Generative Models
оооооооооо

Bayesian Network
ооооооооооооооооооооооооооооооооо

Naive Bayes
ооооооо

# CS540 Introduction to Artificial Intelligence
## Lecture 8

Young Wu

Based on lecture slides by Jerry Zhu, Yingyu Liang, and Charles Dyer

July 13, 2022

Generative Models
●○○○○○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Rock Paper Scissors Game

Quiz

$$P \to \begin{pmatrix} 5 \\ 1 \\ 2 \end{pmatrix} \implies \frac{5}{8}$$

$R \to \implies \frac{1}{8}$

$S \to \implies \frac{2}{8}$

$Q1$

- Rock ($R$) beats Scissors ($S$) beats Paper ($P$) beats Rock ($R$).
- You: SPSRRSSSPRPPRPPSPSPRPPPRRSPPRS?
- AI: SPRPPSPRRPSSPPRPPSPSRRRSPRPSRS?

$\frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{3}$

- If AI uses the bigram model to guess your next action and chooses the best response to that action, which action should you choose next to win?

$\to S$

- $B : R, C : P, D : S$

game $\neq$ classification

**Generative Models**
○●○○○○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○●○○○○

# Regrade Request
## Admin

- If you submit $P1, M1$ etc late, please submit a regrade request.

- Please do not post solutions before I announce the assignment or after the due date of the assignment.

- You can participate in $D1, D1$ too, $D2, D2$ too, $D3, D3$ too group discussions until the midterm.

- No office hours this Friday.

- Replied to the notes in $Q11$. Thanks for the feedback!

**Generative Models**
○○●○○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○●○○

# Discriminative Model vs Generative Model

Motivation

$$\hat{y} = f(x)$$

$$\hat{y} = \underset{y \in \{0, 1, \dots\}}{\text{argmax}} \; Pr\{Y = y \mid X = x\}$$

$$0, 1$$

Cat, dog     image

$$f(x)$$

$$t$$
$$\Rightarrow MLE$$

$$\hat{Pr}\{X = x \mid Y = y\},$$

image     cat, dog

**Bayes Rule**

$$Pr\{Y = y \mid X = x\} = \frac{Pr\{Y = y, X = x\}}{Pr\{X = x\}}$$

$$= \frac{Pr\{X = x \mid Y = y\} \; Pr\{Y = y\}}{\sum_{y'} Pr\{X = x, Y = y'\}}$$

Generative Models
○○○○●○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Generative Models

Motivation

$$\Pr\{Y=y \mid X=x\} = \frac{\Pr[X=x \mid Y=y]\,\Pr\{Y=y\}}{\sum_{y'}\Pr[X=x \mid Y=y']\,\Pr\{Y=y'\}}$$

$$\hat{y} = \arg\max_{y} \Pr\{Y=y \mid X=x\}$$

- In probability terms, discriminative models are estimating $\mathbb{P}\{Y|X\}$, the conditional distribution. For example, $a_i \approx \mathbb{P}\{y_i = 1|x_i\}$ and $1 - a_i \approx \mathbb{P}\{y_i = 0|x_i\}$.

- Generative models are estimating $\mathbb{P}\{Y, X\}$, the joint distribution.

- Bayes rule is used to perform classification tasks.

$$\mathbb{P}\{Y|X\} = \frac{\mathbb{P}\{Y,X\}}{\mathbb{P}\{X\}} = \frac{\mathbb{P}\{X|Y\}\,\mathbb{P}\{Y\}}{\mathbb{P}\{X\}}$$

**Generative Models**
○○○○○●○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# 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.

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

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

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

Generative Models
○○○○○○●○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# 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.

