

Joint Distribution

Motivation

- The joint distribution of X_j and $X_{j'}$ provides the probability of $X_j = x_j$ and $X_{j'} = x_{j'}$ occur at the same time.

$$\mathbb{P}\{X_j = x_j, X_{j'} = x_{j'}\}$$

classification

$\hat{z}_t = \underset{z}{\operatorname{argmax}} \Pr\{z | z_{t-1}\} \cdot \frac{c_{z,t}}{c_{z,t-1}}$

training
 $\hat{p}_t\{z_t | z_{t-1}\}$

- The marginal distribution of X_j can be found by summing over all possible values of $X_{j'}$.

$$\mathbb{P}\{X_j = x_j\} = \sum_{x \in X_{j'}} \mathbb{P}\{X_j = x_j, X_{j'} = x\}$$

$z_0 \rightarrow z_1 \rightarrow z_2$
 $\rightarrow z_c.$

$$\Pr\{X_j = x_j\} = \Pr\{X_j = x_j, X_{j'} = 0\} + \Pr\{X_j = x_j, X_{j'} = 1\}$$

Conditional Distribution

Motivation

- Suppose the joint distribution is given.

$$\mathbb{P} \{ X_j = x_j, X_{j'} = x_{j'} \}$$

- The conditional distribution of X_j given $X_{j'} = x_{j'}$ is ratio between the joint distribution and the marginal distribution.

$$\mathbb{P} \{ X_j = x_j | X_{j'} = x_{j'} \} = \frac{\mathbb{P} \{ X_j = x_j, X_{j'} = x_{j'} \}}{\mathbb{P} \{ X_{j'} = x_{j'} \}}$$

Notation

Motivation

- The notations for joint, marginal, and conditional distributions will be shortened as the following.

$$\mathbb{P}\{x_j, x_{j'}\}, \mathbb{P}\{x_j\}, \mathbb{P}\{x_j|x_{j'}\}$$

- When the context is not clear, for example when $x_j = a, x_{j'} = b$ with specific constants a, b , subscripts will be used under the probability sign.

$$\mathbb{P}_{x_j, x_{j'}}\{a, b\}, \mathbb{P}_{x_j}\{a\}, \mathbb{P}_{x_j|x_{j'}}\{a|b\}$$

$$P_r(H|A) \cdot P_r(A) + P_r(H|B) \cdot P_r(B)$$

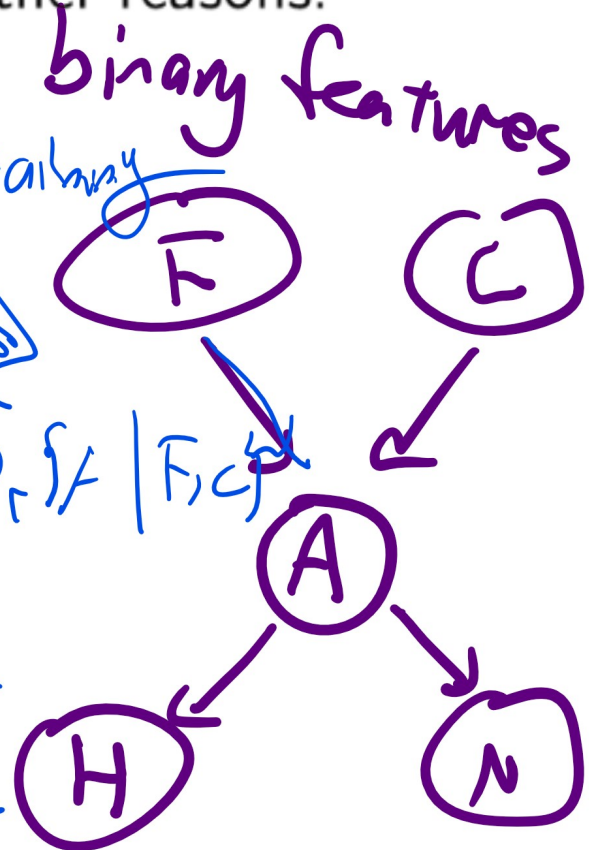
Bayesian Network Diagram

Definition

- Story: You are travelling far from home. There may be a Fire problem or a Cat problem at home. Either problem might trigger an Alarm. Then your neighbors Nick or Happy or both might call you because of the alarm or for other reasons.

