

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.

$$\mathbb{P}\{X_j = x_j, X_{j'} = x_{j'}\} \in [0, 1]$$

\uparrow $0, 1, 2$ \uparrow $0, 1, 2$ \leftarrow 1

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

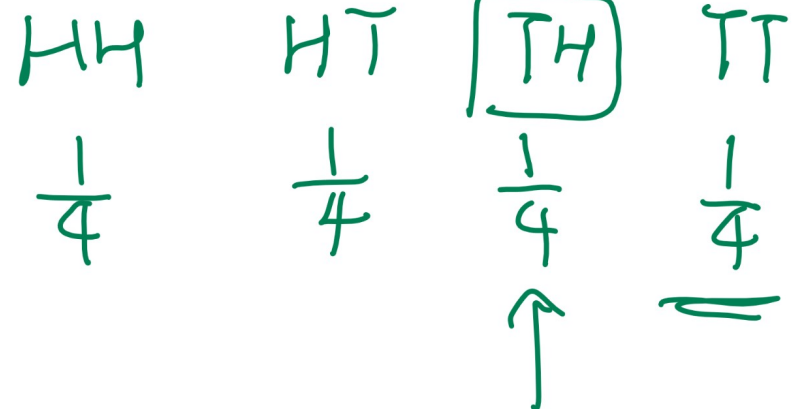
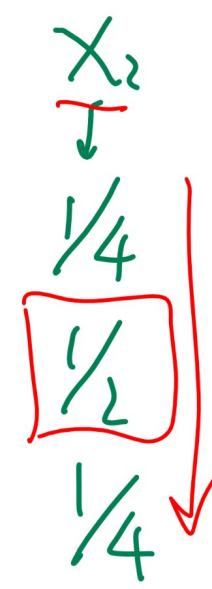
$$\mathbb{P}\{X_j = x_j\} = \sum_{x \in X_{j'}} \mathbb{P}\{X_j = x_j, X_{j'} = x\}$$

Distribution Example

first coin

		T	H
$X_1 =$	0	1	
$\rightarrow X_2 = 0$	$\frac{1}{4}$	$\frac{0}{4}$	
# H	1	$\frac{1}{4}$	$\frac{1}{4}$
2 whs	2	0	$\frac{1}{4}$
$X_1 \rightarrow$	$\frac{1}{2}$	$\frac{1}{2}$	

Motivation



$$P\{X_1=1 | X_2=1\} = \frac{1/4}{1/2} = \frac{1}{2}$$

$$P\{X_2=1 | X_1=1\} = \frac{1/4}{1/2} = \frac{1}{2}$$

Notation

Motivation

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

$$\begin{array}{ccc}
 \mathbb{P}\{x_j, x_{j'}\}, \mathbb{P}\{x_j\}, \mathbb{P}\{x_j|x_{j'}\} & & \mathbb{P}_r\{X_1=x_1, X_2=x_2\} \\
 \hline
 \mathbb{P}_c\{I, I\} & & \mathbb{P}_r\{x_1, x_2\}
 \end{array}$$

(Note: In the original image, blue arrows point from the handwritten $\mathbb{P}_c\{I, I\}$ and $\mathbb{P}_r\{x_1, x_2\}$ to the corresponding terms in the top row.)

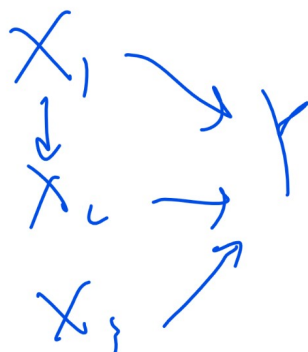
- When the context is not clear, for example when $x_j = a, x_{j'} = b$ with specific constants a, b , subscripts will be used under the probability sign.

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

(Note: In the original image, the subscripts $x_j, x_{j'}$ in the first term are circled in blue.)

