

IS590DT: Data Mining Applications

WEEK 11: NB, BN, SP and OCR

Spring 2019

Yingjun Guan

yingjun2@illinois.edu

<http://ischool.illinois.edu/people/yingjun-guan>



ILLINOIS

School of
Information Sciences
The iSchool at Illinois

Contents today



- Naïve Bayes
- Bayesian Network
- OCR application {Reproducibility}
- Structural Prediction
- Summary

Let's start with how classification works...



- The classification of classification algorithms (How indeed do we make the final classification?)
- Let's check the following machine learning algorithm summarizing figure first.

Which one is classification?





Bayesian Theorem...

- Conditional probability.

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

Probability of event A occurred
and event B occurred

Probability of event A
given B has occurred

Probability of event B

- EG: $P(\text{healthy} | \text{positive result})$

Truth	Test result		Total
	Positive	Negative	
Renal disease	44	23	67
Healthy	10	60	70
Total	54	83	137

Bayesian Theorem...



- X sometimes **buys an ice cream** after DT classes, depending on the the assignment Yingjun leaves.
- If it is an **easy** assignment (50% to happen), 90% likelihood to buy.
- If it is a **hard** assignment (50% to happen), 30% likelihood to buy.
- After today's class, X buys an ice cream, what do you think of today's assignment? **Easy? Hard?**