$P_r(H | \text{no } A)$
 $P_r(F | H, N)?$
 training
 see
 $P_r(C)$

	x_1	x_2	x_3	x_4	x_5
	F	C	A	H	N
day 1	0	0	0	1	0
day 2	0	1	0	0	0
day 3	0	0	0	1	1
	1	0	0	0	1
	0	0	1	1	0
	0	0	1	0	1
	0	0	1	1	1
	0	0	1	1	1



inference
 simulation

generate new data

simulation

Bayesian Network


Definition

- A Bayesian network is a directed acyclic graph (DAG) and a set of conditional probability distributions.
- Each vertex represents a feature X_j .
- Each edge from X_j to $X_{j'}$ represents that X_j directly influences $X_{j'}$.
- No edge between X_j and $X_{j'}$ implies independence or conditional independence between the two features.

Conditional Independence

Definition

- Recall two events A, B are independent if:

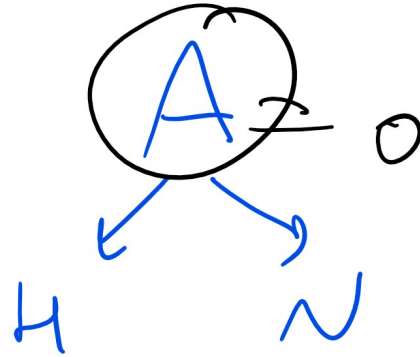
$$\mathbb{P}\{A, B\} = \mathbb{P}\{A\} \mathbb{P}\{B\} \quad \text{or} \quad \mathbb{P}\{A|B\} = \mathbb{P}\{A\}$$


- In general, two events A, B are conditionally independent, conditional on event C if:

$$\mathbb{P}\{A, B|C\} = \mathbb{P}\{A|C\} \mathbb{P}\{B|C\} \quad \text{or} \quad \mathbb{P}\{A|B, C\} = \mathbb{P}\{A|C\}$$


Common Cause

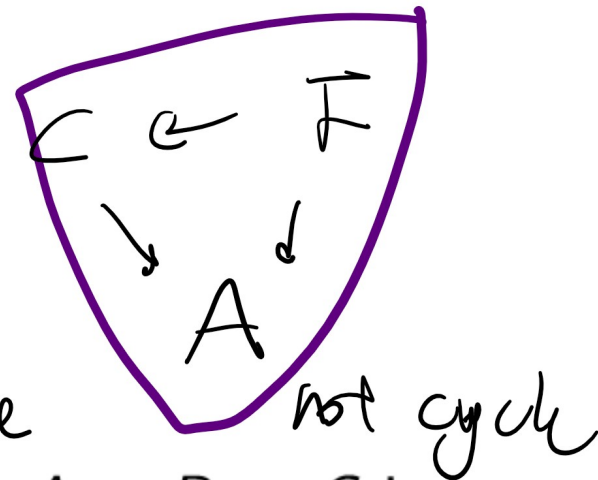
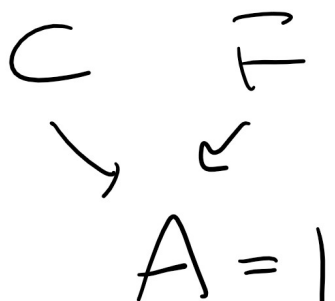
Definition



- For three events A, B, C , the configuration $A \leftarrow B \rightarrow C$ is called common cause.
- In this configuration, A is not independent of C , but A is conditionally independent of C given information about B .
- Once B is observed, A and C are independent.

Common Effect

Definition



- For three events A, B, C , the configuration $A \rightarrow B \leftarrow C$ is called common effect.
- In this configuration, A is independent of C but A is not *correct* conditionally independent of C given information about B .
- Once B is observed, A and C are not independent.

