

Lecture 2 – Probability Basics for Data Analytics

Dr. Jing Li

Department of Computing

The Hong Kong Polytechnic University

Jan 18&20, 2022

Looking Back (Last Lecture)

- Course arrangement plan and structure
- Introduction to Data Analytics (An Introduction)
 - Key concepts, *data*, *data analytics*, etc.
 - Applications, e.g., *search engine*, *product recommendation*, *spam detection*, etc.
 - Advance Technology, e.g., *classification*, *regression*, *clustering*, *matching*, *time-series analysis*, etc.
- Some basic concepts of *data mining* and *big data*.
 - Good to have an impression about what they are.
 - You'll have chances to learn more about them later!

Probability VS. Data Analytics

- In analytics process, we usually use random variables to describe the data. **Why?**
 - Mathematical process to solve real-life problem
 - To make data computable.
 - To discover the patterns and trends behind data.
 - Allows the analysis of distribution and statistics.
- Probability helps predict the how likely that an event will happen.
 - E.g., *weather prediction, product recommendation.*

What is probability?

- *Probability* is a numerical description of *how likely an event is to occur* and or *how likely that a proposition is true*. --- From [Wikipedia](#).
- We run a random experiment n times, during which an event A occurs m times, then we say the *frequency* of A 's occurrence is $f_A = \frac{m}{n}$.
- When n is large enough, f_A will be very close to a value p , which is defined as the probability of A to occur, i.e., $\lim_{n \rightarrow +\infty} f_A \equiv P(A) = p$
 - When we toss a coin, the probability of “heads up” is 0.5

Roadmap

- Basic Concepts of Probability
- Conditional Probability and Bayes' Formula
- Data Analytics Application: Naïve Bayes
 - Text Classification
 - Naïve Bayes Formulation

Roadmap

- Basic Concepts of Probability
- Conditional Probability and Bayes' Formula
- Data Analytics Application: Naïve Bayes
 - Text Classification
 - Naïve Bayes Formulation

What is *sample space*?

- The set of all possible outcomes of an experiment
- Example: *coin flipping*
- What are the possible outcomes?
 - The coin lands *Heads up* - H
 - The coin lands *Tails up* - T
- The sample space: $S = \{H, T\}$
- For fair coins
 - What are the chances of H and T happening?
 - Empirically, H and T have 50-50 chance to happen.
 - The *probability* for H to happen is 50% (0.5), so does T



What is an *event*?



- Subsets of the sample space
- Example: *rolling a die once*
 - The sample space $S = \{1,2,3,4,5,6\}$
 - An example event $E = \{1,3,4\}$
 - If we *roll the die once*, and the it lands with 1 or 3 or 4 up, then we say the event E *occurs*.
- The probability that E *occurs* is:
 - $P(E) = \frac{1+1+1}{6} = \frac{3}{6} = \frac{1}{2}$

What is *probability*?

- A *probability function* P that assigns a real number (the probability of E) to every event $E \subset S$
- P *must* satisfy the following basic properties:
 - $0 \leq P(E) \leq 1$
 - $P(S) = 1$
- **An important property to speed up computing:**
 - For any *disjoint events*, E_i ($i = 1, 2, \dots, n$), we have
$$P(E_1 \cup E_2 \cup \dots \cup E_n) = P(E_1) + P(E_2) + \dots + P(E_n)$$

Some Notions

- $P(E \cup F)$: the probability that E or F occurs
- $P(E, F)$ or $P(EF)$: the probability that both E and F occurs.
- $P(E^c)$: the probability that E does not occur.

Revisit the Coin Flipping Example

- **Scenario 1:** if the coin is flipped twice, what is the probability of two *heads*?
 - $P(H, H) = \frac{1}{4}$
- **Scenario 2:** if the coin is flipped twice, what is the probability of two *heads*, given that *we know the first toss gave a head*.
 - $P(H, H | H \text{ at the first toss}) = \frac{1}{2}$