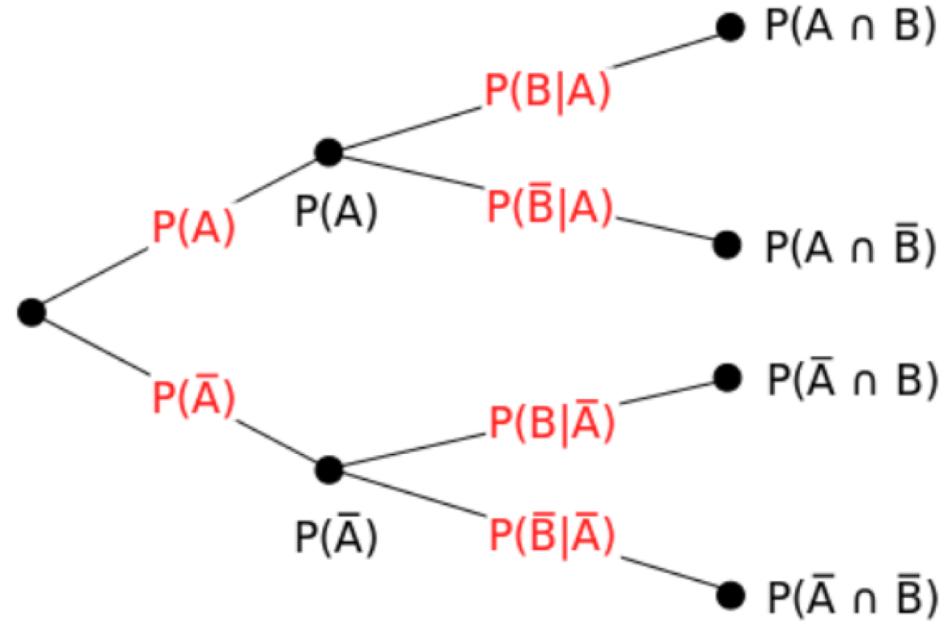
Formula and Tree structure



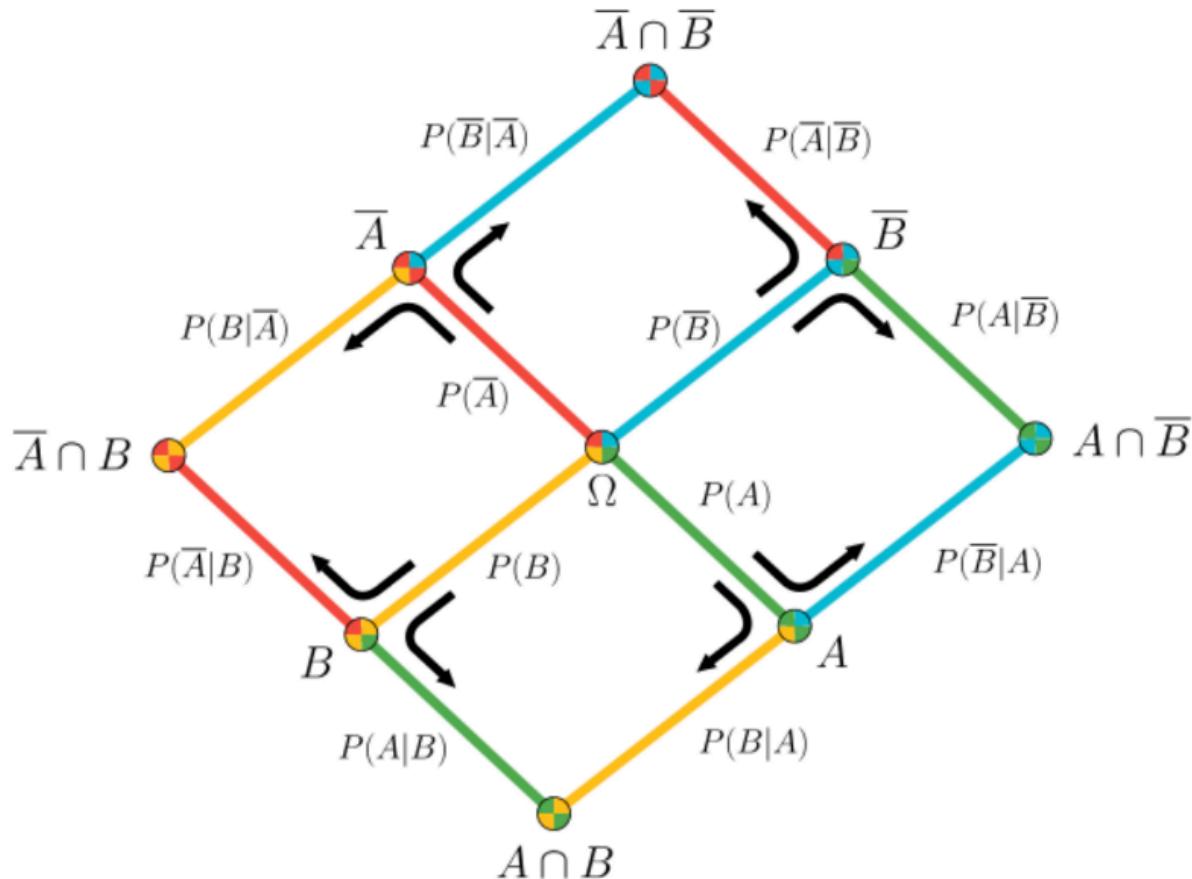
- Formula

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

- Tree Structure



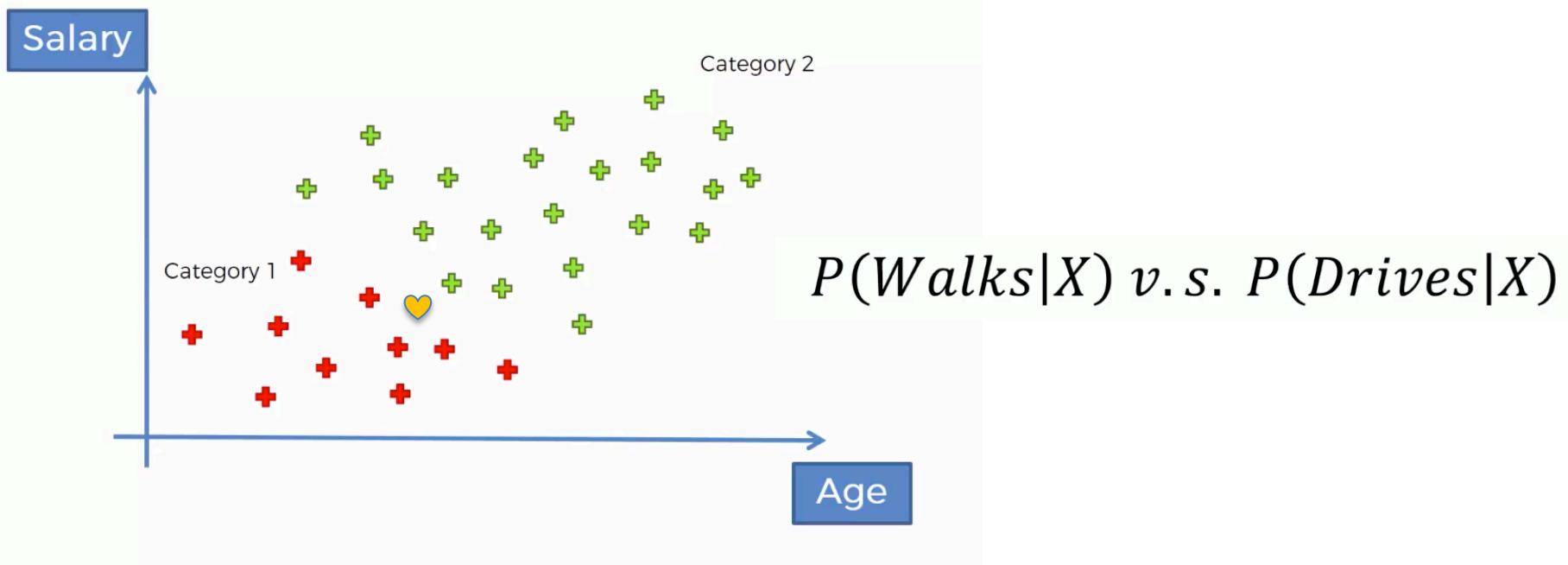
Another Visualizations of Bayes' theorem



Example on naïve bayes algorithm



- What happens if we add a new observation, how do we classify it? 「walk? Drive?」



Let's calculate $P(\text{Walks}|X)$ first.



