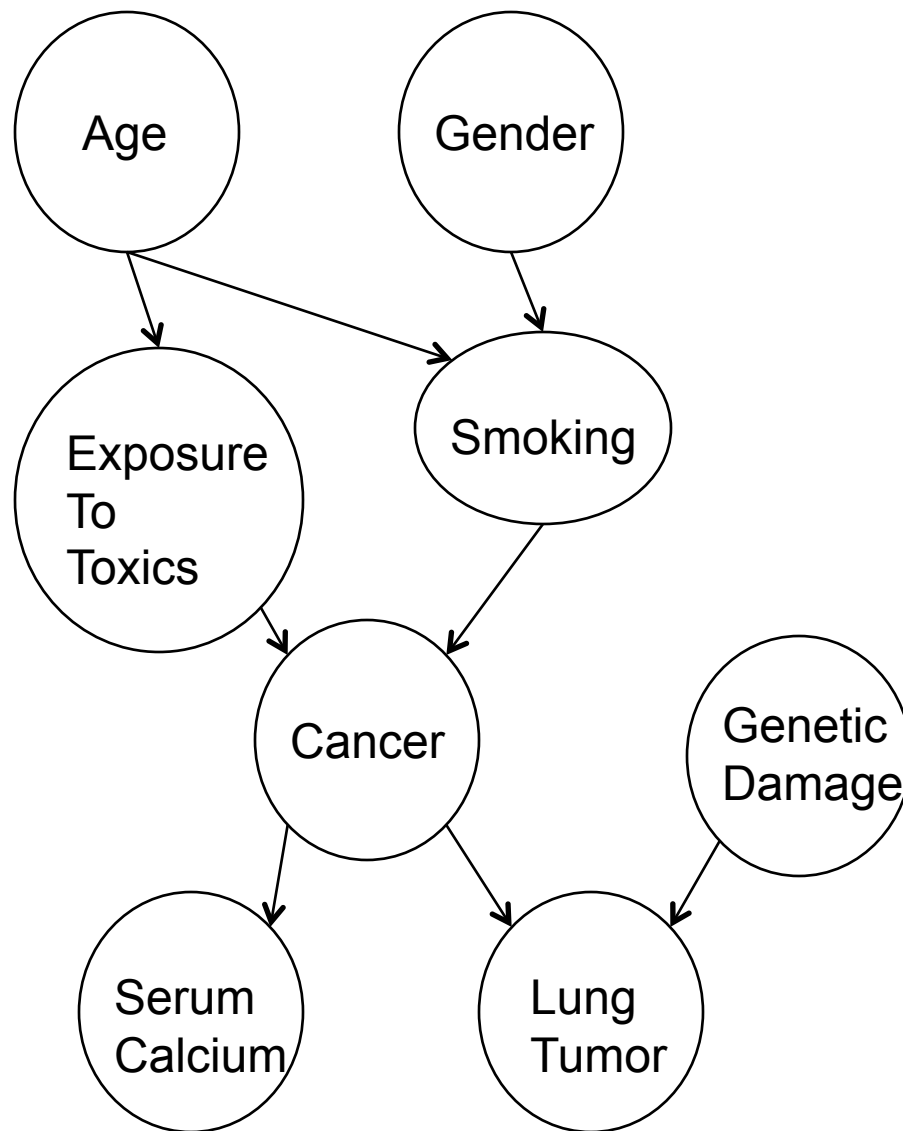
If we knew the graph structure a priori then we could simplify some probability calculations



2. Importance of Conditional Independence

- Important concept for probability distributions
- Role in PR and ML
 - Simplifying structure of model
 - Computations needed to perform inference and learning
- Role of graphical models
 - Testing for conditional independence from an expression of joint distribution is time consuming
 - Can be read directly from the graphical model
 - Using framework of *d-separation* (“directed”)