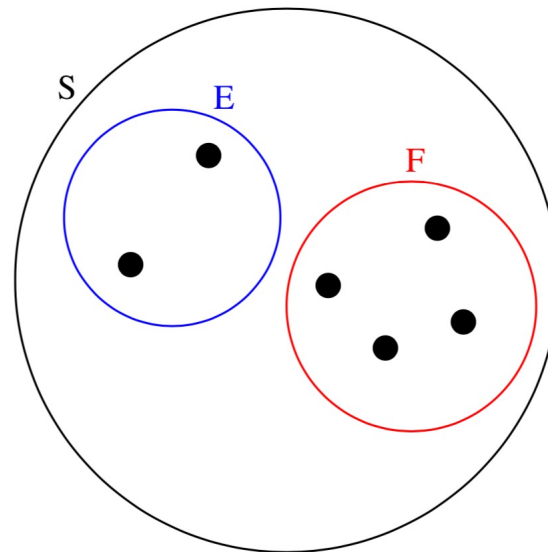
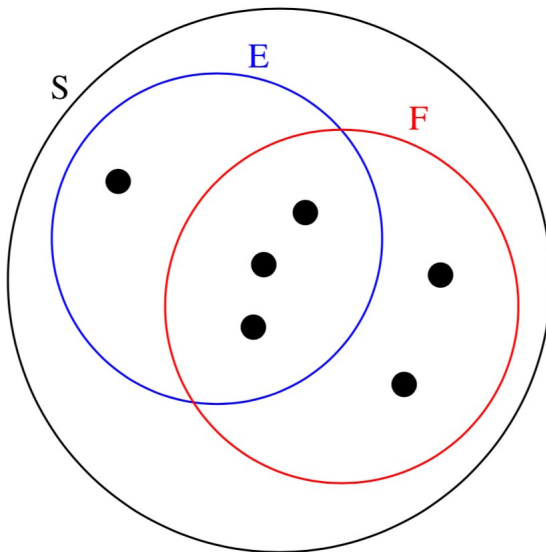


Roadmap

- Basic Concepts of Probability
- Conditional Probability and Bayes' Formula
- Data Analytics Application: Naïve Bayes
 - Text Classification
 - Naïve Bayes Formulation

Conditional Probability

- If E and F are events, then $P(E|F)$ is the *conditional probability* of E , given F .
- $P(E|F) \equiv \frac{P(E,F)}{P(F)}$, assuming that $P(F) \neq 0$



What is
 $P(E|F)$ for the
two cases?

Example: Conditional Probability

- Suppose we draw a card from a shuffled set of 52 playing cards.
- What is the probability of drawing a *Queen*, given that the card drawn is of suits *Hearts*.

- $$P(Q|H) = \frac{P(Q,H)}{P(H)} = \frac{1/52}{1/4} = \frac{1}{13}$$

- What is the probability of drawing a *Queen*, given that the card drawn is a *Face* card?

- $$P(Q|F) = \frac{P(Q,F)}{P(F)} = \frac{P(Q)}{P(F)} = \frac{4/52}{12/52} = \frac{1}{3}$$



Discussion: Teddy and Charlie

- It is well known that *Uncle Bob* has two beautiful children, *Teddy* and *Charlie*. However, no one knows whether they are *sons* or *daughters*.
- One day, you met *Bob* in the park. He told you that he has at least a *daughter*. Can you estimate the probability both *Teddy* and *Charlie* are daughters?
- On the other day, you met *Uncle Bob* again. He was with his beautiful daughter and introduced that she was *Teddy*. Now, can you estimate again the probability both *Teddy* and *Charlie* are daughters?



Law of Total Probability

- Sometimes, the computation of $P(E)$ will be easier if we condition E on another event F , namely, from
- $P(E) = P(E(F \cup F^c)) = P(E, F) + P(E, F^c)$
- Also, $P(E, F) = P(E|F)P(F)$
and $P(E, F^c) = P(E|F^c)P(F^c)$
- $P(E) = P(E|F)P(F) + P(E|F^c)P(F^c)$



IMPORTANT

Example: Law of Total Probability

- An insurance company holds the following data concerning the probability of an *insurance claim*:
 - For people under age 30, the probability is 4%
 - For people over age 30, the probability is 2%
- And it is known that 30% of the targeted population is under age 30.
- What is the probability of an insurance claim for a randomly chosen person?

Example: Law of Total Probability

- $S = \{\text{all persons under consideration}\}$
- $C = \{\text{persons filing a claim}\}$
- $U = \{\text{persons under age 30}\}$
- Thus, $P(C) = P(C|U) P(U) + P(C|U^c) P(U^c)$

$$\begin{aligned} &= \frac{4}{100} \frac{3}{10} + \frac{2}{100} \frac{7}{10} \\ &= \frac{26}{1000} = 2.6\% . \end{aligned}$$

Bayes' formula

- Sometimes, we need a formula that *inverts conditioning*, e.g., predicts an event conditioned on some observations.