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:

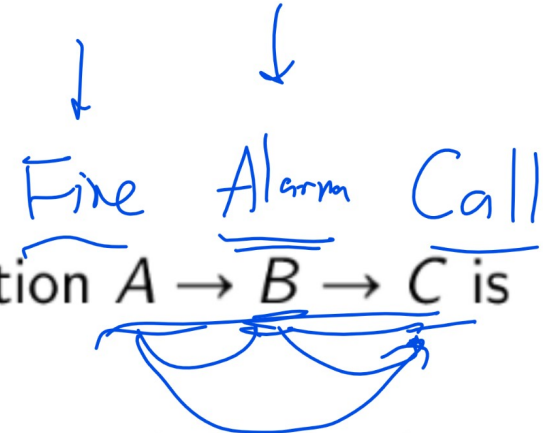
$$\mathbb{P}\{A, B\} = \mathbb{P}\{A\} \mathbb{P}\{B\} \text{ or } \mathbb{P}\{A|B\} = \mathbb{P}\{A\}$$

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

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

Causal Chain

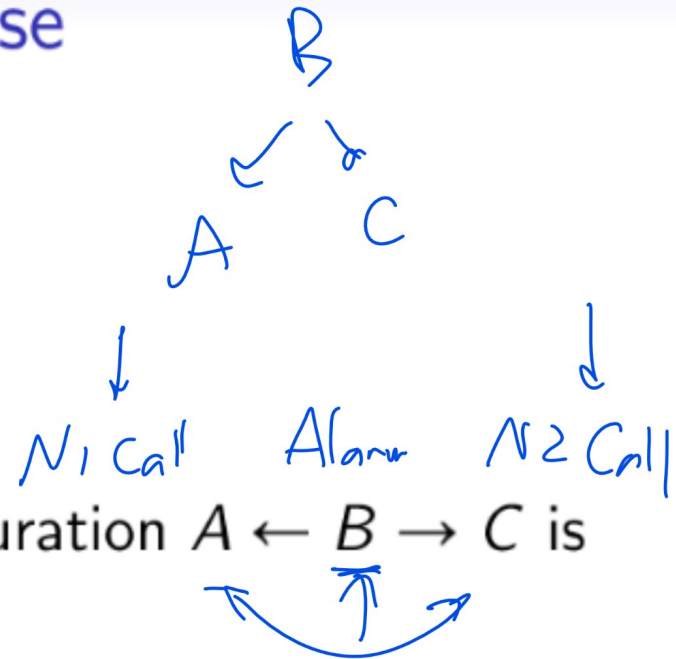
Definition



- For three events A, B, C , the configuration $A \rightarrow B \rightarrow C$ is called causal chain.
- 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 Cause

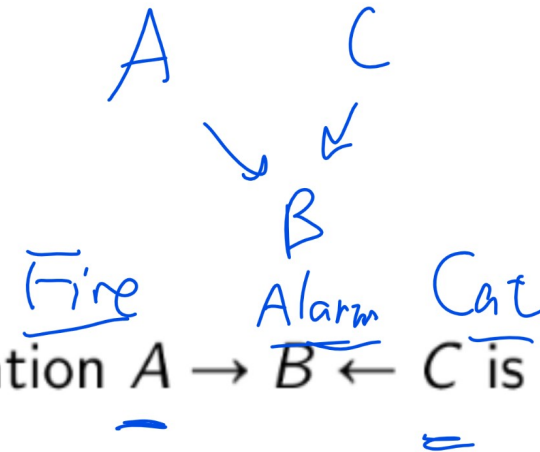
Definition



- 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 conditionally independent of C given information about B .
- Once B is observed, A and C are not independent.