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

Generative Models
○○○○○○○●○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayes Rule Example 1
## Quiz

- Two documents $A$ and $B$. Suppose $A$ contains 1 "Groot" and 9 other words, and $B$ contains 8 "Groot" and 2 other words. One document is taken out $A$ with probably $\frac{2}{3}$ and $B$ with probably $\frac{1}{3}$, and one word is picked out at random with equal probabilities. The word is "Groot". What is the probability that the document is $A$?

$$\left[ \begin{array}{l} Pr\{ G | A \} = \frac{1}{10} \\ Pr\{ G | B \} = \frac{8}{10} \end{array} \right.$$

$$Pr\{ \neg G | A \} = \frac{9}{10}$$

$$Pr\{ B \} = \frac{1}{3}$$

$$Pr\{ A \} = \frac{2}{3}$$

$$Pr\{ A | G \} = \frac{Pr\{ G | A \} \cdot Pr\{ A \}}{Pr\{ G | B \} Pr\{ B \} + Pr\{ G | A \} Pr\{ A \}}$$

$$\frac{1}{10} \qquad \frac{2}{3}$$

**Generative Models**
○○○○○○○●○

Bayesian Network
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayes Rule Example 1 Distribution

Quiz

join

$\Pr\{G|A\}, \Pr\{SA\},$

marginal

|     | A | B |
|-----|---|---|
| G | $\frac{2}{30}$ | $\frac{8}{30}$ |
| ¬G | $\frac{18}{30}$ | $\frac{2}{30}$ |

$\frac{10}{30}$

$\frac{20}{30}$

marginal

$\frac{20}{30}$   $\frac{10}{30}$

Generative Models
○○○○○○○○○●

Bayesian Network
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayes Rule Example 2

## Quiz

Q3

- Two documents $A$ and $B$. Suppose $A$ contains 1 "Groot" and 9 other words, and $B$ contains 8 "Groot" and 2 other words. One document is taken out at random (with equal probability), and one word is picked out at random (all words with equal probability). The word is "Groot". What is the probability that the document is $A$?

- $A : \dfrac{1}{9}$ , B: $\dfrac{1}{20}$ , C: $\dfrac{2}{5}$ , D: $\dfrac{9}{20}$ , E: I don't understand

$$\Pr\{A \mid G\} = \frac{\frac{1}{10} \cdot \frac{1}{2}}{\frac{1}{10} \cdot \frac{1}{2} + \frac{8}{10} \cdot \frac{1}{2}} = \frac{1}{9}$$

Generative Models
○○○○○○○○○

Bayesian Network
●○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# 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.

Generative Models
ooooooooo

Bayesian Network
o●ooooooooooooooooooooooooooooo

Naive Bayes
ooooooo

# 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:

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

Generative Models
○○○○○○○○○

Bayesian Network
○○●○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# 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.

Generative Models
○○○○○○○○○

Bayesian Network
○○○●○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# 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.

Generative Models
○○○○○○○○○

Bayesian Network
○○○○●○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Common Effect

## Definition

$A \quad C$

$B$

- 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.

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○●○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# 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 $P(X_j)$.

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

MLE

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

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

count in training set

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○●○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayesian Network Diagram

## Quiz

- Story: either Amber $(H)$ or Johnny's dog $(D)$ stepped on a bee, and put something on Johnny's bed $(B)$, and given there is something on Johnny's bed $(B)$, Johnny $(J)$ and Amber $(A)$ can be unhappy.

| H | D | B | J | A |
|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 |

*(handwritten annotations: day 1, day 2, day 3, Johnny, Amber)*

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○●○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayesian Network Diagram CPT Count

Quiz

$$Pr\{B=1|H=0\} = 1 - \boxed{\begin{array}{l} Pr\{B=0|H=0\} \\ Pr\{B=0|H=1\} \end{array}}$$

$$Pr\{B=1|H=1\}_{MLE}$$

$M \qquad D$

$B$

$J \qquad A$

2   $Pr\{B|H\} = \dfrac{C_{BH}}{C_H}$

2   $Pr\{B|D\}$

1   $Pr\{D\}$

1   $Pr\{H\}$

2   $Pr\{J|B\}$

2   $Pr\{A|B\}$

Conditional prob table

$\underline{\underline{CPT}}$

joint prob table.

$$Pr\{\underline{H}, D, B, J, A\} = 2^5 \text{ to store}$$
$$\quad\; 2 \quad 2 \quad 2 \quad 2 \quad 2$$