$P(A|B)$
 \downarrow
 A is not conditionally indep of C //

$$P(A|C, B) = \frac{P(A, C, B)}{P(B, C)} = \frac{P(C)P(A, B)}{P(B, C)} \neq \frac{P(C)P(A, B)}{P(B)P(C)}$$



Storing Distribution

Definition

$$P(\{A, B, C\}) = P(A | C) \text{ if } A, B \text{ are independent}$$

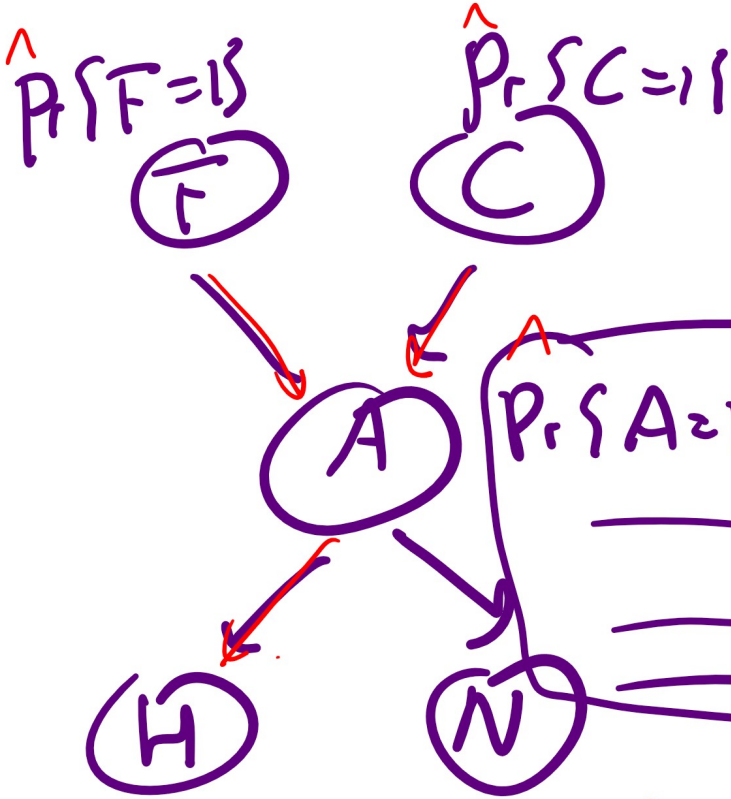
- If there are m binary variables with k edges, there are 2^m joint probabilities to store.
- There are significantly less conditional probabilities to store. For example, if each node has at most 2 parents, then there are less than $4m$ conditional probabilities to store.
- Given the conditional probabilities, the joint probabilities can be recovered.

$$\begin{aligned} P(\{F, C, A, H, N\}) &= P(F) \cdot P(C, A, H, N | F) \\ &= P(F) \cdot P(C | F) \cdot P(A, H, N | C, F) \\ &= P(F) \cdot P(C | \cancel{F}) \cdot P(A | C, F) \cdot P(H | \cancel{A, F}) \end{aligned}$$

$P_r(SN) / \{A, H, N\}$

Conditional Probability Table Diagram

Definition



$P_r(A=1 F=1, C=1)$	
	0 1
	1 0
	0 0

	F	C	A	H	N
P_r	0	0	0	0	0
	0	0	0	0	1
	0	0	0	1	0
	0	0	0	1	1

$\hat{P}_r(H=1 | A=0)$

$\hat{P}_r(N=1 | A=0)$

$\hat{P}_r(H=1 | A=1)$

$\hat{P}_r(N=1 | A=1)$

$2^5 - 1 = 31$

10 complexity

$P_r(F=1, C=1, A=1, H=1, N=1)$

$= \prod_{j=1}^m P_r(x_j | \text{Parent } x_i)$

Training Bayes Net

Definition

- Training a Bayesian network given the DAG is estimating the conditional probabilities. Let $P(X_j)$ denote the parents of the vertex X_j , and $p(X_j)$ be realizations (possible values) of $P(X_j)$.

$$\mathbb{P}\{x_j | p(X_j)\}, p(X_j) \in P(X_j)$$

- It can be done by maximum likelihood estimation given a training set.

$$\hat{\mathbb{P}}\{x_j | p(X_j)\} = \frac{c_{x_j, p(X_j)}}{c_{p(X_j)}}$$

Bayes Net Training Example, Training, Part I

Definition

- Given a network and the training data.