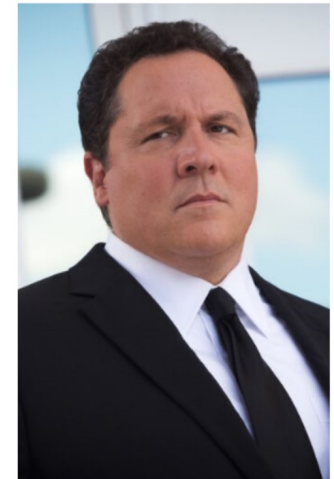
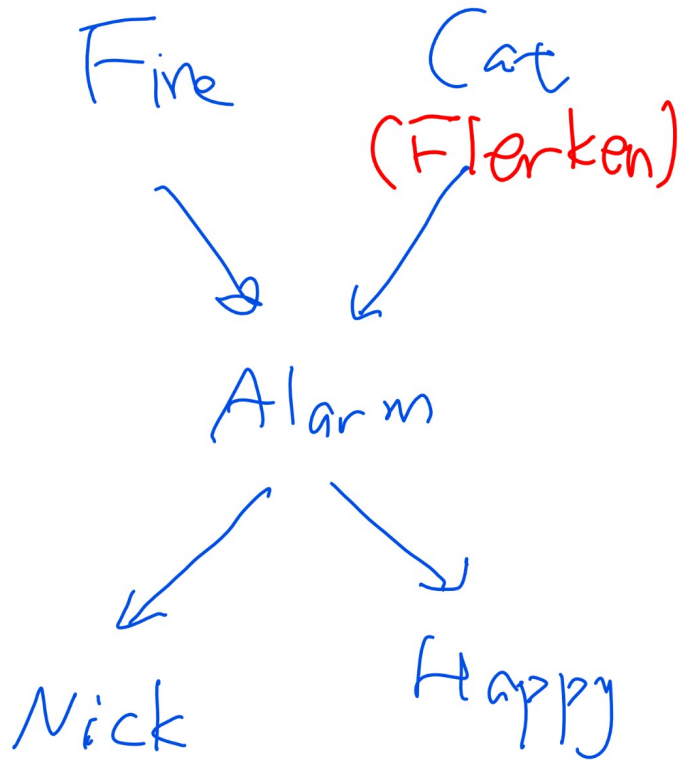
Storing Distribution

Definition

- 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.

Conditional Probability Table Diagram

Definition



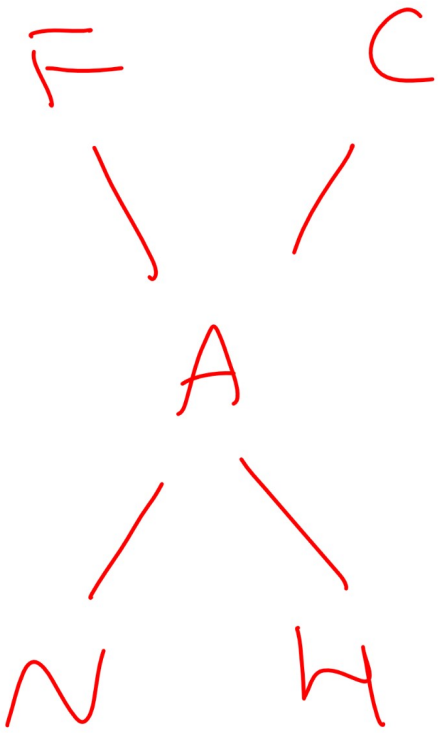
Conditional Probability Table Example

Definition

$$P_0 \{ \overline{F=0}, \overline{C=0}, A=0, N=0, H=0 \}$$

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

$$2^{\Sigma} = 32 - 1$$



- $P_0 \{ F \}$
- ~~$P_0 \{ F \}$~~
- $P_0 \{ C \}$
- $P_0 \{ A | FC \}, P_0 \{ A | \neg FC \}$
- $P_0 \{ A | F \neg C \}, P_0 \{ A | \neg F \neg C \}$
- $P_0 \{ N | A \}, P_0 \{ N | \neg A \}, P_0 \{ H | A \}, P_0 \{ H | \neg A \}$