Step 1

$$P(\text{Walks}|X) = \frac{P(X|\text{Walks}) * P(\text{Walks})}{P(X)}$$

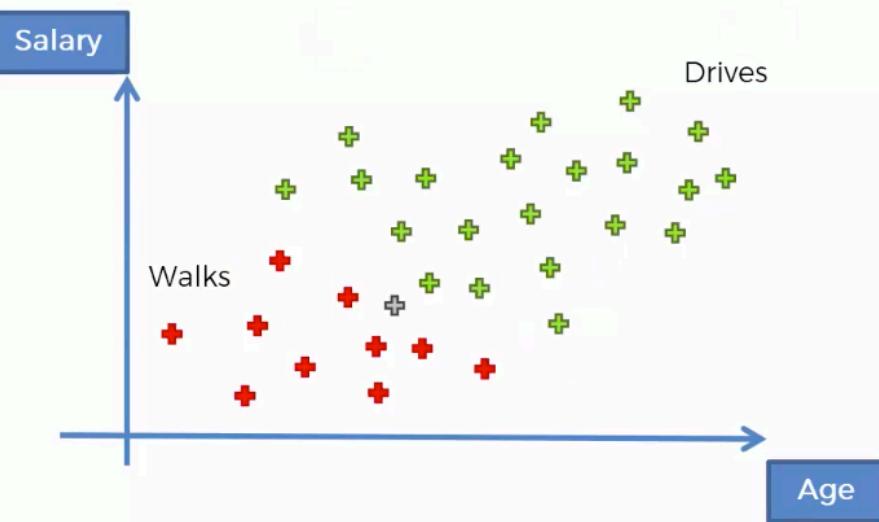
The diagram illustrates the components of the Bayes' theorem formula:

- #4 Posterior Probability
- #3 Likelihood
- #1 Prior Probability
- #2 Marginal Likelihood

Arrows point from each numbered label to its corresponding term in the formula:

- #4 points to $P(\text{Walks}|X)$
- #3 points to $P(X|\text{Walks})$
- #1 points to $P(\text{Walks})$
- #2 points to $P(X)$

P(Walks|X)



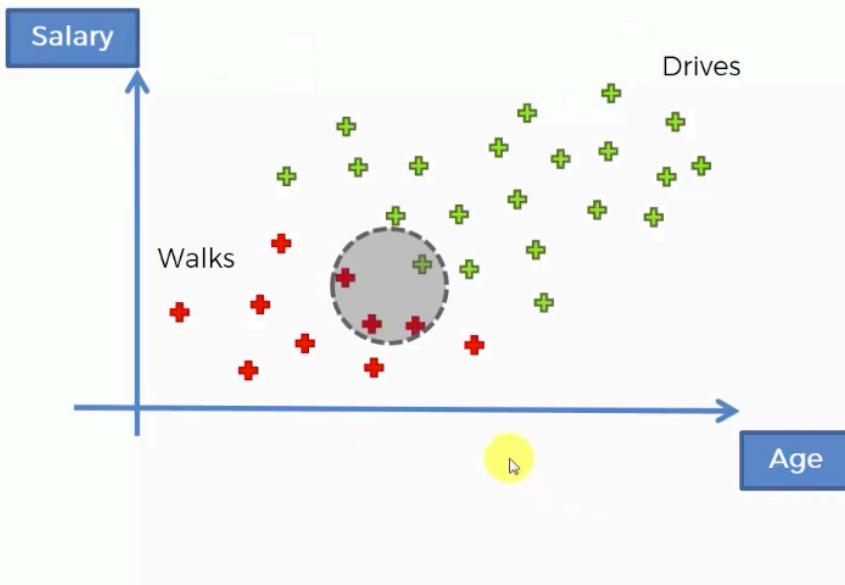
#1. P(Walks)

$$P(\text{Walks}) = \frac{\text{Number of Walkers}}{\text{Total Observations}}$$

$$P(\text{Walks}) = \frac{10}{30}$$



P(Walks|X)



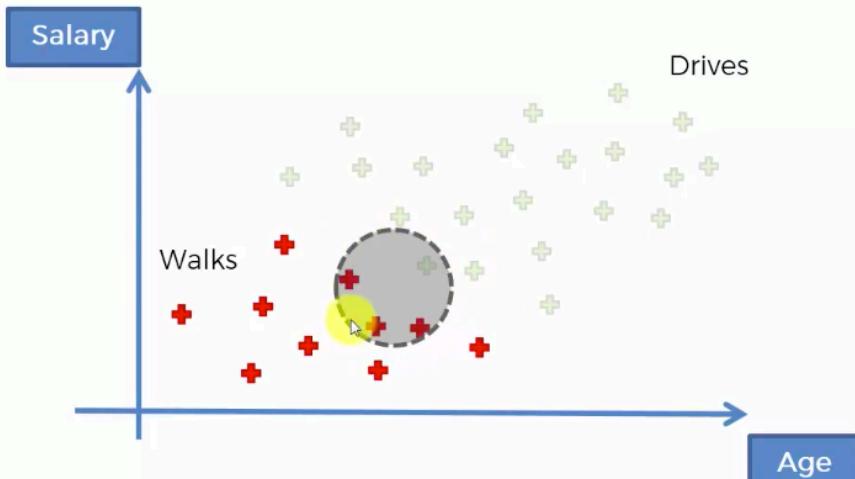
#2. P(X)