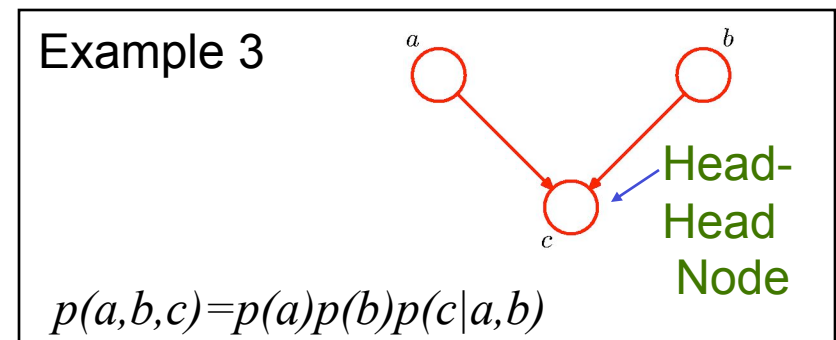
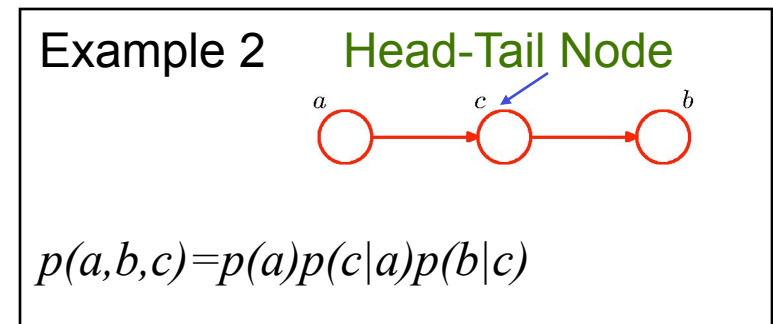
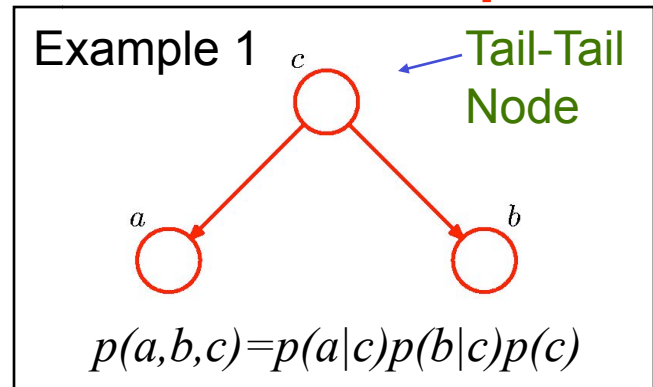
A Causal Bayesian Network



Cancer is independent of Age and Gender given Exposure to toxics and Smoking

3. Conditional Independence from Graphs

- Three Example Graphs
- Each has just three nodes
- Together they illustrate concept of *d* (directed)-separation



First example *without* conditioning

- Joint distribution is

$$p(a,b,c) = p(a|c)p(b|c)p(c)$$

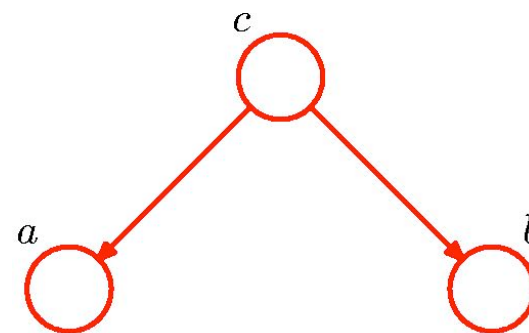
- If none of the variables are observed, we can investigate whether a and b are independent
 - marginalizing both sides wrt c gives

$$p(a,b) = \sum_c p(a|c)p(b|c)p(c)$$

- In general this does not factorize to $p(a)p(b)$, and so

$$a \not\perp b | \phi$$

- Meaning conditional independence property does not hold given the empty set ϕ



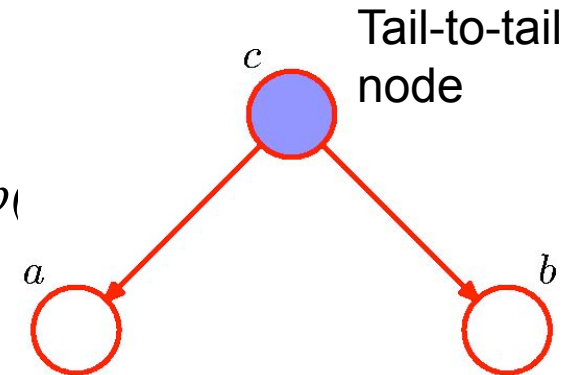
Note: May hold for a particular distribution due to specific numeral values. but does not follow in general from graph structure ¹¹

First example *with* conditioning

- Condition variable c

- 1. By product rule $p(a,b,c)=p(a,b|c)p(c)$
- 2. According to graph $p(a,b,c)=p(a|c)p(b|c)p(c)$
- 3. Combining the two

$$p(a,b|c)=p(a|c)p(b|c)$$



- So we obtain the conditional independence property

$$a \perp b | c$$

- Graphical Interpretation

- Consider path from node a to node b
 - Node c is *tail-to-tail* since it is connected to the tails of the two arrows
 - Causes these nodes to be dependent