~~$P_0 \{ A | FC \}$~~



10

Conditional Probability Table Larger Example

Definition

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, Part I

Definition

day
day
day

	F	C	A	H	N
1	0	0	0	1	0
2	0	1	0	0	0
3	0	0	0	1	1
	1	0	0	0	0
	0	0	1	1	0
	0	0	1	0	1
	0	0	1	1	1
	0	0	1	1	1

x_1
 x_2

$$Pr\{F\} = \frac{1}{8}$$

$$Pr\{H | A\} = \frac{3}{4}$$

$$Pr\{H | \neg A\} = \frac{2}{4} = \frac{1}{2}$$

$$Pr\{N | \neg A\} = \frac{2}{4} = \frac{1}{2}$$

$$Pr\{A | \neg F, \neg C\} = \frac{4}{6} = \frac{2}{3}$$

$$Pr\{A | F, C\} = \frac{0}{0} = ?$$

Laplace Smoothing

Definition

- Recall that the MLE estimation can incorporate Laplace smoothing.

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

$P. \{A | F, C\}$
 $= \frac{0 + 1}{0 + 2} = \frac{1}{2}$

- 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.

$\frac{0}{0}$

Bayes Net Inference

Definition

- Given the conditional probability table, the joint probabilities can be calculated using conditional independence.

$$\mathbb{P}\{x_1, x_2, \dots, x_m\} = \prod_{j=1}^m \mathbb{P}\{x_j | x_{j+1}, x_{j+2}, \dots, x_m\}$$

$$= \prod_{j=1}^m \mathbb{P}\{x_j | p(x_j)\}$$

- Given the joint probabilities, all other marginal and conditional probabilities can be calculated using their definitions.

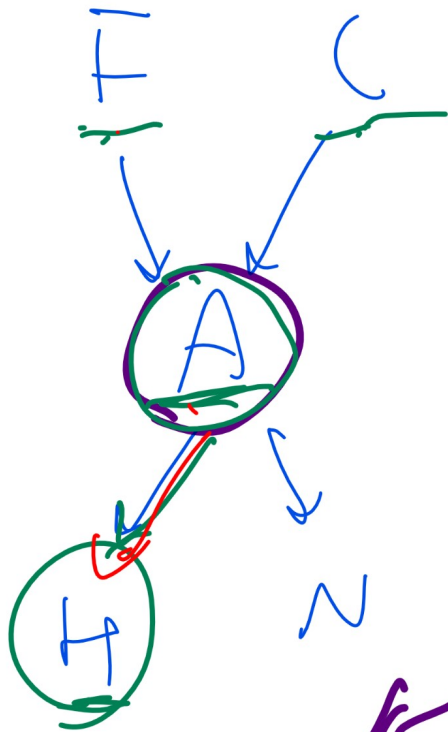
$$\mathbb{P}\{x_j | x_{j'}, x_{j''}, \dots\} = \frac{\mathbb{P}\{x_j, x_{j'}, x_{j''}, \dots\}}{\mathbb{P}\{x_{j'}, x_{j''}, \dots\}}$$

$$\mathbb{P}\{x_j, x_{j'}, x_{j''}, \dots\} = \sum_{x_k: k \neq j, j', j'', \dots} \mathbb{P}\{x_1, x_2, \dots, x_m\}$$

$$\mathbb{P}\{x_{j'}, x_{j''}, \dots\} = \sum_{x_k: k \neq j', j'', \dots} \mathbb{P}\{x_1, x_2, \dots, x_m\}$$

Bayes Net Inference Example, Part I

Definition



$$P_0 \{ H \mid \neg F, \neg C \}$$

$$\frac{P_1 \{ H, \neg F, \neg C \}}{P_0 \{ \neg F, \neg C \}}$$

$$= P_0 \{ \neg F \} \cdot P_0 \{ \neg C \} \\ = (1 - P_1 \{ F \}) (1 - P_1 \{ C \})$$

