Probability Distributions
OOOO

Bayesian Network
OOOOOOOOOOOOOOOOOOOO

Naive Bayes
OOOO

# CS540 Introduction to Artificial Intelligence
# Lecture 10

Young Wu

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

June 15, 2020

**Probability Distributions**
●○○○

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

Naive Bayes
○○○○

# Bayes Rule Example 1

## Quiz

Q1

( any answer is okay ).

- 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: $\dfrac{1}{2}$ , B: $\dfrac{1}{3}$ , C: $\dfrac{1}{4}$ , D: $\dfrac{1}{8}$ , E: $\dfrac{1}{9}$

$\dfrac{1}{9}$ $\dfrac{1}{9}$

**Probability Distributions**
○●○○

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

Naive Bayes
○○○○

# Bayes Rule Example 1 Distribution

Quiz

$$Y = G \qquad Y = N$$

joint probability

$$
\begin{array}{c|cc}
 & Y=G & Y=N \\
0.5 \ X=A & \boxed{0.05} & 0.45 \\
0.5 \ X=B & 0.4 & 0.1
\end{array}
$$

joint distribution

$$\boxed{0.45} \qquad 0.55 \leftarrow \text{marginal distribution of } Y.$$

marg. of $X$

$$\Pr\{Y = G \mid X = A\} = \frac{1}{10} = \frac{\Pr\{Y=G, X=A\}}{\Pr\{X=A\}}$$

$$\Pr\{X = A \mid Y = G\} = \frac{\Pr\{Y=G, X=A\}}{\Pr\{Y=G\}} = \frac{0.05}{0.5} = 0.1$$

$$= \frac{0.05}{0.45} = \boxed{\frac{1}{9}}$$

**Probability Distributions**
○○●○

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

**Naive Bayes**
○○○●

# Bayes Rule Example 2
## Quiz

Q2

- 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 $\frac{1}{3}$ and $B$ with probably $\frac{2}{3}$, and one word is picked out at random with equal probabilities. The word is "Groot". What is the probability that the document is $A$?

- A: $\frac{1}{9}$ , B: $\frac{1}{10}$ , C: $\frac{1}{16}$ , D: $\frac{1}{17}$ , E: $\frac{1}{25}$

$\frac{1}{17}$

**Probability Distributions**
○○○●

Bayesian Network
○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○○○

# Bayes Rule Example 2 Distribution

Quiz

$$\Pr\{X=A \mid Y=G\} = \frac{\Pr\{X=A, Y=G\}}{\Pr\{Y=G\}}$$

posterior

and

$$= \frac{\Pr\{Y=G \mid X=A\} \cdot \Pr\{X=A\}}{\Pr\{Y=G, X=A\} + \Pr\{Y=G, X=B\}}$$

likelihood    prior

marginalization

Bayes Rule

$$= \frac{\Pr\{Y=G \mid X=A\} \cdot \Pr\{X=A\}}{\Pr\{Y=G \mid X=A\} \cdot \Pr\{X=A\} + \Pr\{Y=G \mid X=B\} \cdot \Pr\{X=B\}}$$

$$= \frac{0.1 \cdot \frac{1}{3}}{0.1 \cdot \frac{1}{3} + 0.8 \cdot \frac{2}{3}} = \frac{1}{1 + 8 \cdot 2} = \frac{1}{17}.$$

Probability Distributions
○○○○

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

Naive Bayes
○○○○

# Bayesian Network
## Definition

$$A$$
$$B \to C$$

- 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.

$$Pr\{X_j = a, X_{j'} = b \mid X_{j''} = c\}$$
$$\overset{indep}{=} Pr\{X_j = a\} \cdot Pr\{X_{j'} = b\}$$
$$\mid X_{j''} \qquad\qquad \mid X_{j''}$$

Probability Distributions
○○○○

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 $P(X_j)$.