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○○●○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayes Net Training Example, Training

## Quiz

Q1
piek anything

- Given a network and the training data.
  $H \to B, D \to B, B \to J, B \to A.$

day 1
day 2

| H | D | B | J | A |
|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 |

H        D

B

J        A

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○○○●○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayes Net Training Example, Training 1
## Quiz

- Compute $\hat{\mathbb{P}}\{D = 1\}$.

| H | D | B | J | A |
|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 |

$$D \quad H$$
$$\searrow \quad \swarrow$$
$$B$$
$$\swarrow \quad \searrow$$
$$J \quad A$$

CPT

$\hat{\mathbb{P}}\{\underline{D} = 1\}$

$= \dfrac{\#D=1}{\#items} = \dfrac{2}{8}$

$= \dfrac{1}{4}$

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○○○○●○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayes Net Training Example, Training 2

## Quiz

- Compute $\hat{\mathbb{P}}\{J = 1|B = 1\}$

$\Pr\{J = 1 | B = 0\}$

$= 1 - \Pr\{J = 0 | B = 1\}$

$\sim \dfrac{\#_{J=1, B=1}}{\# B=1}$

$= \dfrac{3}{4}$

| H | D | B | J | A |
|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 |

B

J ↙ ↘ A

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○○○○○●○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayes Net Training Example, Training 3

## Quiz

- What is the conditional probability $\hat{\mathbb{P}}\{J = 1 | B = 0\}$?

- $A$ : I don't understand, B: $\dfrac{1}{4}$ , C: $\dfrac{1}{2}$ , D: $\dfrac{3}{4}$ , E: 1

| H | D | B | J | A |
|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 |

$Q2$

$$\frac{\#_{J=1,\,B=0}}{\#_{B=0}}$$

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

Generative Models
ooooooooo

Bayesian Network
ooooooooooooo●ooooooooooooooooo

Naive Bayes
ooooooo

# Bayes Net Training Example, Training 4

## Quiz

- Compute $\hat{\mathbb{P}}\{B = 1 | H = 0, D = 1\}$.

$$0 \quad 1$$
$$\boxed{0 \quad 1 \, 1}$$

$$H \quad D$$
$$\searrow \swarrow$$
$$B$$

$$\hat{p}_r\{B | \neg H, D\}$$
$$\Downarrow \quad \Downarrow \quad \Downarrow$$
$$B = 1 \quad H = 0 \quad D = 1$$

| H | D | B | J | A |
|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 |

$$\frac{\# B, \neg H, D}{\# \neg H, D}$$

$$= \frac{1}{2}$$

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○○○○○○○○●○○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayes Net Training Example, Training 5

## Quiz

- What is the conditional probability $\hat{\mathbb{P}} \{B = 1 | H = 0, D = 0\}$?

- $A$ : I don't understand, B: $\frac{1}{4}$ , C: $\frac{1}{2}$ , D: $\frac{3}{4}$ , E: 1

| H | D | B | J | A |
|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 |

Handwritten annotations:

Pr {B | H, D}

¬H, D
H, ¬D
¬H, ¬D

Q3

$\dfrac{\# B, \neg H, \neg D}{\# \neg H, \neg D}$

$\Rightarrow \dfrac{2}{4} = \dfrac{1}{2}$

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○○○○○○○○●○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayes Net Training Example, Training 5

## Quiz

- What is the conditional probability $\hat{\mathbb{P}}\{A = 0 | H = 1, D = 1\}$?

- $A$ : I don't understand, B: 0 , C: $\frac{1}{2}$ , D: 1 , E: NA

| H | D | B | J | A |
|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 |
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 |

$$\frac{0}{0}$$

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○○○○○○○○○●○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Laplace Smoothing

## Definition

- Recall that the MLE estimation can incorporate Laplace smoothing.

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

- 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.