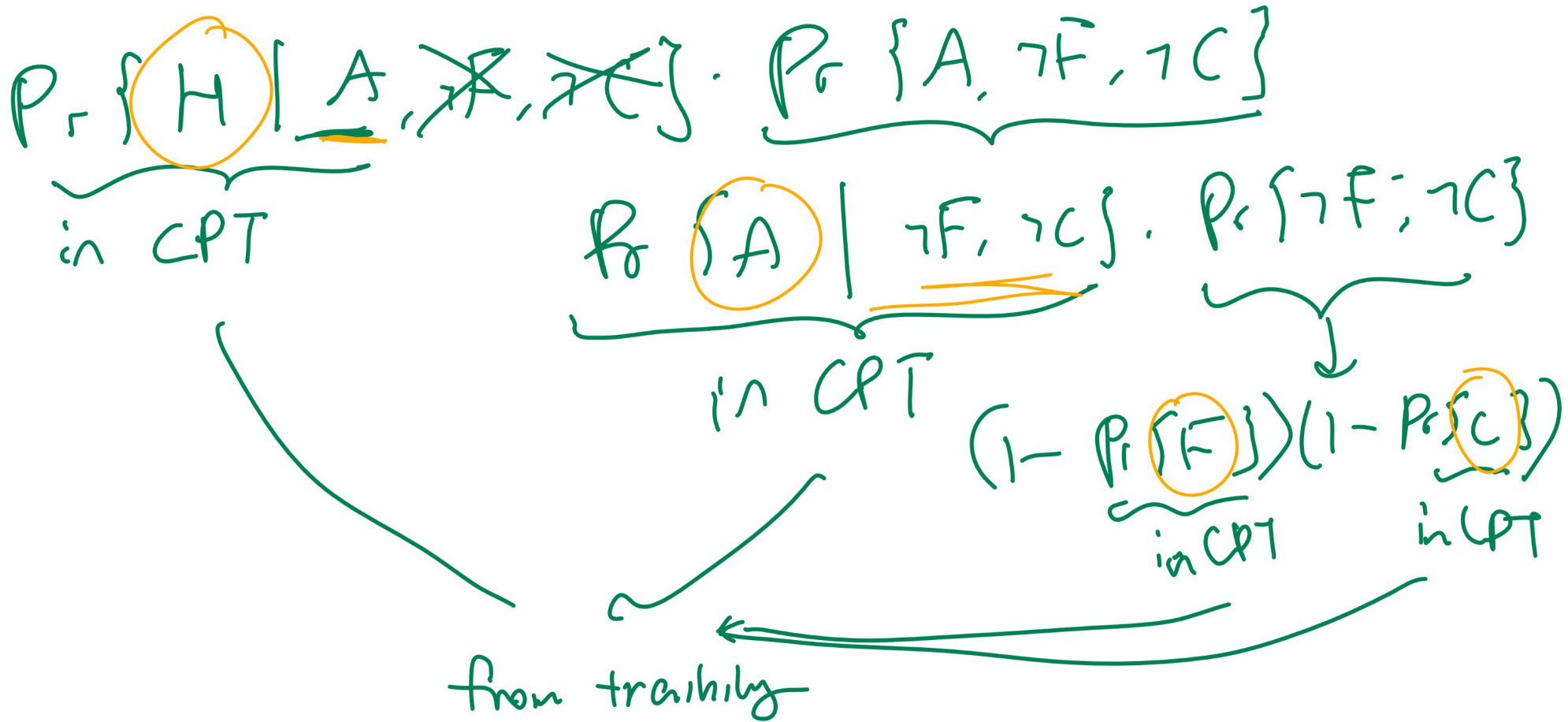
in CPT
in CPT

$$P_1 \{ H, A, \neg F, \neg C \} + P_1 \{ H, \neg A, \neg F, \neg C \}$$

$$P_0 \{ H \mid \neg A \} \cdot P_0 \{ \neg A \mid \neg F, \neg C \} \cdot P_0 \{ F \} \cdot P_0 \{ C \}$$

Bayes Net Inference Example, Part II

Definition



Bayes Net Inference Example, Part III

Definition

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

- 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.