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

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

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

$$\frac{\# \quad \text{parent}}{\text{Pr}\{x_j, x_{j'}\}} \leftarrow \text{val of } x_j$$

$$\frac{}{\text{Pr}\{x_{j'}\}}$$

$$\#$$

Probability Distributions
○○○○

Bayesian Network
○○●○○○○○○○○○○○○○○○○

Naive Bayes
○○○○

# Bayesian Network Diagram

## Quiz

- 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 (Fury) or Happy or both might call you because of the alarm or for other reasons.

weights
parameter

cond prob
of children
given parents.

day 1
day 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 |

F    C
 ↘  ↙
   A
 ↙  ↘
N      H

$2^5 = 32$

Probability Distributions
OOOO

Bayesian Network
OOO●OOOOOOOOOOOOOOOOO

Naive Bayes
OOOO

# Bayes Net Training Example, Training 1

## Quiz

- Compute $\hat{\mathbb{P}}\{C = 1\}$. $= \dfrac{1}{8}$

| 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 |

Probability Distributions
○○○○

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

Naive Bayes
○○○○

# Bayes Net Training Example, Training 2

## Quiz

- Compute $\hat{\mathbb{P}}\{N = 1 | A = 1\}. = \dfrac{C_{N=1, A=1}}{C_{A=1}} = \dfrac{3}{4}$

| 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 |

Probability Distributions
○○○○

Bayesian Network
○○○○○●○○○○○○○○○○○○○

Naive Bayes
○○○○

# Bayes Net Training Example, Training 3

## Quiz

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

$$= \frac{C_{A, \neg F, C}}{C_{\neg F, C}}$$

| 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 |

$$= \frac{0}{1}$$

$$= 0$$

Probability Distributions
○○○○

Bayesian Network
○○○○○○●○○○○○○○○○○○○

Naive Bayes
○○○○

# Bayes Net Training Example, Training 4

## Quiz

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

Q3

| 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 |

Probability Distributions
oooo

Bayesian Network
oooooooo●oooooooooooo

Naive Bayes
oooo

# Bayes Net Training Example, Training 5

## Quiz

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

- A: 0 , B: $\dfrac{1}{3}$ , C: $\dfrac{1}{2}$ , D: $\dfrac{2}{3}$ , E: 1

$A = 1$

$1 - 0 = 1$

Q4

| 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 |

Probability Distributions
OOOO

Bayesian Network
OOOOOOOO●OOOOOOOOOOO

Naive Bayes
OOOO

# 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.

Probability Distributions
oooo

Bayesian Network
ooooooooo●oooooooooo

Naive Bayes
oooo

# Bayes Net Inference 1
## Definition

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

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

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

$X_1$

Probability Distributions
oooo

Bayesian Network
ooooooooooo●ooooooooo

Naive Bayes
oooo

# Bayes Net Inference 2
## Definition

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

$$\mathbb{P}\left\{x_j | x_{j'}, x_{j''}, ...\right\} = \frac{\mathbb{P}\left\{x_j, x_{j'}, x_{j''}, ...\right\}}{\mathbb{P}\left\{x_{j'}, x_{j''}, ...\right\}}$$

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

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

Probability Distributions
OOOO

Bayesian Network
OOOOOOOOOOO●OOOOOOOO

Naive Bayes
OOOO

# Bayes Net Inference Example 1

## Quiz

$$\overline{F} \qquad C$$

$$A$$

$$H \qquad N$$

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

$$\hat{\mathbb{P}}\{F = 1\} = 0.001, \hat{\mathbb{P}}\{C = 1\} = 0.001$$
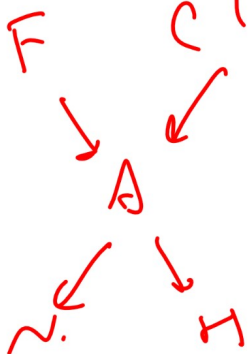
$$\hat{\mathbb{P}}\{A = 1 | F = 1, C = 1\} = 0.95, \hat{\mathbb{P}}\{A = 1 | F = 1, C = 0\} = 0.94$$