*(handwritten annotations)*

# possible values of: $w_t$?

$\forall r \{w_t \mid w_{t-1}\}$

$\frac{\# w_t w_{t-1} + 1}{\# w_{t-1} + m}$

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○●○○○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayes Net Inference 1

### Definition

$$2^m$$

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

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

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

CPT

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○●○○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayes Net Inference 2
## Definition

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

$$\mathbb{P}\left\{x_j \mid x_{j'}, x_{j''}, \ldots\right\} = \frac{\mathbb{P}\left\{x_j, x_{j'}, x_{j''}, \ldots\right\}}{\mathbb{P}\left\{x_{j'}, x_{j''}, \ldots\right\}}$$

$$\mathbb{P}\left\{x_j, x_{j'}, x_{j''}, \ldots\right\} = \sum_{X_k : k \neq j, j', j'', \ldots} \mathbb{P}\left\{x_1, x_2, \ldots, x_m\right\}$$

$$\mathbb{P}\left\{x_{j'}, x_{j''}, \ldots\right\} = \sum_{X_k : k \neq j', j'', \ldots} \mathbb{P}\left\{x_1, x_2, \ldots, x_m\right\}$$

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○●○○○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayes Net Inference Example 1

## Quiz

- Assume the network is trained on a larger set with the following CPT. Compute $\hat{\mathbb{P}}\{H = 0, D = 1 | J = 1, A = 0\}$?

$$\hat{\mathbb{P}}\{H = 1\} = 0.001, \hat{\mathbb{P}}\{D = 1\} = 0.001$$

$$\hat{\mathbb{P}}\{B = 1 | H = 1, D = 1\} = 0.95, \hat{\mathbb{P}}\{B = 1 | H = 1, D = 0\} = 0.94$$

$$\hat{\mathbb{P}}\{B = 1 | H = 0, D = 1\} = 0.29, \hat{\mathbb{P}}\{B = 1 | H = 0, D = 0\} = 0.00$$

$$\hat{\mathbb{P}}\{J = 1 | B = 1\} = 0.9, \hat{\mathbb{P}}\{J = 1 | B = 0\} = 0.05$$

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

CPT

day n+1

| J | A | B | H | D |
|---|---|---|---|---|
| 1 | 0 | ? | 0 | 1 |

Generative Models
oooooooo

Bayesian Network
ooooooooooooooooooooo●ooooooooo

Naive Bayes
ooooooo

# Bayes Net Inference Example 1 Computation 1

Quiz

$$Pr\{\neg H, D \mid J, \neg A\}$$

not in CPT

$$= \frac{Pr\{\neg H, D, J, \neg A\}}{Pr\{J, \neg A\}}$$

$$Pr\{\neg H, D, J, \neg A\} = Pr\{\neg H, D, J, \neg A, B\}$$
$$+ Pr\{\neg H, D, J, \neg A, \neg B\}$$

$$\underset{0.999}{Pr\{\neg H\}} \cdot \underset{0.001}{Pr\{D\}} \cdot \underset{0.9}{Pr\{J \mid B\}} \cdot \underset{0.3}{Pr\{\neg A \mid B\}} \cdot \underset{0.29}{Pr\{B \mid \neg H, D\}}$$

$$Pr\{\neg H\} \cdot Pr\{D\} \cdot Pr\{J \mid \neg B\} \cdot Pr\{\neg A \mid \neg B\} \cdot Pr\{\neg B \mid \neg H, D\}$$

Generative Models
○○○○○○○○○

**Bayesian Network**
○○○○○○○○○○○○○○○○○○○○○●○○○○○○○○○

Naive Bayes
○○○○○○○

# Bayes Net Inference Example 1 Computation 2

Quiz

$$\Pr\{J, \neg A\} = \quad \Pr\{J, \neg A, \quad H, D, B$$

$$\nearrow \quad \prod_j \Pr\{x_j | pc(x_j)\}$$

$$+ \quad \underline{\quad} \qquad H, D, \neg B$$

$$H, \neg D, B$$

$$H, \neg D, \neg B$$

$$\vdots$$

Bayes Ball

Generative Models
ooooooooo

Bayesian Network
oooooooooooooooooooooo●ooooooo

Naive Bayes
ooooooo

# Bayes Net Inference Example 2

## Quiz

$$D \qquad H$$

$$\searrow \swarrow$$

$$B$$

$$\hat{\mathbb{P}} \{D, \neg H\}$$

$$\overline{\mathbb{P}\{\neg H\}}$$

- Compute $\hat{\mathbb{P}}\{D = 1 | H = 0\}$?