577

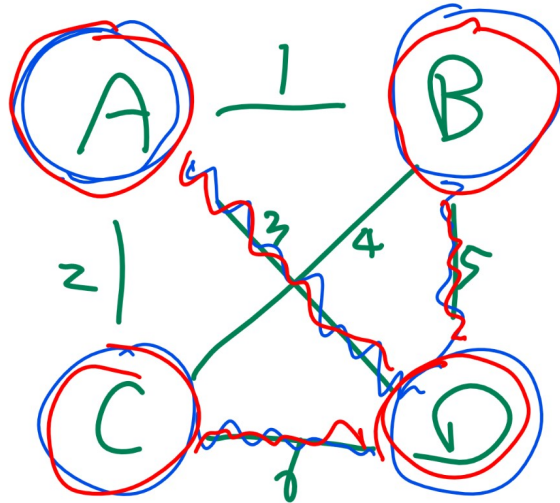
Aside: Prim's Algorithm

Discussion

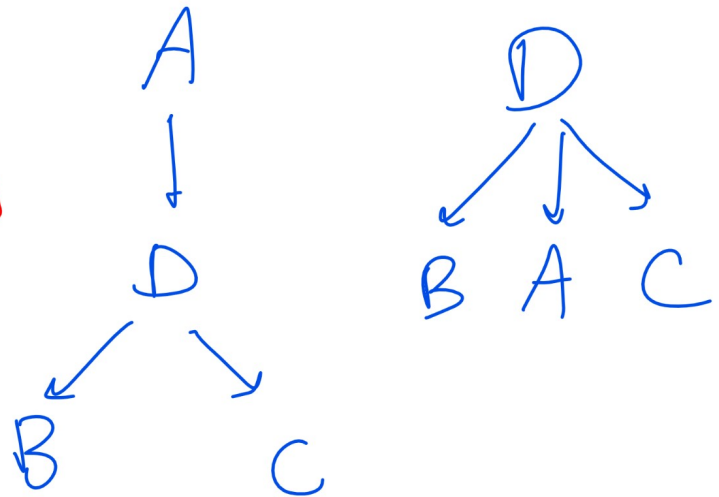
- To find the maximum spanning tree, start with an arbitrary vertex, a vertex set containing only this vertex, V , and an empty edge set, E .
- Choose an edge with the maximum weight from a vertex $v \in V$ to a vertex $v' \notin V$ and add v' to V , add an edge from v to v' to E
- Repeat this process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

Aside: Prim's Algorithm Diagram

Discussion



$$\text{Info Gain}(X_i, X_j) = \text{Info Gain}(X_j, X_i)$$



Classification Problem

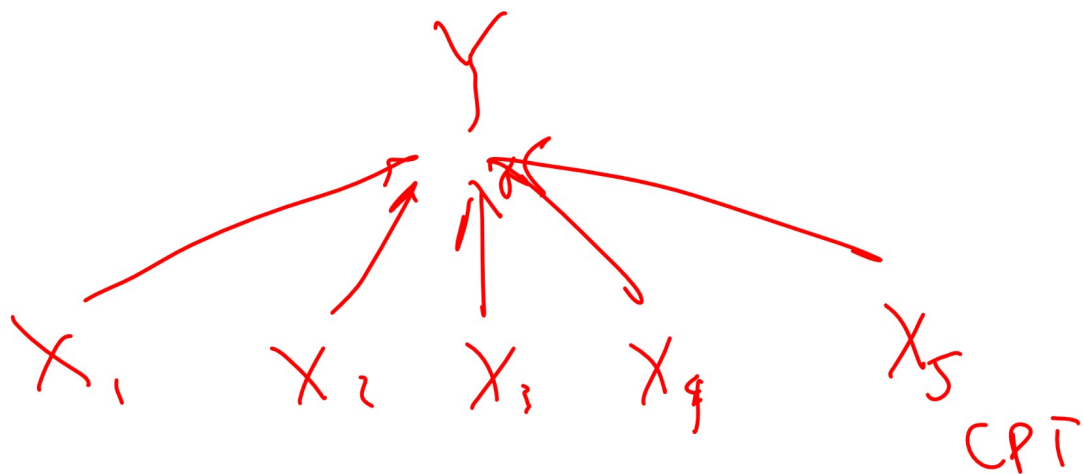
Discussion

- 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 rule is used to compute $\mathbb{P}\{Y = y|X = x\}$, and the prediction \hat{y} is y that maximizes the conditional probability.