$F \rightarrow A, C \rightarrow A, A \rightarrow H, A \rightarrow N.$

F	C	A	H	N
0	0	0	1	0
0	1	0	0	0
0	0	0	1	1
1	0	0	0	1
0	0	1	1	0
0	0	1	0	1
0	0	1	1	1
0	0	1	1	1

Bayes Net Training Example, Training, Part II

Definition

- Compute $\hat{\mathbb{P}}\{F = 1\} = \frac{C_{F=1}}{n} = \frac{1}{8}$

F	C	A	H	N
0	0	0	1	0
0	1	0	0	0
0	0	0	1	1
1	0	0	0	1
0	0	1	1	0
0	0	1	0	1
0	0	1	1	1
0	0	1	1	1

Bayes Net Training Example, Training, Part III

Definition

- Compute $\hat{\mathbb{P}} \{ \underline{H = 1} | \underline{A = 0} \}$

$$\frac{C_{H=1, A=0}}{C_{A=0}} = \frac{2}{4} = \frac{1}{2}$$

F	C	A	H	N
0	0	0	1	0
0	1	0	0	0
0	0	0	1	1
1	0	0	0	1
0	0	1	1	0
0	0	1	0	1
0	0	1	1	1
0	0	1	1	1

Bayes Net Training Example, Training, Part V

Definition

- Compute $\hat{\mathbb{P}} \{A = 1 | F = 0, C = 1\}$.

F	C	A	H	N
0	0	0	1	0
0	1	0	0	0
0	0	0	1	1
1	0	0	0	1
0	0	1	1	0
0	0	1	0	1
0	0	1	1	1
0	0	1	1	1

$$\frac{C_{A=1, F=0, C=1}}{C_{F=0, C=1}} = \frac{0}{2} = 0$$

Bayes Net Training Example, Training, Part VI

Quiz (Graded)

- What is the conditional probability $\hat{\mathbb{P}}\{A = 1 | F = 0, C = 0\}$?
- A: 0 , B: $\frac{1}{3}$, C: $\frac{1}{2}$, D: $\frac{2}{3}$, E: 1

Q8

F	C	A	H	N
0	0	0	1	0
0	1	0	0	0
0	0	0	1	1
1	0	0	0	1
0	0	1	1	0
0	0	1	0	1
0	0	1	1	1
0	0	1	1	1

$$\frac{4}{6} = \frac{2}{3}$$

Laplace Smoothing

Definition

- Recall that the MLE estimation can incorporate Laplace smoothing.

$$\hat{\mathbb{P}}\{x_j | p(X_j)\} = \frac{c_{x_j, p(X_j)} + 1}{c_{p(X_j)} + |X_j|}$$

$$P_r\{A | F, C\} = \frac{c_{AFC} + 1}{c_{FC} + 4}$$

- Here, $|X_j|$ is the number of possible values (number of categories) of X_j .
- Laplace smoothing is considered regularization for Bayesian networks because it avoids overfitting the training data.

Bayes Net Inference Example, Part I

Definition

- Assume the network is trained on a larger set with the following CPT. Compute $\hat{\mathbb{P}}\{F = 1|H = 1, N = 1\}$?

$$\hat{\mathbb{P}}\{F = 1\} = 0.001, \hat{\mathbb{P}}\{C = 1\} = 0.001$$

$$\hat{\mathbb{P}}\{A = 1|F = 1, C = 1\} = 0.95, \hat{\mathbb{P}}\{A = 1|F = 1, C = 0\} = 0.94$$

$$\hat{\mathbb{P}}\{A = 1|F = 0, C = 1\} = 0.29, \hat{\mathbb{P}}\{A = 1|F = 0, C = 0\} = 0.00$$

$$\hat{\mathbb{P}}\{H = 1|A = 1\} = 0.9, \hat{\mathbb{P}}\{H = 1|A = 0\} = 0.05$$

$$\hat{\mathbb{P}}\{N = 1|A = 1\} = 0.7, \hat{\mathbb{P}}\{N = 1|A = 0\} = 0.01$$