$$\hat{\mathbb{P}}\{H = 1\} = 0.001, \hat{\mathbb{P}}\{D = 1\} = 0.001$$

$$\hat{\mathbb{P}}\{B = 1 | H = 1, D = 1\} = 0.95, \hat{\mathbb{P}}\{B = 1 | H = 1, D = 0\} = 0.94$$

$$\hat{\mathbb{P}}\{B = 1 | H = 0, D = 1\} = 0.29, \hat{\mathbb{P}}\{B = 1 | H = 0, D = 0\} = 0.00$$

- $A : 0,$ B: $0.001,$ C: $0.0094,$ D: $0.0095,$ E: $1$

# Bayes Net Inference Example 2 Derivation

## Quiz

Generative Models
ooooooooo

Bayesian Network
ooooooooooooooooooooooooooo●oooooo

Naive Bayes
ooooooo

# Bayes Net Inference Example 3

## Quiz

$P_r\{\neg H, D, B\}$ ~ $J, A$

$P_r\{\neg H\} \cdot P_r\{D\} \cdot P_r\{B \mid \neg H, D\}$

$$\frac{P_r\{\neg H, D, B\}}{P_r\{B\}}$$

- Compute $\hat{\mathbb{P}}\{H = 0, D = 1 \mid B = 1\}$?

$$\hat{\mathbb{P}}\{H = 1\} = 0.001, \hat{\mathbb{P}}\{D = 1\} = 0.001$$

$$\hat{\mathbb{P}}\{B = 1 \mid H = 1, D = 1\} = 0.95, \hat{\mathbb{P}}\{B = 1 \mid H = 1, D = 0\} = 0.94$$

$$\hat{\mathbb{P}}\{B = 1 \mid H = 0, D = 1\} = 0.29, \hat{\mathbb{P}}\{B = 1 \mid H = 0, D = 0\} = 0.00$$

$\neg H, D$
$\neg H, \neg D$
$H, \neg D$
$H D$

- $A : 0$, B: $0.001$, C: $0.0094$, D: $0.0095$, E: $1$

# Bayes Net Inference Example 3 Derivation

## Quiz

Generative Models
ooooooooo

Bayesian Network
ooooooooooooooooooooooooooo●ooo

Naive Bayes
ooooooo

# Bayes Net Inference Example 4
## Quiz

- Compute $\hat{\mathbb{P}}\{B=1|J=1, A=0\}$?

$$\hat{\mathbb{P}}\{J=1|B=1\}=0.9, \hat{\mathbb{P}}\{J=1|B=0\}=0.05$$

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

Given

$$\mathbb{P}\{B=1\}=0.001 \cdot 0.001 \cdot 0.95 + 0.001 \cdot 0.999 \cdot (0.94 + 0.29).$$
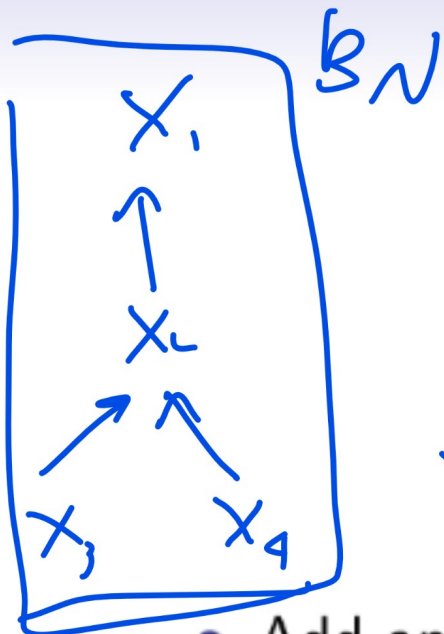
- $A:0$, B: 0.001, C: 0.0094, D: 0.0095, E: 1