- When conditioned on node c ,

- blocks path from a to b
- causes a and b to become conditionally independent

Second Example *without* conditioning

- Joint distribution is

$$p(a,b,c) = p(a)p(c|a)p(b|c)$$

- If no variable is observed

- Marginalizing both sides wrt c

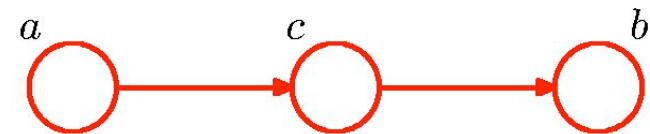
$$p(a,b) = p(a) \sum_c p(c|a)p(b|c)$$

- It does not generalize to $p(a)p(b)$ and so

$$a \not\perp b | \phi$$

- Graphical Interpretation

- Node c is *head-to-tail* wrt path from node a to node b
- Such a path renders them dependent



head-to-tail
node

Second Example *with* conditioning

- Condition on node c
 - 1. By product rule $p(a,b,c)=p(a,b|c)p(c)$
 - 2. According to graph $p(a,b,c)=p(a)p(c|a)p(b|c)$
 - 3. Combining the two

$$p(a,b|c) = p(a)p(c|a)p(b|c)/p(c)$$

- 4. By Bayes rule

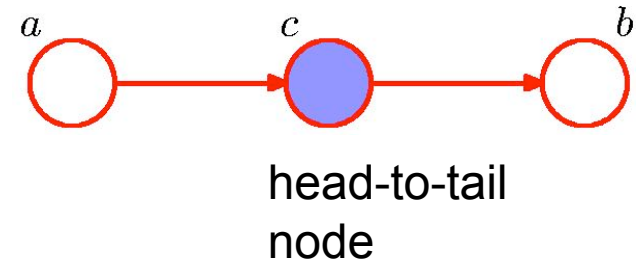
$$p(a,b|c) = p(a|c) p(b|c)$$

- Therefore

$$a \perp b | c$$

- Graphical Interpretation

- If we observe c
 - Observation blocks path from from a to b
 - Thus we get conditional independence



Third Example *without* conditioning

- Joint distribution is

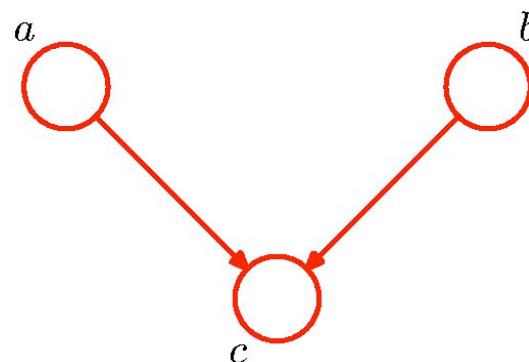
$$p(a,b,c) = p(a)p(b)p(c|a,b)$$

- Marginalizing both sides gives

$$p(a,b) = p(a)p(b)$$

- Thus a and b are independent in contrast to previous two examples, or

$$a \perp b \mid \phi$$

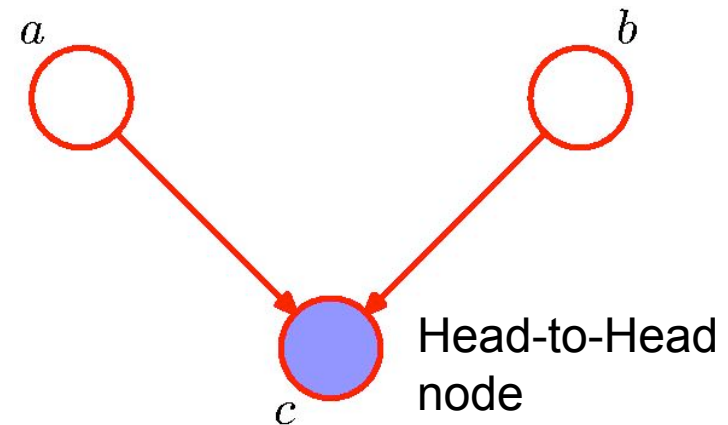


Third Example *with* conditioning

- By product rule
- $p(a,b|c)=p(a,b,c)/p(c)$
- Joint distribution is $p(a,b,c)=p(a)p(b)p(c|a,b)$
- Does not factorize into product $p(a)p(b)$

$$a \not\perp b | c$$

- Therefore, example has opposite behavior from previous two cases
- Graphical Interpretation
 - Node c is head-to-head wrt path from a to b
 - It connects to the heads of the two arrows
 - When c is unobserved it blocks the path and the variables are independent
 - Conditioning on c unblocks the path and renders a and b dependent