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

Naive Bayes Diagram

Discussion



P3.

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.

$$\mathbb{P}\{X_j = x|Y = y\} = p_x$$

$$\hat{p}_x = \frac{c_{x,y}}{c_y}$$

1
2

Gaussian Naive Bayes

Discussion

- If the features are not categorical, continuous distributions can be estimated using MLE as the conditional distribution.
- Gaussian Naive Bayes is used if $X_j | Y = y$ is assumed to have the normal distribution.

$$\lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} \mathbb{P} \{x < X_j \leq x + \varepsilon | Y = y\} = \frac{1}{\sqrt{2\pi}\sigma_y^{(j)}} \exp\left(-\frac{(x - \mu_y^{(j)})^2}{2(\sigma_y^{(j)})^2}\right)$$

PDF

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 .

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

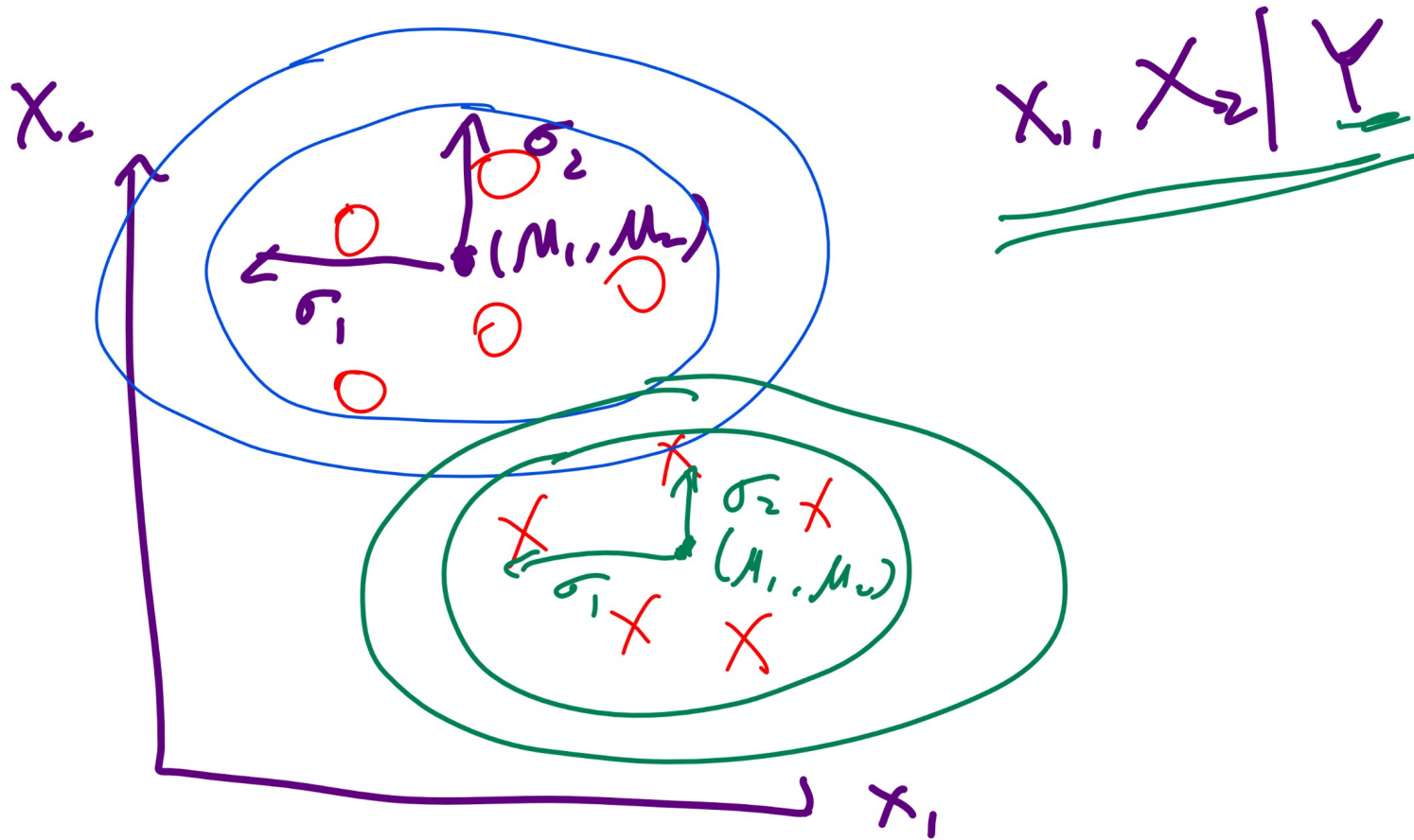
$(\hat{\sigma}_y^{(j)})^2 = \frac{1}{n_y} \sum_{i=1}^n (x_{ij} - \hat{\mu}_y^{(j)})^2 \mathbb{1}_{\{y_i=y\}}$