# Bayes Net Inference Example 4 Derivation
## Quiz

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○●○

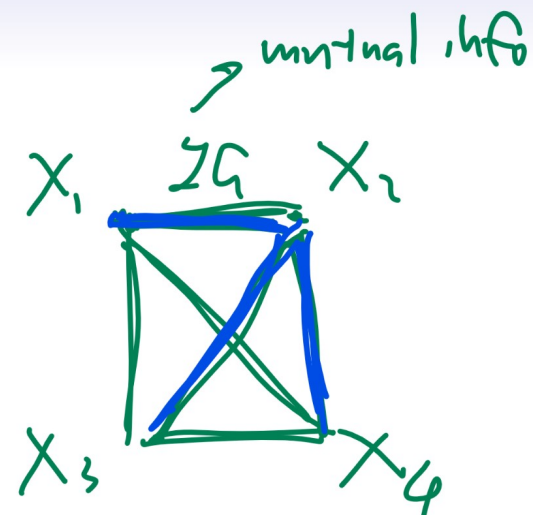Naive Bayes
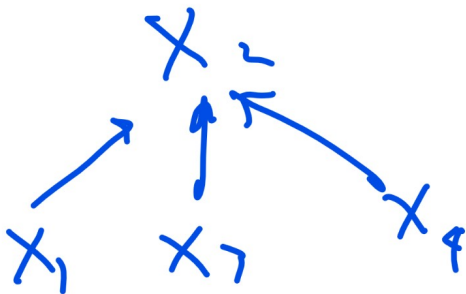○○○○○○○

# 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.

Generative Models
OOOOOOOOO

Bayesian Network
OOOOOOOOOOOOOOOOOOOOOOOOOOOOOOO●

Naive Bayes
OOOOOOO

# 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.

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
●○○○○○○

# Classification Problem
## Discussion

- Bayesian networks do not have a clear separation of the label $Y$ and the features $X_1, X_2, ..., X_m$.
- The Bayesian network with a tree structure and $Y$ as the root and $X_1, X_2, ..., X_m$ as the leaves is called the Naïve 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.
$$\hat{y}_i = \operatorname*{argmax}_y \mathbb{P}\{Y = y | X = x_i\}$$

CPT $\quad$ $Pr\{X_i | Y\}$. $\quad$ Bayes Rule $\quad$ Naive Bayes.

$Y$

$X_1 \quad X_2 \quad X_3 \quad X_4$

Generative Models
ooooooooo

Bayesian Network
ooooooooooooooooooooooooooooooo

Naive Bayes
o●oooooo

# Naive Bayes Diagram

## Discussion

$$\Pr\{Y=y \mid X_j = x\} = \frac{\prod_j \Pr\{X_j = x \mid Y = y\} \cdot \Pr\{Y=y\}}{\sum_{y'} \prod_j \Pr\{X_j = x \mid Y = y'\} \cdot \Pr\{Y=y'\}}$$

CPT    CPT

PS

Generative Models
○○○○○○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○●○○○○

# Multinomial Naive Bayes

### Discussion

$Y \sim 0, 1$

$X_j \sim 0, 1$

- 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.

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

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

$\frac{c_{x,y}}{c_{py}}$

Generative Models
ooooooooo

Bayesian Network
ooooooooooooooooooooooooooooo

Naive Bayes
ooo●ooo

# Gaussian Naive Bayes

## Discussion

$\hat{P}\{x, y\}$

- 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.

PDF

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

Generative Models
○○○○○●○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○●○○

# Gaussian Naive Bayes Training
## Discussion

- Training involves estimating $\mu_y^{(j)}$ and $\sigma_y^{(j)}$ since they completely determine the distribution of $X_j|Y = y$.

- The maximum likelihood estimates of $\mu_y^{(j)}$ and $\left(\sigma_y^{(j)}\right)^2$ are the sample mean and variance of the feature $j$.

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

MLE

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

MLE

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

# Tree Augmented Network Algorithm

### Discussion

- It is also possible to create a Bayesian network with all features $X_1, X_2, ..., X_m$ 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