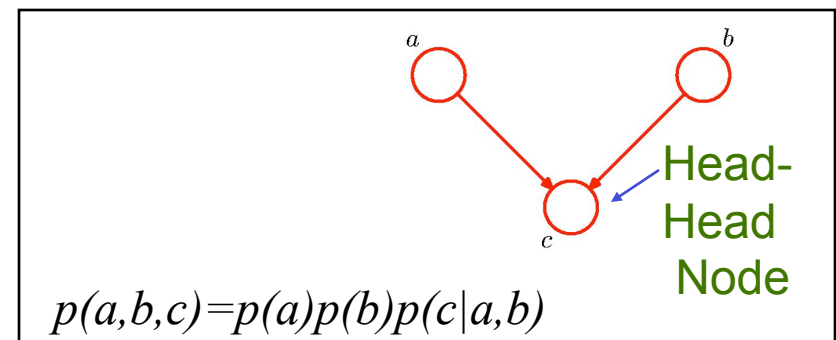
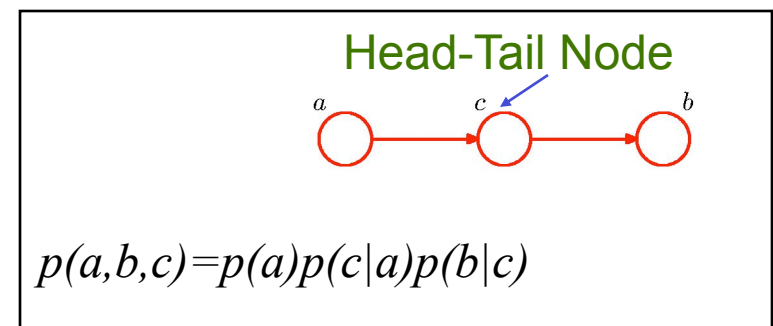
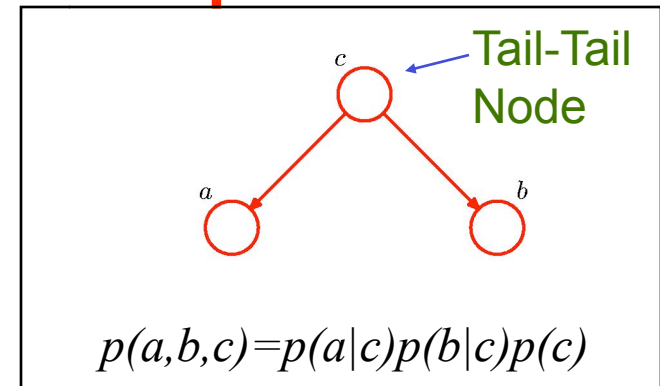


General Rule

- Node y is a descendent of Node x if there is a path from x to y in which each step follows direction of arrows.
- Head-to-head path becomes unblocked if either node *or any of its descendants* is observed¹⁶

Summary of three examples

- Conditioning on **tail-to-tail** node renders *conditional independence*
- Conditioning on **head-to-tail** node renders *conditional independence*
- Conditioning on **head-to-head** node renders *conditional dependence*



4. Example to understand Case 3

- Three binary variables

- Battery: B

- Charged ($B=1$) or Dead ($B=0$)

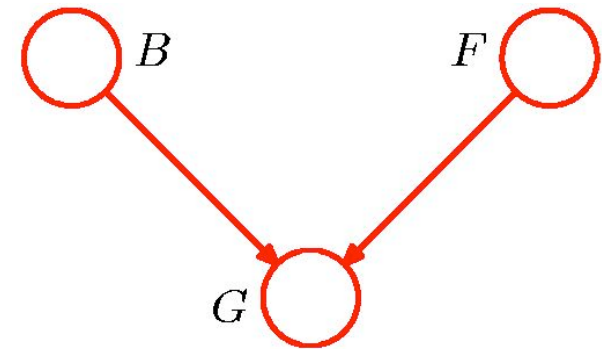
- Fuel Tank: F

- Full ($F=1$) or Empty ($F=0$)

- Guage, Electric Fuel: G

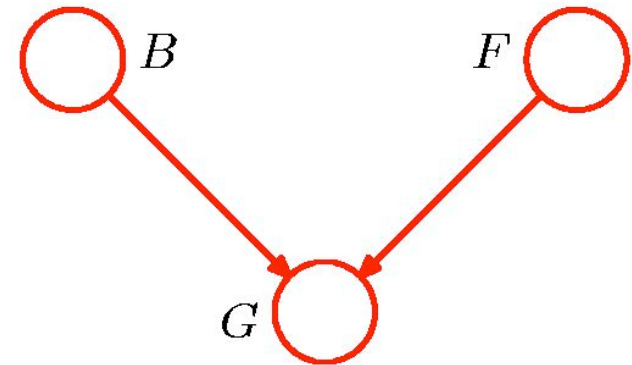
- Indicates Full($G=1$) or Empty($G=0$)