Bayes Net Inference Example, Part II

Definition

$$\frac{Pr\{F, H, N\}}{Pr\{H, N\}} = \frac{Pr\{F, H, N\}}{Pr\{F, H, N\} + Pr\{\text{not } F, H, N\}}$$

$$Pr\{F, H, N\} = Pr\{F, H, N, A=0, C=0\}$$

- Compute $\hat{P}\{F = 1 | H = 1, N = 1\}$?

$$Pr\{F, H, N, \text{not } A, \text{not } C\} = \prod_{j=1}^n Pr\{x_j | \text{Parents}(x_j)\}$$

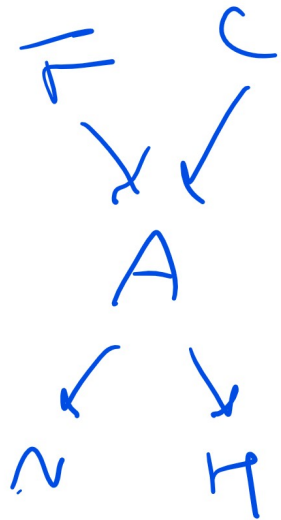
$$= Pr\{F\} \cdot Pr\{H | \text{not } A\} \cdot Pr\{N | \text{not } A\} \cdot Pr\{\text{not } A | F, C\} \cdot Pr\{\text{not } C\}$$

Bayes Net Inference Example, Part III

Definition

0,001

$$P(\underline{F}, \underline{H}, \underline{N}, \underline{\text{not } A}, \underline{C}) = 0.001 \cdot 0.05 \cdot 0.01 \cdot 0.05$$



$$\hat{P}\{F\} = 0.001, \hat{P}\{C\} = 0.001$$

$$\hat{P}\{A|F, C\} = 0.95, \hat{P}\{A|F, \neg C\} = 0.94$$

$$\hat{P}\{A|\neg F, C\} = 0.29, \hat{P}\{A|\neg F, \neg C\} = 0.00$$

$$\hat{P}\{H|A\} = 0.9, \hat{P}\{H|\neg A\} = 0.05$$

$$\hat{P}\{N|A\} = 0.7, \hat{P}\{N|\neg A\} = 0.01$$

Bayes Net Inference Example, Part IV

Definition

- Which of the following probabilities (multiple) are not required to compute $\hat{\mathbb{P}}\{C = 1|H = 1, N = 1\}$?
- A: $\hat{\mathbb{P}}\{A = 1|F = 1, C = 1\} = 0.95$
- B: $\hat{\mathbb{P}}\{A = 1|F = 1, C = 0\} = 0.94$
- C: $\hat{\mathbb{P}}\{A = 1|F = 0, C = 1\} = 0.29$
- D: $\hat{\mathbb{P}}\{A = 1|F = 0, C = 0\} = 0.00$
- E: none of the above.

Common Cause Example, Part I

Quiz (Graded)

- 2005 Fall Final Q20, 2006 Fall Final Q20
- Suppose A is the common cause of B and C . All variables are binary. What is $\mathbb{P}\{C = 1|B = 1\}$?

$$\mathbb{P}\{A = 1\} = 0.4, \mathbb{P}\{B = 1|A = 1\} = 0.9, \mathbb{P}\{B = 1|A = 0\} = 0.8$$

$$\mathbb{P}\{C = 1|A = 1\} = 0.5, \mathbb{P}\{C = 1|A = 0\} = 0.2$$

Common Cause Example, Part II

Quiz (Graded)

- What is $\mathbb{P}\{B = 1|C = 1\}$?