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

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

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

$$1 \, 0$$

$$F \qquad C$$

$$A$$

$$N \qquad H$$

Probability Distributions
○○○○

Bayesian Network
○○○○○○○○○○○○●○○○○○○

Naive Bayes
○○○○

# Bayes Net Inference Example 1 Computation 1

$$R=1 \qquad C=1 \qquad H=0 \qquad N=0 \text{ Quiz}$$

$$P_r\{F, C, | \neg H, \neg N\} = \left\{ \frac{P_r\{F, C, \neg H, \neg N\}}{\boxed{P_r\{\neg H, \neg N\}}} \right.$$

$$\hookrightarrow P_r\{F, C, \neg H, \neg N, A\} + \boxed{P_r\{F, C, \neg H, \neg N, \neg A\}}$$

marginal

$$= P_r\{\cancel{F}\} \cdot P_r\{\cancel{C}\} \cdot P_r\{\cancel{\neg H} | A\} \cdot P_r\{\neg N | A\} \cdot P_r\{A | F, C\}$$

$$0.001 \qquad 0.001 \qquad (1 - 0.5) \qquad (1 - 0.7) \cdot 0.95$$

$$+ \ 0.001 \cdot 0.001 \cdot (1 - 0.05) \ (1 - 0.01) \cdot (1 - 0.95)$$

Probability Distributions
oooo

Bayesian Network
oooooooooooooo●oooooo

Naive Bayes
oooo

# Bayes Net Inference Example 1 Computation 2

## Quiz

$$Pr\{\neg M, \neg N\} = Pr\{F, C, A, \neg M, \neg N\}$$

$$\bar{F} \quad C, \neg A$$
$$\bar{F} \quad \neg C, A$$
$$F \quad \neg C, \neg A$$
$$\neg F \quad C \quad A$$
$$\neg F \quad C \quad \neg A$$
$$\neg F \quad \neg C \quad A$$
$$\neg F \quad \neg C \quad \neg A$$

Probability Distributions

0000

Bayesian Network

0000000000000●00000

Naive Bayes

0000

# Bayes Net Inference Example 1 Computation 3

## Quiz

Probability Distributions
oooo

Bayesian Network
ooooooooooooooo●oooo

Naive Bayes
oooo

# Bayes Net Inference Example 2

Quiz

Q5 (last)

$P_r \{\neg H, \neg N, A\}$

$= P_r \{\neg H, \emptyset A\}$

$P_s \{\neg N | A\}$

$P_r \{A | \begin{matrix} F, C \\ \neg F, \neg C \\ F, \neg C \\ \neg F, C \end{matrix}\}$

- Compute $\hat{\mathbb{P}}\{C = 1 | F = 1\}$?

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

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

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

$P_r \{A\} < \dfrac{\#A}{\#n}$

$\dfrac{C\ F\ A + \neg F\ \neg A}{F}$

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

$$\frac{P_r \{C, \bar{F}\}}{P_n \{\bar{F}\}} = \frac{P_r \{C\} P_r \{\bar{F}\}}{P_r \{\bar{F}\}} = 0.001$$

Probability Distributions
○○○○

Bayesian Network
○○○○○○○○○○○○○○○●○○○

Naive Bayes
○○○○

# Bayes Net Inference Example 2 Computation
## Quiz

- Compute $\hat{\mathbb{P}}\{C = 1|F = 1\}$?

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

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

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

Probability Distributions
○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○●○○

Naive Bayes
○○○○

# Bayes Net Inference Example 3

## Quiz

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

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

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

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

- A: $0.001 \cdot 0.001$, B: $0.001 \cdot 0.001 \cdot 0.95$,

- C: $\dfrac{0.001}{0.001 \cdot 0.95 + 0.999 \cdot (0.94 + 0.29)}$

- D : $\dfrac{0.001 \cdot 0.001}{0.001 \cdot 0.95 + 0.999 \cdot (0.94 + 0.29)}$

- E : $\dfrac{0.001 \cdot 0.95}{0.001 \cdot 0.95 + 0.999 \cdot (0.94 + 0.29)}$

Probability Distributions
○○○○

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

Naive Bayes
○○○○

# Bayes Net Inference Example 3 Computation
## Quiz

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

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

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

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

Probability Distributions
○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○●

Naive Bayes
○○○○

# 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.

# 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 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.

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

# Naive Bayes Diagram

## Discussion

Probability Distributions
○○○○

Bayesian Network
○○○○○○○○○○○○○○○○○○○

Naive Bayes
○○●○

# 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.

Probability Distributions
ooooo

Bayesian Network
ooooooooooooooooooooo

Naive Bayes
oooo●

# Tree Augmented Network Algorithm Diagram
## Discussion