- B and F are independent with prior probabilities



Fuel System: Given Probability Values

- We are given prior probabilities and one set of conditional probabilities



Battery
Prior Probabilities

 $p(B)$

B	$p(B)$
1	0.9
0	0.1

Fuel
Prior Probabilities

 $p(F)$

F	$p(F)$
1	0.9
0	0.1

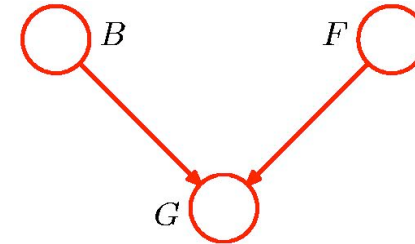
Conditional probabilities of Gauge

 $p(G|B,F)$

B	F	$p(G=1)$
1	1	0.8
1	0	0.2
0	1	0.2
0	0	0.1

Fuel System Joint Probability

- Given prior probabilities and one set of conditional probabilities



Probabilistic structure is completely specified since

$$\begin{aligned}
 p(G, B, F) &= p(B, F | G) p(G) \\
 &= p(G | B, F) p(B, F) \\
 &= p(G | B, F) p(B) p(F)
 \end{aligned}$$

e.g., $p(G) = \sum_{B, F} p(G | B, F) p(B) p(F)$

$$p(G | F) = \sum_B p(G, B | F)$$

by sum rule

= $\sum_B p(G | B, F) p(B)$ prod rule

B	$p(B)$
1	0.9
0	0.1

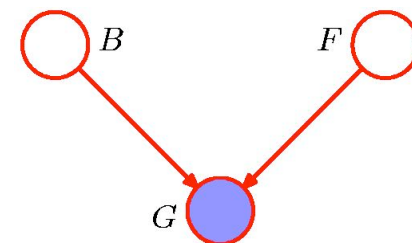
Conditional probs

B	F	$p(G=1)$
1	1	0.8
1	0	0.2
0	1	0.2
0	0	0.1

F	$p(F)$
1	0.9
0	0.1

$$p(F | G, B) = \frac{p(G, B | F) p(F)}{p(G, B)} = \frac{p(G | B, F) p(B | F) p(F)}{p(G | B) p(B)} = \frac{p(G | B, F) p(F)}{\sum_F p(G | B, F) p(F)}$$

Fuel System Example

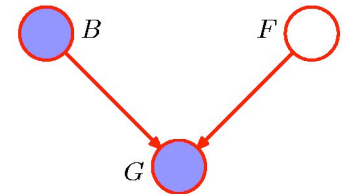


- Suppose guage reads empty ($G=0$)
- We can use Bayes theorem to evaluate fuel tank being empty ($F=0$)

$$p(F = 0 | G = 0) = \frac{p(G = 0 | F = 0)p(F = 0)}{p(G = 0)}$$

- **Where** $p(G = 0) = \sum_{B=0,1} \sum_{F=0,1} p(G = 0 | B, F)p(B)p(F) = 0.315$
 $p(G = 0 | F = 0) = \sum_{B=0,1} p(G = 0 | B, F = 0)p(B) = 0.81$
- **Therefore** $p(F=0|G=0)=0.257 > p(F=0)=0.1$ ²¹

Clamping additional node

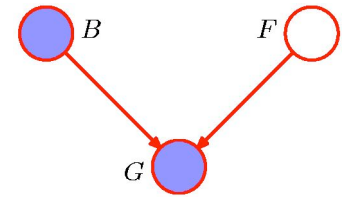


- Observing both fuel guage and battery
- Suppose Guage reads empty ($G=0$) and Battery is dead ($B=0$)
- Probability that Fuel tank is empty

$$p(F = 0 | G = 0, B = 0) = \frac{p(G = 0 | B = 0, F = 0)p(F = 0)}{\sum_{F=0,1} p(G = 0 | B = 0, F)p(F)} = 0.111$$

- Probability has decreased from 0.257 to 0.111

Explaining Away



- Probability has decreased from 0.257 to 0.111
- *Battery is dead explains away the observation that fuel guage reads empty*
- *State of the fuel tank and that of battery have become dependent on each other as a result of observing the fuel guage*
- *This would also be the case if we observed state of some descendant of G*
- *Note: 0.111 is greater than $P(F=0)=0.1$ because*
 - *observation that fuel guage reads zero provides some evidence in favor of empty tank*