$$\mathbb{P}\{A = 1\} = 0.4, \mathbb{P}\{B = 1|A = 1\} = 0.9, \mathbb{P}\{B = 1|A = 0\} = 0.8$$

$$\mathbb{P}\{C = 1|A = 1\} = 0.5, \mathbb{P}\{C = 1|A = 0\} = 0.2$$

- A: $\frac{0.9 \cdot 0.4 \cdot 0.5 \cdot 0.4 + 0.8 \cdot 0.6 \cdot 0.2 \cdot 0.6}{0.4 \cdot 0.5 + 0.6 \cdot 0.2}$
- B: $\frac{0.9 \cdot 0.4 \cdot 0.5 + 0.8 \cdot 0.6 \cdot 0.2}{0.4 \cdot 0.5 + 0.6 \cdot 0.2}$
- C: $\frac{0.9 \cdot 0.5 + 0.8 \cdot 0.2}{0.5 + 0.2}$
- D: $0.9 \cdot 0.4 + 0.8 \cdot 0.6$, E: none of the above

Bayesian Network

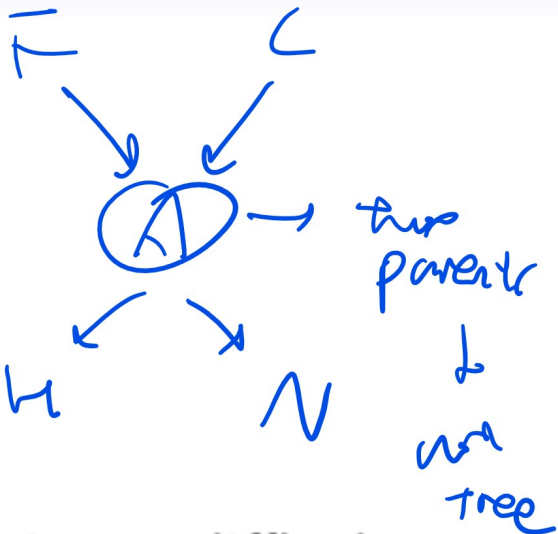
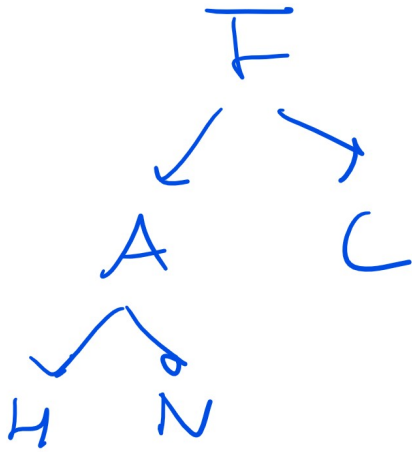
Algorithm

- Input: instances: $\{x_i\}_{i=1}^n$ and a directed acyclic graph such that feature X_j has parents $P(X_j)$.
- Output: conditional probability tables (CPTs): $\hat{\mathbb{P}}\{x_j|p(X_j)\}$ for $j = 1, 2, \dots, m$.
- Compute the transition probabilities using counts and Laplace smoothing.

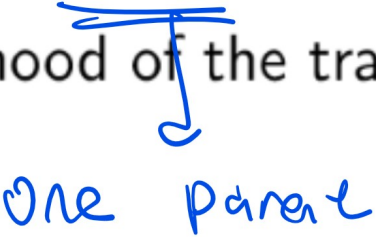
$$\hat{\mathbb{P}}\{x_j|p(X_j)\} = \frac{c_{x_j,p(X_j)} + 1}{c_{p(X_j)} + |X_j|}$$

Network Structure

Discussion

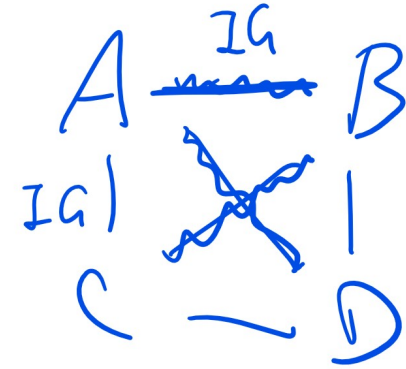


- Selecting from all possible structures (DAGs) is too difficult.
- Usually, a Bayesian network is learned with a tree structure.
- Choose the tree that maximizes the likelihood of the training data.