- Since $P(EF) = P(E|F)P(F)$
and $P(EF) = P(F|E)P(E)$



IMPORTANT

- Then we have

$$P(F|E) = \frac{P(E, F)}{P(E)} = \frac{P(E|F)P(F)}{P(E|F)P(F) + P(E|F^c)P(F^c)}$$

Law of total probability

Example: Bayes' formula

- Suppose 1 in 1,000 persons has a certain disease.
- For 99% of the diseased persons, a test will yield positive results.
- For 5% of the healthy persons, a test will also yield positive results (false alarm).
- What is the probability of a positive test diagnose the disease?
 - $D = \{\textit{Diseased persons}\}$
 - $H = \{\textit{Healthy persons}\}$
 - $+$ = $\{\textit{Persons with positive test results}\}$

Example: Bayes' formula

- What is the probability of a positive test diagnose the disease?
 - $D = \{\textit{Diseased persons}\}$
 - $H = \{\textit{Healthy persons}\}$
 - $+$ = $\{\textit{Persons with positive test results}\}$
- By the given statistics, we have
 - $P(D) = 0.001, P(+|D) = 0.99, P(+|H) = 0.05$
- With Bayes' formula:

$$\begin{aligned} P(D|+) &= \frac{P(+|D) \cdot P(D)}{P(+|D) \cdot P(D) + P(+|H) \cdot P(H)} \\ &= \frac{0.99 \cdot 0.001}{0.99 \cdot 0.001 + 0.05 \cdot 0.999} \cong 0.0194 \end{aligned}$$

Independent Events

- Two events E and F are *independent* if
 - $P(E, F) = P(E)P(F)$
- In this case:
 - $P(E|F) = \frac{P(E, F)}{P(F)} = \frac{P(E)P(F)}{P(F)} = P(E)$, assuming $P(F) \neq 0$
- In other words,
 - *knowing F occurred doesn't change the probability of E .*

Example: Independent Events

- Two numbers are drawn at random from $\{1,2,3,4\}$
- If order is unimportant, then what is the sample space S ?
- Define the following functions on S :
 - $X(\{i,j\}) = i + j$
 - $Y(\{i,j\}) = |i - j|$
- Which of the following pairs are independent?
 - $X = 5$ and $Y = 2$
 - $X = 5$ and $Y = 1$

Roadmap

- Basic Concepts of Probability
- Conditional Probability and Bayes' Formula
- Data Analytics Application: Naïve Bayes
 - Text Classification
 - Naïve Bayes Formulation

Roadmap

- Basic Concepts of Probability
- Conditional Probability and Bayes' Formula
- Data Analytics Application: Naïve Bayes
 - Text Classification
 - Naïve Bayes Formulation

Is this spam?

Subject: Important notice!

From: Stanford University <newsforum@stanford.edu>

Date: October 28, 2011 12:34:16 PM PDT

To: undisclosed-recipients::;

Greats News!

You can now access the latest news by using the link below to login to Stanford University News Forum.

<http://www.123contactform.com/contact-form-StanfordNew1-236335.html>

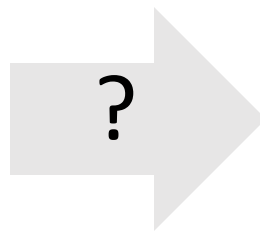
Click on the above link to login for more information about this new exciting forum. You can also copy the above link to your browser bar and login for more information about the new services.

© Stanford University. All Rights Reserved.

What is the subject of this medical article?

MeSH Subject Category Hierarchy

MEDLINE Article



- Antagonists and Inhibitors
- Blood Supply
- Chemistry
- Drug Therapy
- Embryology
- Epidemiology
- ...

Positive or negative movie review?

- + *...zany characters and richly applied satire, and some great plot twists*
- *It was pathetic. The worst part about it was the boxing scenes...*
- + *...awesome caramel sauce and sweet toasty almonds. I love this place!*
- *...awful pizza and ridiculously overpriced...*

Positive or negative movie review?

- + ...*zany* characters and **richly** applied satire, and some **great** plot twists
- It was **pathetic**. The **worst** part about it was the boxing scenes...
- + ...**awesome** caramel sauce and sweet toasty almonds. I **love** this place!
- ...**awful** pizza and **ridiculously** overpriced...

Why sentiment analysis?

- *Movie*: is this review positive or negative?
- *Products*: what do people think about the new iPhone?
- *Public sentiment*: how is consumer confidence?
- *Politics*: what do people think about this candidate or issue?
- *Prediction*: predict election outcomes or market trends from sentiment

Basic Sentiment Classification

- Sentiment analysis is the detection of **attitudes**
- Simple task we focus on in this chapter
 - Is the attitude of this text positive or negative?
- We return to affect classification in later chapters

Summary: Text Classification

- Sentiment analysis
- Spam detection
- Authorship identification
- Language Identification
- Assigning subject categories, topics, or genres
- ...