$$P(X) = \frac{\text{Number of Similar Observations}}{\text{Total Observations}}$$

$$P(X) = \frac{4}{30}$$



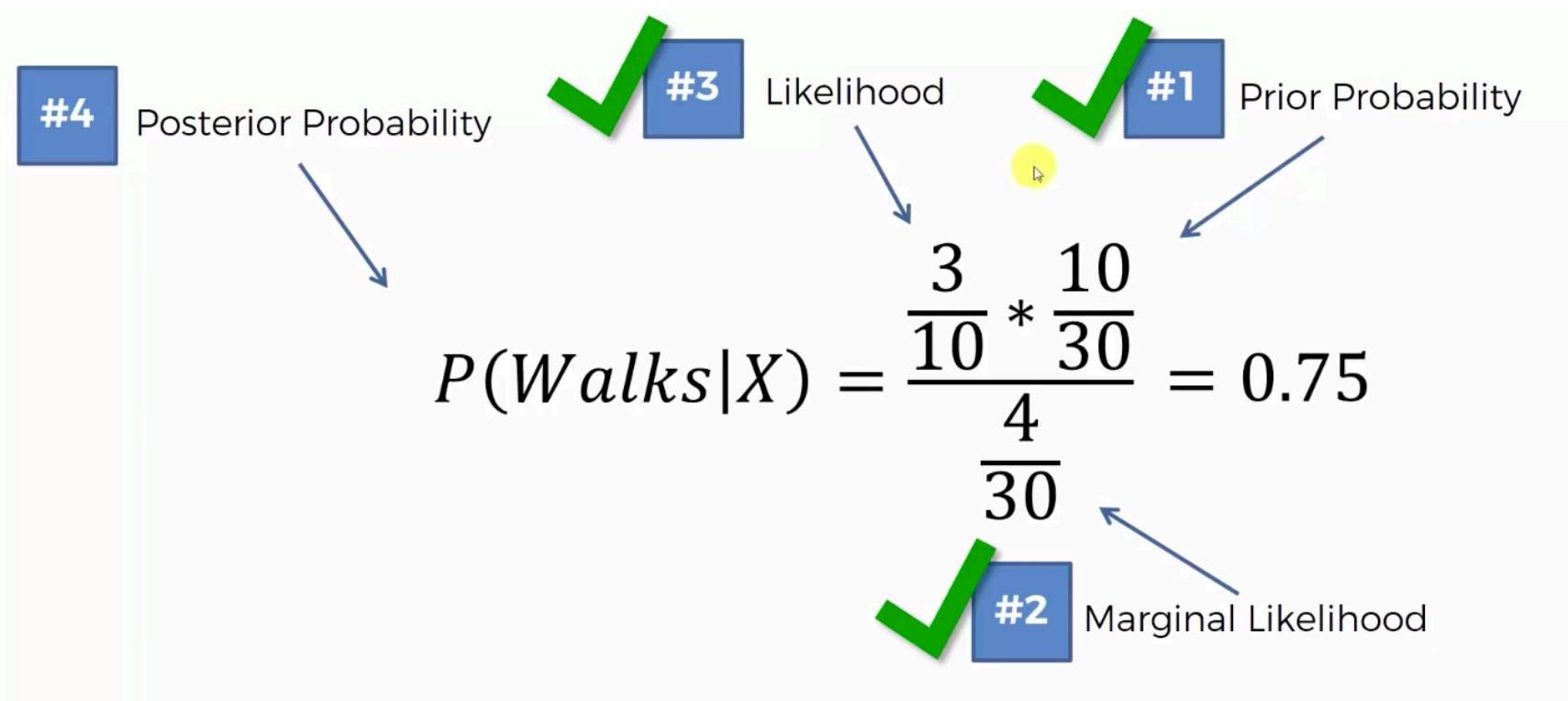
P(Walks|X)