5. D-separation Motivation

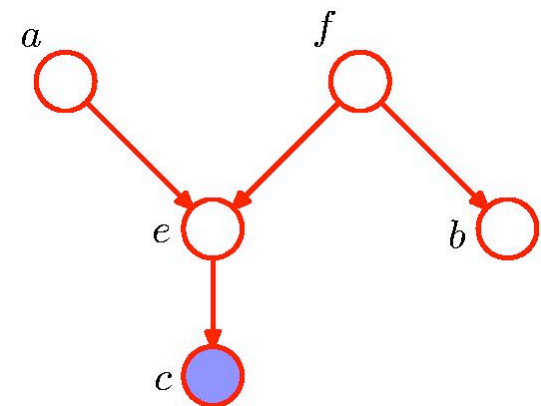
- Consider general directed graph where
 - A, B, C are arbitrary, non-intersecting nodes
 - Whose union could be smaller than complete set of nodes
- Wish to ascertain
 - Whether an independence statement $A \perp B \mid C$ is implied by a directed acyclic graph (no paths from node to itself)

D-separation Definition

- Consider all possible paths from any node in A to any node in B
- If all possible paths are *blocked* then A is d-separated from B by C , i.e.,
 - Joint distribution will satisfy $A \perp B | C$
- A path is blocked if it includes a node such that
 - Arrows in path meet *head-to-tail* or *tail-to-tail* at the node and node is in set C
 - Arrows meet *head-to-head* at the node, and
 - Neither the node or its descendants is in set C

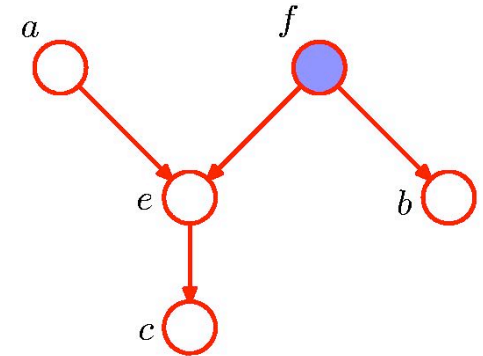
D-separation Example 1

- Path from a to b is
 - Not blocked by f because
 - It is a *tail-to-tail* node
 - It is not observed
 - Not blocked by e because
 - Although it is a head-to-head node, it has a descendent c in the conditioning set
- Thus $a \perp b | c$ does *not* follow from this graph



D-separation Example 2

- Path from a to b is
 - Blocked by f because
 - It is a *tail-to-tail* node that is observed
- Thus $a \perp b \mid f$ is satisfied by any distribution that factorizes according to this graph
- It is also blocked by e
 - It is a *head-to-head* node
 - Neither it nor its descendants is in conditioning set

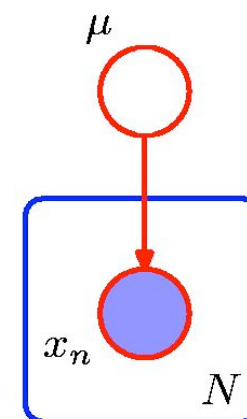
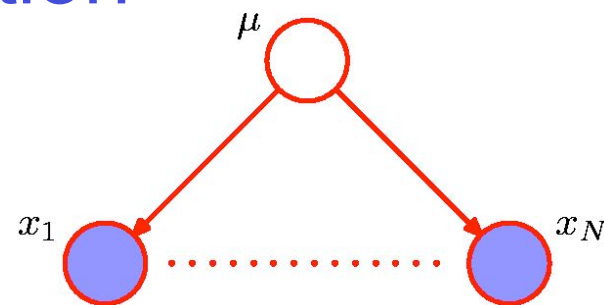


6. Three practical examples of conditional independence and d-separation

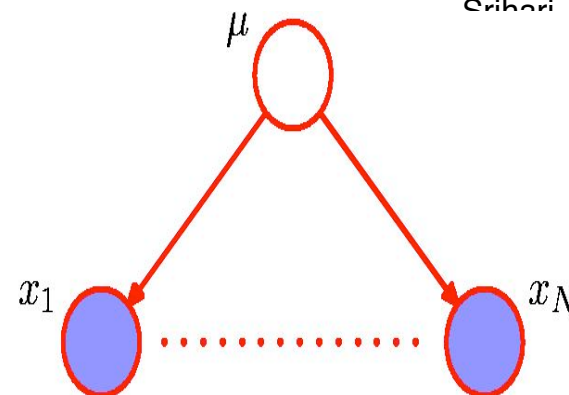
1. Independent, identically distributed samples
 1. in estimating mean of univariate Gaussian
 2. In Bayesian polynomial regression
2. Naïve Bayes classification

I.I.D. data to illustrate Cond. Indep.