Chow Liu Algorithm

Discussion



decision tree

- Add an edge between features X_j and $X_{j'}$ with edge weight equal to the information gain of X_j given $X_{j'}$ for all pairs j, j' .
- Find the maximum spanning tree given these edges. The spanning tree is used as the structure of the Bayesian network.

Aside: Prim's Algorithm

Discussion

- To find the maximum spanning tree, start with an arbitrary vertex, a vertex set containing only this vertex, V , and an empty edge set, E .
- Choose an edge with the maximum weight from a vertex $v \in V$ to a vertex $v' \notin V$ and add v' to V , add an edge from v to v' to E
- Repeat this process until all vertices are in V . The tree (V, E) is the maximum spanning tree.

Aside: Prim's Algorithm Diagram

Discussion

Classification Problem

Discussion

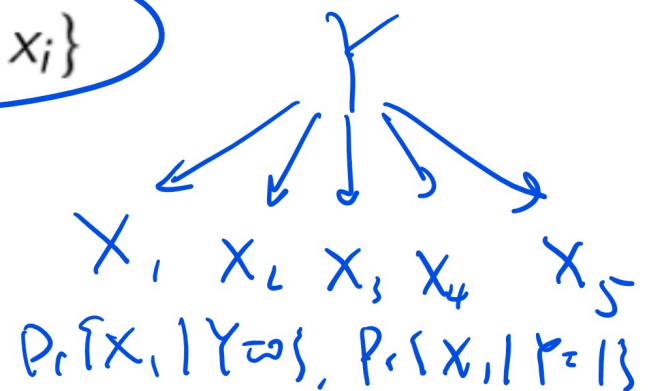
- Bayesian networks do not have a clear separation of the label Y and the features X_1, X_2, \dots, X_m .
- The Bayesian network with a tree structure and Y as the root and X_1, X_2, \dots, X_m as the leaves is called the Naive Bayes classifier.
- Bayes rules is used to compute $\mathbb{P}\{Y = y | X = x\}$, and the prediction \hat{y} is y that maximizes the conditional probability.



$P_r\{Y=1 | x\}$, $P_r\{Y=0 | x\}$

$$\hat{y}_i = \arg \max_y \mathbb{P}\{Y = y | X = x_i\}$$

Naive Bayes



Multinomial Naive Bayes

Discussion

- The implicit assumption for using the counts as the maximum likelihood estimate is that the distribution of $X_j|Y = y$, or in general, $X_j|P(X_j) = p(X_j)$ has the multinomial distribution.

$$\mathbb{P}\{X_j = x|Y = y\} = p_x$$
$$\hat{p}_x = \frac{c_{x,y}}{c_y}$$

Gaussian Naive Bayes Training

Discussion

- Training involves estimating $\mu_y^{(j)}$ and $\sigma_y^{(j)}$ since they completely determines the distribution of $X_j | Y = y$.
- The maximum likelihood estimates of $\mu_y^{(j)}$ and $(\sigma_y^{(j)})^2$ are the sample mean and variance of the feature j .

$$\hat{\mu}_y^{(j)} = \frac{1}{n_y} \sum_{i=1}^n x_{ij} \mathbb{1}_{\{y_i=y\}}, \quad n_y = \sum_{i=1}^n \mathbb{1}_{\{y_i=y\}}$$

$$(\hat{\sigma}_y^{(j)})^2 = \frac{1}{n_y} \sum_{i=1}^n (x_{ij} - \hat{\mu}_y^{(j)})^2 \mathbb{1}_{\{y_i=y\}}$$

sometimes $(\hat{\sigma}_y^{(j)})^2 \approx \frac{1}{n_y - 1} \sum_{i=1}^n (x_{ij} - \hat{\mu}_y^{(j)})^2 \mathbb{1}_{\{y_i=y\}}$