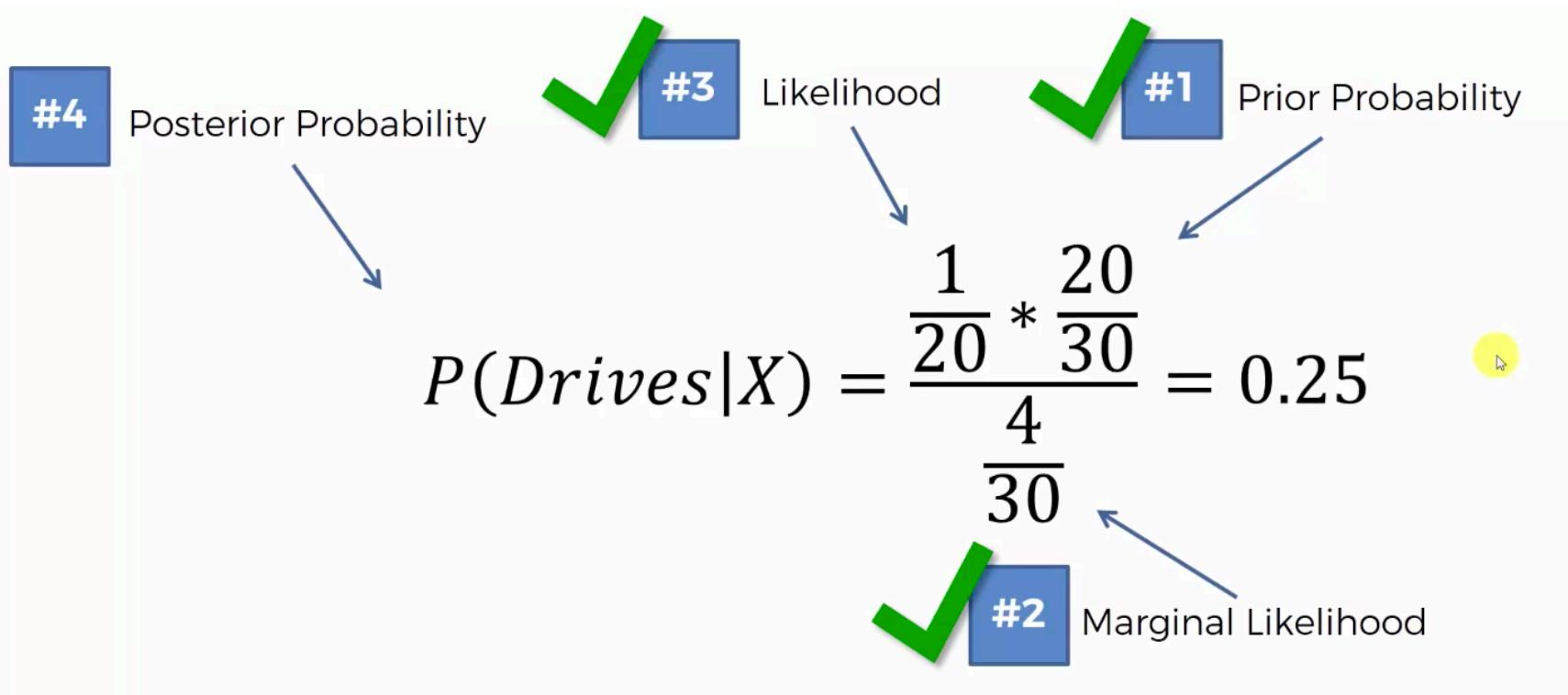
#3. P(X|Walks)

*Number of Similar Observations
Among those who Walk*
$$P(X|Walks) = \frac{\text{Number of Similar Observations}}{\text{Total number of Walkers}}$$

$P(\text{Walks}|X)$



$P(\text{Drives}|X)$



To simplify...



$P(Walks|X)$ v.s. $P(Drives|X)$

$$\frac{P(X|Walks) * P(Walks)}{\cancel{P(X)}} \text{ v.s. } \frac{P(X|Drives) * P(Drives)}{\cancel{P(X)}}$$

When the # of classes are >2 ...

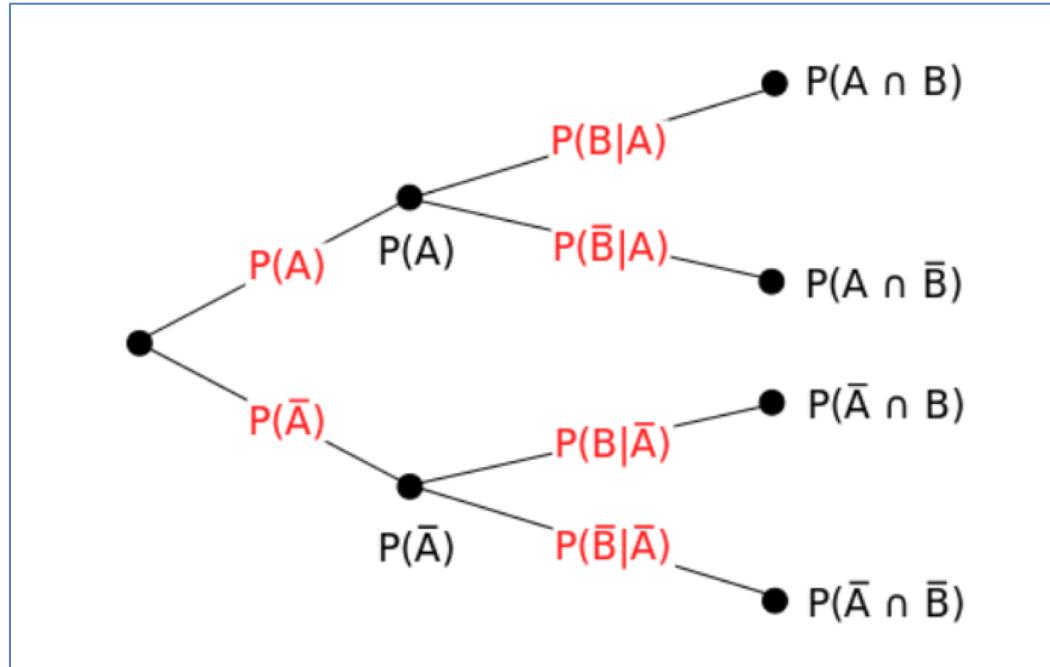


- $P(A_1|X) = P(x|A_1) * P(A_1) / ?$
- $P(A_1|X) = P(x|A_1) * P(A_1) / P(X)$
- $P(A_1|X) = P(x|A_1) * P(A_1)$
- $P(A_2|X) = P(x|A_2) * P(A_2)$
- ...