- Task of finding posterior distribution for mean of univariate Gaussian distribution
- Joint distribution represented by graph
 - Prior $p(\mu)$
 - Conditionals $p(x_n|\mu)$ for $n=1,\dots,N$
- We observe $D=\{x_1,\dots,x_N\}$
- Goal is to infer μ



Inferring Gaussian Mean



- Suppose we condition on μ
- Consider joint distribution of observations
- Using d-separation,
 - There is a unique path from any x_i to any other x_j
 - Path is tail-to-tail wrt observed node μ
 - Every such path is blocked
 - So observations $D = \{x_1, \dots, x_N\}$ are independent given μ so that

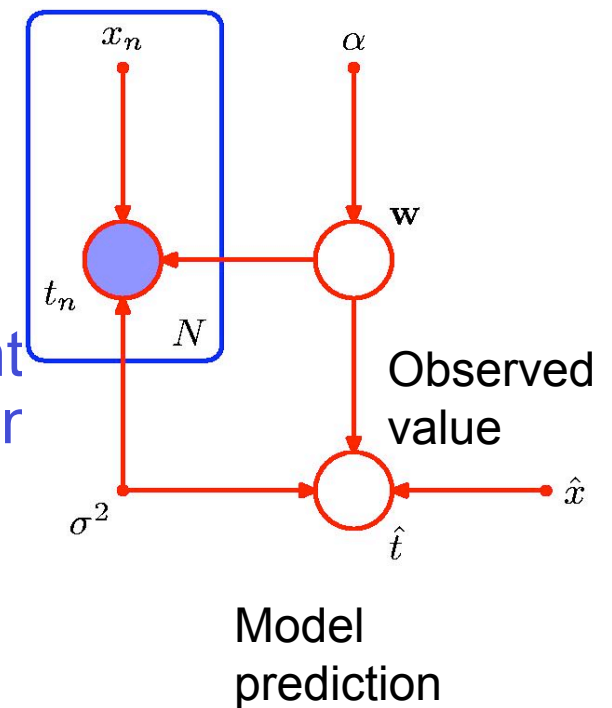
$$p(D | \mu) = \prod_{n=1}^N p(x_n | \mu)$$

- However if we integrate over μ , the observations are no longer independent

$$p(D) = \int_0^{\infty} p(D | \mu) p(\mu) d\mu \neq \prod_{n=1}^N p(x_n)$$

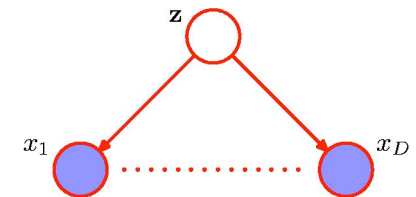
Bayesian Polynomial Regression

- Another example of i.i.d. data
- Stochastic nodes correspond to t_n , w and t^\wedge
- Node for w is tail-to-tail wrt path from t^\wedge to any one of nodes t_n
- So we have conditional independence property
- Thus conditioned on polynomial coefficient w predictive distribution of t^\wedge is independent of training data $\{t_n\}$
- Therefore we can first use training data to determine posterior distribution over coefficients w
 - and then discard training data and use posterior distribution for w to make predictions of t^\wedge for new input x^\wedge



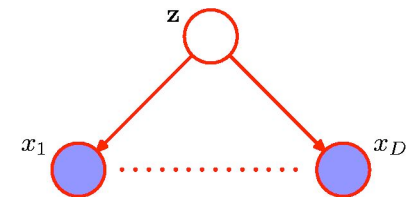
Naïve Bayes Model

- Conditional independence assumption is used to simplify model structure
- Observed variable $\mathbf{x} = (x_1, \dots, x_D)^T$
- Task is to assign \mathbf{x} to one of K classes
- Use 1 of K coding scheme
 - Class represented by K -dimensional binary vector z
- Generative model
 - multinomial prior $p(z|\mu)$ over class labels where $\mu = \{\mu_1, \dots, \mu_K\}$ are class priors
 - Class conditional distribution $p(\mathbf{x}|z)$ for observed vector \mathbf{x}
- Key assumption
 - Conditioned on class z
 - distribution of input variables x_1, \dots, x_D are independent
 - Leads to graphical model shown