Gaussian Naive Bayes Diagram

Discussion

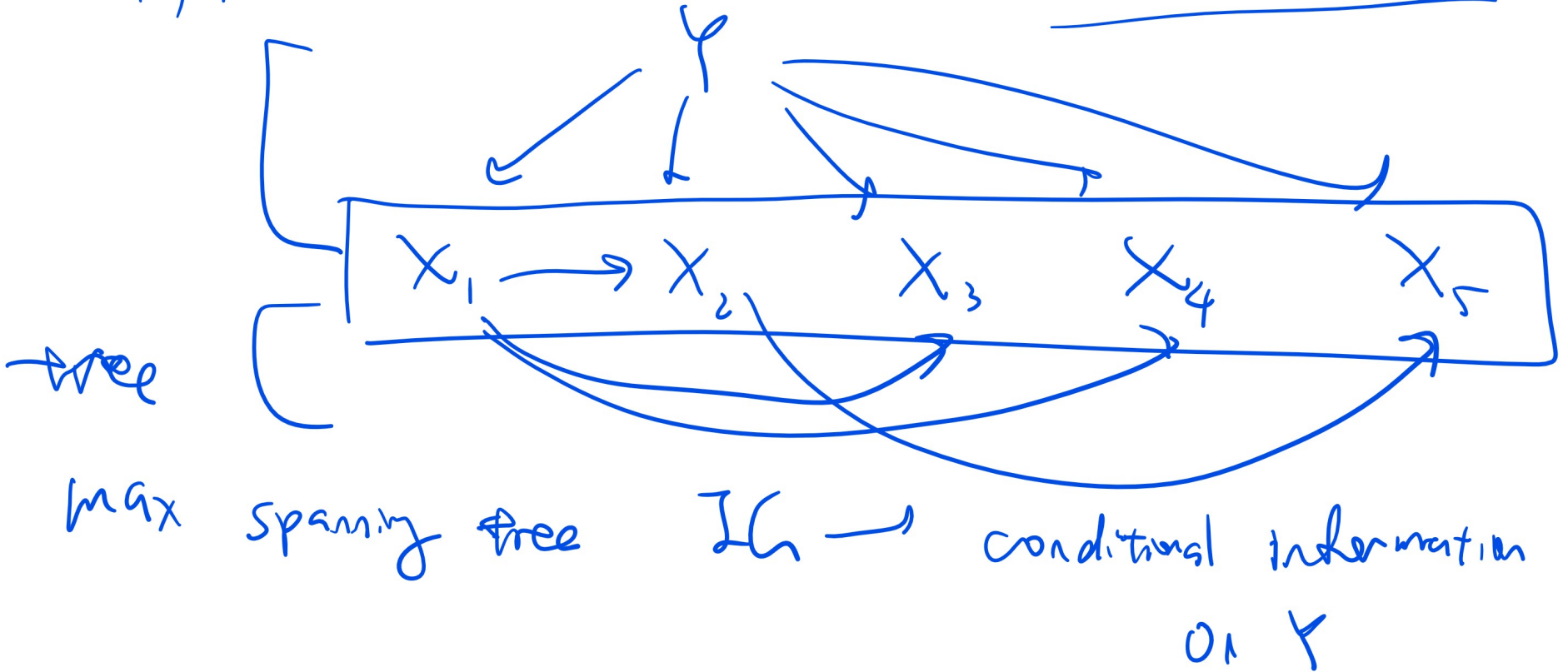
Tree Augmented Network Algorithm

Discussion

- It is also possible to create a Bayesian network with all features X_1, X_2, \dots, X_m connected to Y (Naive Bayes edges) and the features themselves form a network, usually a tree (MST edges).
- Information gain is replaced by conditional information gain (conditional on Y) when finding the maximum spanning tree.
- This algorithm is called TAN: Tree Augmented Network.

Tree Augmented Network Algorithm Diagram

TAN Discussion $P(Y | X_1 \dots X_5)$



End of Supervised learning

end of midterm coverage