Why is it called Naïve Bayes?



- It holds the assumption that *the attributes are conditionally independent given the class*.



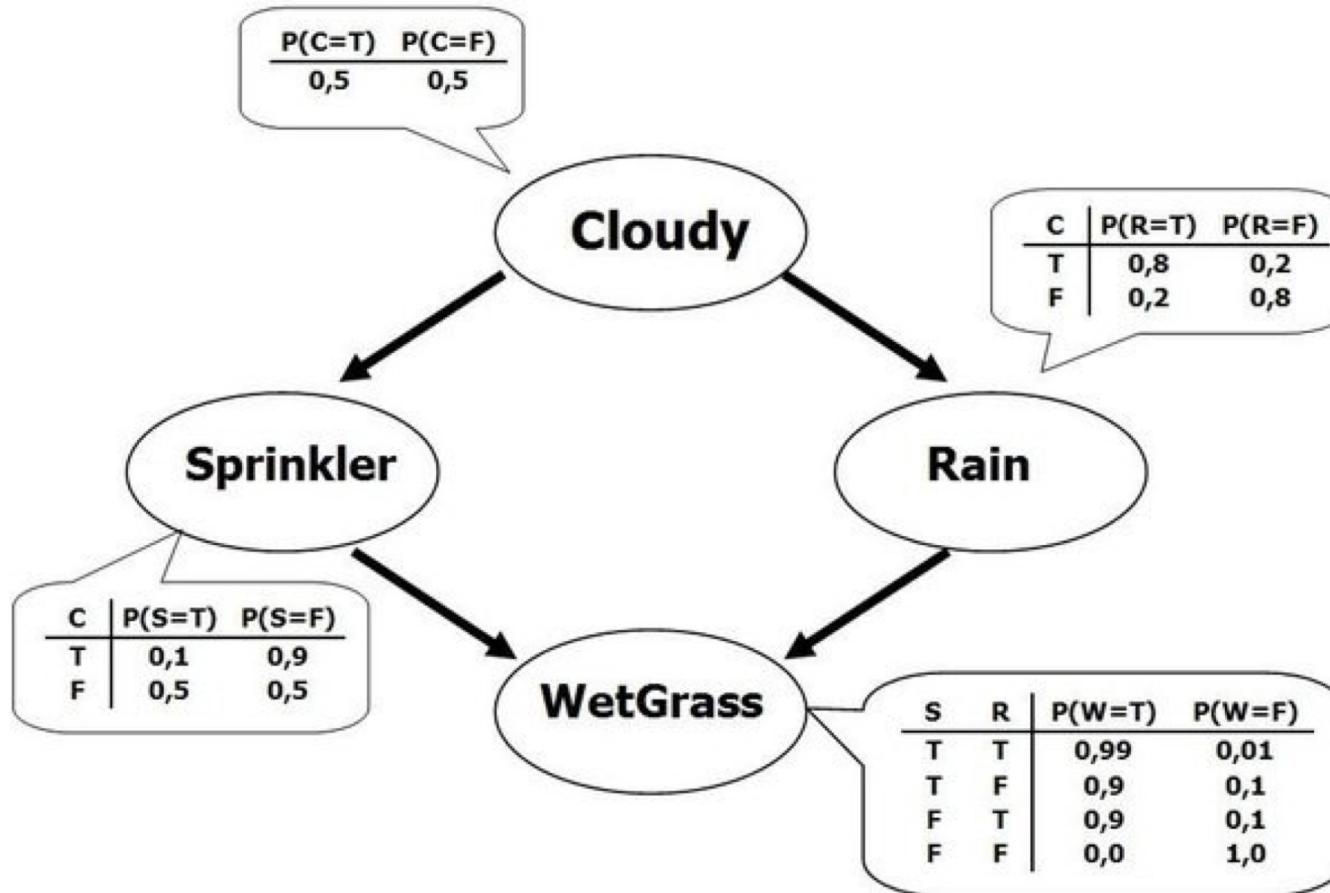
NB to BN...



- Bayesian network is designed to eliminate the naïve assumption.