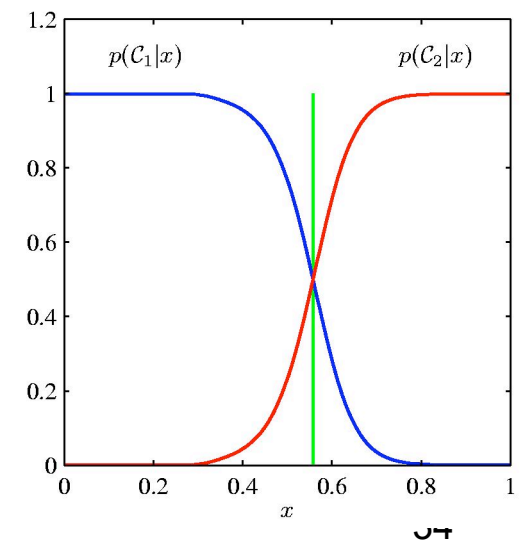
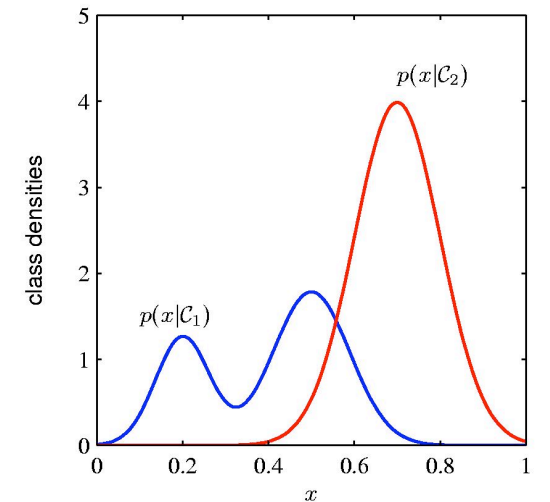
Naïve Bayes Graphical Model

- Observation of z blocks path between x_i and x_j
 - Because such paths are tail-to-tail at node z
- So x_i and x_j are conditionally independent given z
- If we marginalize out z (so that z is unobserved)
 - Tail-to-tail path is no longer blocked
 - Implies marginal density $p(\mathbf{x})$ will not factorize w.r.t. components of \mathbf{x}



When is Naïve Bayes good?

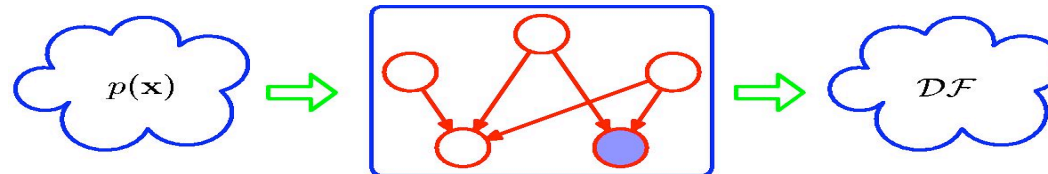
- Helpful when:
 - Dimensionality of the input space is high making density estimation in D -dimensional space challenging
 - When input vector contains both discrete and continuous variables
 - Some can be Bernoulli, others Gaussian, etc
- Good classification performance in practice since
 - Decision boundaries can be insensitive to details of class conditional densities



7. Directed graph as filter

- A directed graph represents
 - a specific decomposition of joint distribution into a product of conditional probabilities
 - Also expresses conditional independence statements obtained through d-separation criteria
- Can think of directed graph as a filter
 - Allows a distribution to pass through iff it can be expressed as the factorization implied

Directed Factorization



- Of all possible distributions $p(\mathbf{x})$ the subset passed through are denoted \mathcal{DF}
- Alternatively graph can be viewed as filter that lists all conditional independence properties obtained by applying d-separation criterion
 - Then allow a distribution if it satisfies all those properties
- Passes through a whole family of distributions
 - Discrete or continuous
 - Will include distributions that have additional independence properties beyond those described by graph

8. Markov Blanket

- Joint distribution $p(\mathbf{x}_1, \dots, \mathbf{x}_D)$ represented by directed graph with D nodes

- Distribution of \mathbf{x}_i conditioned on rest is

$$p(\mathbf{x}_i | \mathbf{x}_{\{j \neq i\}}) = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_D)}{\int p(\mathbf{x}_1, \dots, \mathbf{x}_D) d\mathbf{x}_i} \quad \text{from product rule}$$

$$= \frac{\prod_k p(\mathbf{x}_k | \text{pa}_k)}{\int \prod_k p(\mathbf{x}_k | \text{pa}_k) d\mathbf{x}_i} \quad \text{from factorization property}$$

- Factors not depending on \mathbf{x}_i can be taken outside integral over \mathbf{x}_i and cancel between numerator and denominator
- Remaining will depend only on parents, children and *co-parents*

Markov Blanket of a Node

- Set of parents, children and co-parents of a node
- The conditional distribution of x_i conditioned on all the remaining variables in the graph is dependent on only the variables in the Markov blanket
- Also called *Markov boundary*