Roadmap

- Basic Concepts of Probability
- Conditional Probability and Bayes' Formula
- Data Analytics Application: Naïve Bayes
 - Text Classification
 - Naïve Bayes Formulation

Text Classification: definition

- *Input*:
 - a document d
 - a fixed set of classes $C = \{c_1, c_2, \dots, c_J\}$
- *Output*: a predicted class $c \in C$

Classification Methods:

Hand-coded rules

- Rules based on combinations of words or other features
 - spam: black-list-address OR (“dollars” AND “you have been selected”)
- Accuracy can be high
 - If rules carefully refined by expert
- But building and maintaining these rules is expensive

Classification Methods:

Supervised Machine Learning

- *Input:*
 - a document d
 - a fixed set of classes $C = \{c_1, c_2, \dots, c_J\}$
 - A training set of m hand-labeled documents $(d_1, c_1), \dots, (d_m, c_m)$
- *Output:*
 - a learned classifier $\gamma: d \rightarrow c$

Classification Methods: Supervised Machine Learning

- Any kind of classifier
 - Naïve Bayes
 - Logistic regression
 - Neural networks
 - k-Nearest Neighbors
 - ...

Naive Bayes Intuition

- Simple (“naive”) classification method based on Bayes rule (Bayes’ formula)
- Usually written as Naïve Bayes (in French).
- Relies on very simple representation of document
 - **Bag of words**

The Bag of Words Representation

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

The bag of words representation

$Y($

seen	2
sweet	1
whimsical	1
recommend	1
happy	1
...	...

$) = C$




Bayes' Rule Applied to Documents and Classes

- For a document d and a class c

$$P(c | d) = \frac{P(d | c)P(c)}{P(d)}$$

Our goal is to maximize $P(c | d)$ with a $c \in \mathcal{C}$

Naive Bayes Classifier (I)

argmax: the arguments of the maximum (e.g., find the optimal value of c which maximizes $P(c|d)$)

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(c | d)$$

$$= \operatorname{argmax}_{c \in C} \frac{P(d | c)P(c)}{P(d)}$$

$$= \operatorname{argmax}_{c \in C} P(d | c)P(c)$$

MAP is “maximum a posteriori” = most likely class

Bayes' Formula

Dropping the denominator

Naive Bayes Classifier (II)

"Likelihood"

"Prior"

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(d | c)P(c)$$

$$= \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c)P(c)$$

Document d represented as
features x_1, x_2, \dots, x_n

Naive Bayes Classifier (III)

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c) P(c)$$

$O(|X|^n \cdot |C|)$ parameters

How often does this class occur?

Could only be estimated if a very, very large number of training examples was available.

We can just count the relative frequencies in a corpus

Multinomial Naive Bayes

Independence Assumptions

$$P(x_1, x_2, \dots, x_n | c)$$

- **Bag of Words assumption:** Assume position doesn't matter
- **Conditional Independence:** Assume the feature probabilities $P(x_i | c_j)$ are independent given the class c .

$$P(x_1, \dots, x_n | c) = P(x_1 | c) \cdot P(x_2 | c) \cdot P(x_3 | c) \cdot \dots \cdot P(x_n | c)$$

Multinomial Naive Bayes Classifier

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c) P(c)$$

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c_j) \prod_{x \in X} P(x | c)$$

$$X = \{x_1, x_2, \dots, x_n\}$$

Applying Multinomial Naive Bayes Classifiers to Text Classification

positions \leftarrow all word positions in the test document

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i \in \text{positions}} P(x_i | c_j)$$

Problems with multiplying lots of probabilities

- There's a problem with this:

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i \in \text{positions}} P(x_i | c_j)$$

- Multiplying lots of probabilities can result in floating-point underflow!
- Luckily, $\log(ab) = \log(a) + \log(b)$
- Let's sum logs of probabilities instead of multiplying probabilities!

Put them in in the log space

Instead of this:

This:

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i \in \text{positions}} P(x_i | c_j)$$
$$c_{NB} = \operatorname{argmax}_{c_j \in C} \left[\log P(c_j) + \sum_{i \in \text{positions}} \log P(x_i | c_j) \right]$$

This is ok since log doesn't change the ranking of the classes (class with highest prob still has highest log prob)

Model is now just max of sum of weights: a *linear* function of the inputs

So naive bayes is a *linear classifier*

One Slide to Takeaway

- What is a **sample space, event, probability**?
- What is **conditional probability**?
- What is **Bayes' formula**?
- What are **independent events**?
- What is **text classification**?
- How to formulate a **Naïve Bayes classifier**?