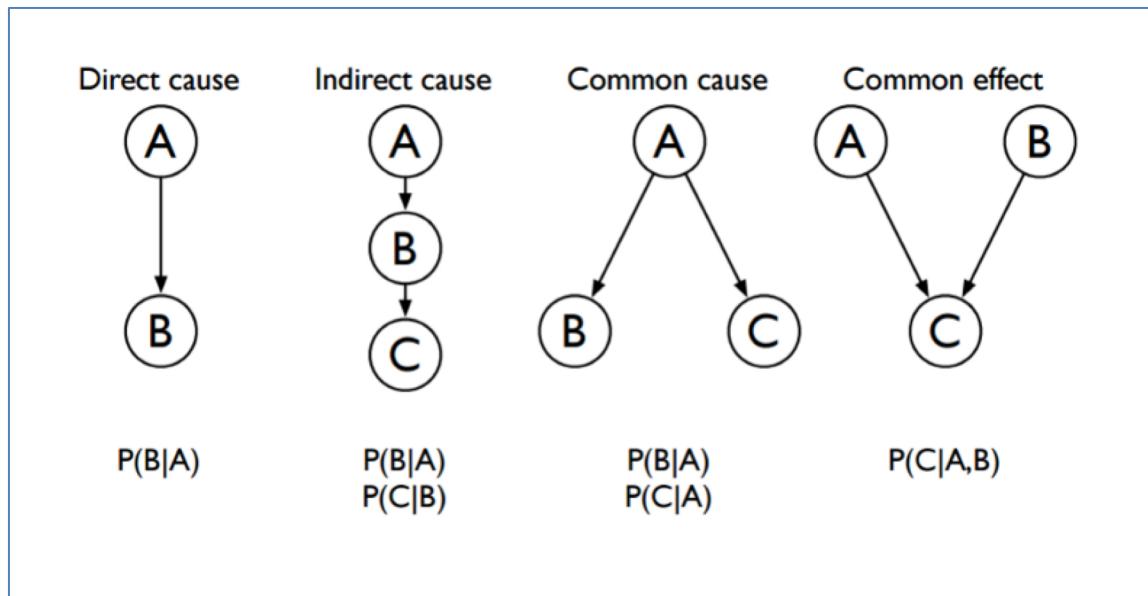


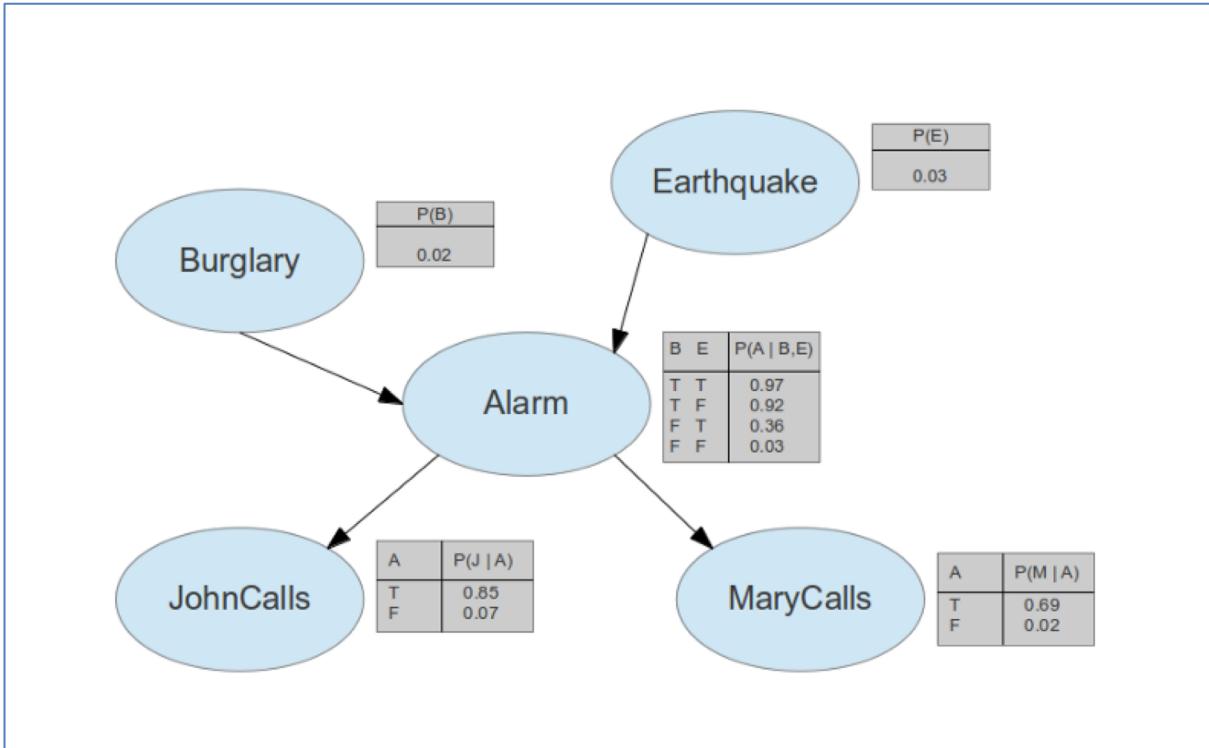
Another example





- Conditional dependency relations (arcs) from *node A* to another *node B* represent that *node B* is a child of *node A* or put it differently the *node A* is a parent of *node B*.





$$p(X_1, \dots, X_n) = \prod_{i=1}^n p(X_i | Parents(X_i))$$



- A Bayesian network is a **directed acyclic graph** in which each edge corresponds to a conditional dependency, and each node corresponds to a unique random variable.
- DAG is used in Markov chain, Deep learning, Structural prediction, etc.