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

(Handwritten red annotations: circles around $\hat{\mu}_y^{(j)}$, x_{ij} , and $n_y - 1$; arrows pointing to $(\hat{\sigma}_y^{(j)})^2$; a large bracket on the right labeled 'MLE')

Gaussian Naive Bayes Diagram

Discussion



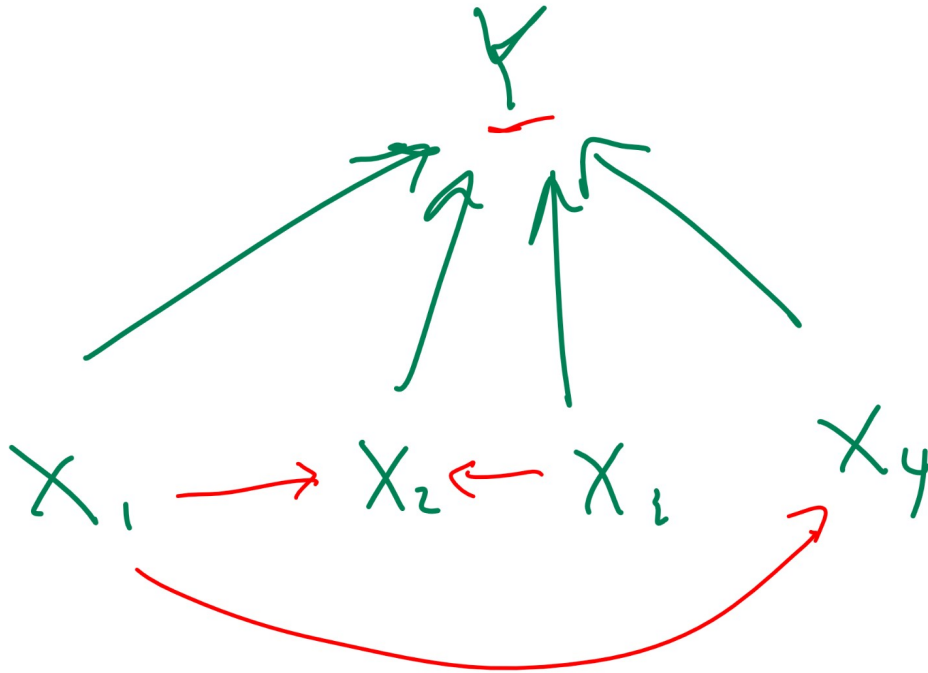
Tree Augmented Network Algorithm

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

Discussion



$Y \mid X_1, X_2, X_3, X_4$