OCR case



- How to recognize OCR letters???
- How do you know it is an A?
- What info do you need to recognize it is an A?



A A A A A A A
B B B B B B B
C C C C C C C
E E E E E E E
F F F F F F F
K K K K K K K
S S S S S S S
X X X X X X X

Figure 1. Examples of the character images generated by "warping" parameters.



- Let's read a paper and see how they deal with the problem.

Letter recognition using Holland-style adaptive classifiers. PW Frey, DJ Slate - Machine Learning, 1991 – Springer.

Reading time



Letter recognition using Holland-style adaptive classifiers.
PW Frey, DJ Slate - Machine Learning, 1991 – Springer.

Read paper: try to identify

- characterization of images (attributes)
- methods used (classifiers)
- performance measures.



- Guess: what might be the most difficult letter to recognize?
- What might be the most difficult pair?
- Check your results with weka.

Let's play with the data.



Let's go to UCI website.

- Intro to UCI archive website.
- <https://archive.ics.uci.edu/ml/datasets/Letter+Recognition>
- Also abel.lis.illinois.edu/data for downloading arff file.

Reproducibility



- Why is it hard to reproduce other research work?
- Replicate vs reproduce.
- Why is it important to make reproducible research outcome?
- What are the good habits to make good reproducibility?

Structural Prediction



- why is structure/are correlations useful?

v z
|
q

why is structure/are correlations useful?



Q **V** **I** **Z**

Q V I Z

Correlations not taken into account

Q **V** **I** **Z**

Q U I Z

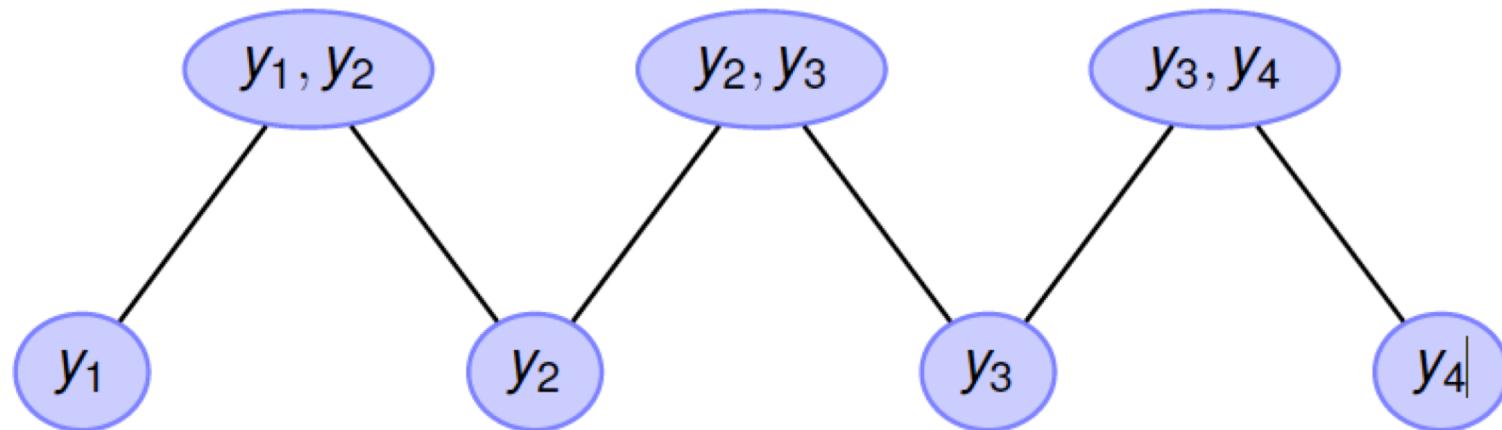
Correlations taken into account



- Decomposition and visualization

Q	V	I	Z
Q	U	I	Z

$$\begin{aligned}F(\mathbf{w}, x, y_1, \dots, y_4) = & f_1(\mathbf{w}, x, y_1) + f_2(\mathbf{w}, x, y_2) + f_3(\mathbf{w}, x, y_3) + f_4(\mathbf{w}, x, y_4) \\& + f_{1,2}(\mathbf{w}, x, y_1, y_2) + f_{2,3}(\mathbf{w}, x, y_2, y_3) + f_{3,4}(\mathbf{w}, x, y_3, y_4)\end{aligned}$$



Structural Prediction



- Why is this helpful?
- Formulate it as prediction of all four letter words (multiclass prediction):

$$y \in \mathcal{Y} = \{1, \dots, 26^4\}$$

- Problem: Really large output space A.
- Using correlations (eg: dictionary), the output space is much smaller.

Why and when do we need SP?



- Another example: denoising.



Predictions from neighboring pixel are useful.



- SP is very powerful in image processing.
- You might have heard of CNN, RNN, deep learning, etc.
- SP is their very good foundation.
- Limitation of SP.

Let's go through the assignment



- Another OCR problem.

Mid-term feedback



- Complete the feedback in moodle for the course.
- Q & A time



-
- That's all for today.