



Bayesian Classifiers

(Put together from many sources)

Sharma Chakravarthy

Department of Computer Science and Engineering
University of Texas at Arlington
Fall 2019

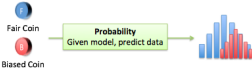
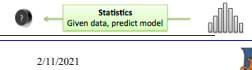





2/11/2021 Sharma Chakravarthy 1

Probability and Statistics

- We need a little bit of **probability** and **statistics** before going into Bayes approach to classification
- So, what is the difference between probability and statistics?
 - **Probability** is starting with an animal, and figuring out what footprints it will make.
 - **Statistics** is seeing a footprint, and guessing the animal.

Probability & Statistics



2/11/2021 © Sharma_Chakravarthy

Basics

- Consider rolling two dice (a blue and yellow one) and noting down the numbers on each roll (alternatively, drawing from a deck of cards)
- Let X and Y denote the number of dots we get on the blue and yellow dice, respectively, and consider the meaning of $P(X+Y=6)$ in the large!
- Sample space is the list of Possible outcomes as shown
- We place a weight of $\frac{1}{36}$ on each Outcome reflecting they are **equally likely (important!)**
- Outcomes (1,5), (2,4),..., (5,1) have a total weight of $\frac{5}{36}$

1,1	1,2	1,3	1,4	1,5	1,6
2,1	2,2	2,3	2,4	2,5	2,6
3,1	3,2	3,3	3,4	3,5	3,6
4,1	4,2	4,3	4,4	4,5	4,6
5,1	5,2	5,3	5,4	5,5	5,6
6,1	6,2	6,3	6,4	6,5	6,6

Sample Space for the Dice Example






2/11/2021 © Sharma_Chakravarthy

Intuitively

- **Sample space** may be difficult to comprehend for some problems
- If you are doing an experiment recording the outcome in a note book, it may look like →
- After many, many repetitions, Approximately $\frac{5}{36}$ of the lines Will have a 'yes' value!
- If you did the experiment 1008 times, approximately $X+Y=6$ will occur $\frac{5}{36} \times 1008 = 140$ times
- This is what probability really is: In **what fraction** of the lines does the **event of interest** happen? **It sounds simple, but if you think about this "as counting the outcome of an experiment", probability problems are a lot easier to understand.** And it is the fundamental basis of computer simulation.

notebook line	outcome	blue+yellow = 6?
1	blue 2, yellow 6	No
2	blue 3, yellow 1	No
3	blue 1, yellow 1	No
4	blue 6, yellow 2	Yes
5	blue 1, yellow 1	No
6	blue 3, yellow 4	No
7	blue 3, yellow 1	Yes
8	blue 3, yellow 6	No
9	blue 2, yellow 5	No

2/11/2021 © Sharma_Chakravarthy

Basics

- We assume experiments are **repeatable!**
- We assume that experiments are performed a **large number of times** (**in the large is important**)
- An **event** A has a Boolean outcome (yes/no) of the experiment. Examples:
 - $X+Y = 6$
 - $X = 1$
 - $Y = 3$
 - $X - Y = 4$
- A **random variable** (X or Y) is a numerical outcome of the experiment: $X+Y$, $2*X*Y$,



2/11/2021



© Sharma,Chakravarthy

Definition of $P(A \wedge B)$ Boolean AND

- For any event of interest A , imagine a column on A in the notebook. The k^{th} line (or row) in the notebook, $k = 1, 2, 3, \dots$, will say Yes or No, depending on whether A occurred or not during the k^{th} repetition of the experiment. For instance, we can have such a column in our table above, for the event {blue + yellow = 6}.
- For any event of interest A , we define $P(A)$ to be the **long-run fraction of lines with Yes entries**.
- For any events $A \wedge B$ (Boolean AND), imagine a new column in our notebook, labeled " $A \wedge B$." In each line, this column will say Yes **if and only if there are Yes entries in columns A and B** .
- $P(A \wedge B)$ is then defined to be the **long-run fraction of lines with Yes entries in the new column labeled " $A \wedge B$."**



2/11/2021



© Sharma,Chakravarthy

Definition of $P(A \vee B)$

- For any events A, B imagine a new column in our notebook, labeled " $A \vee B$." In each line, this column will say Yes if and only if at least one of the entries for A or B says Yes.
- $P(A \vee B)$ is then defined to be the long-run fraction of lines with Yes entries in the new column labeled " $A \vee B$."



2/11/2021



© Sharma,Chakravarthy

Definition

- For any events A, B imagine a new column in our notebook, labeled " $A | B$ " and pronounced " **A given B** ." In each line:
 - This new column will say "NA" ("not applicable") if the B entry is No.
 - If it is a line in which the B column says Yes, then this **new column** will say Yes or No, depending on whether the A column says Yes or No.
- Then $P(A | B)$ means the **long-run fraction of lines in the notebook in which the $A | B$ column says Yes—among the lines which do NOT say NA.**
 - Conditional probability



2/11/2021



© Sharma,Chakravarthy

Common mistake

- A common mistake is to confuse $P(A \wedge B)$ and $P(A|B)$. This is important to understand. Compare the values of $P(X = 1 \wedge S = 6)$ and $P(X = 1 | S = 6)$, where $S = X + Y$

- $P(X = 1 \wedge S = 6)$ is $1/36$ and comes from (1,5) out of possible 36 outcomes (35 no and 1 yes)

- $P(X = 1 | S = 6)$ is $1/5$ comes from (1,5) out of possible 5 relevant outcomes.

outcome	S=6	X=1?
(1,5)	yes	yes
(2,4)	yes	no
(3,3)	yes	no
(4,2)	yes	no
(5,1)	yes	no

Please understand this clearly!

All other cases NA hence, not used!

HW: What is $P(X=2 \vee S=6)$ where $S = X + Y$? What is $P(X=2 \vee Y=4)$?



2/11/2021



© Sharma, Chakravarthy

More Definitions

- Suppose A, B are events such that it is impossible for them to occur at the same time. They are said to be **disjoint** or “mutually exclusive” events. “being freshman” and “being sophomore” are disjoint events

- If A and B are disjoint events, then

$P(A \vee B) = P(A) + P(B)$. Can be generalized to

$$P(A_1 \vee A_2, \dots, \vee A_k) = \sum_{i=1}^k P(A_i)$$

- if A and B are not disjoint, then

$$P(A \vee B) = P(A) + P(B) - P(A \wedge B)$$

- To compute $P(X=3 \text{ or } Y=4)$, first check whether they are disjoint!
 $1/6 + 1/6 - 1/36$ (we have counted (3, 4) twice! Hence the subtraction)
- No compact generalization for \forall !

1,1	1,2	1,3	1,4	1,5	1,6
2,1	2,2	2,3	2,4	2,5	2,6
3,1	3,2	3,3	3,4	3,5	3,6
4,1	4,2	4,3	4,4	4,5	4,6
5,1	5,2	5,3	5,4	5,5	5,6
6,1	6,2	6,3	6,4	6,5	6,6

Yes $6/36 + 1/36$

What about “female” or “freshman”?



2/11/2021



© Sharma, Chakravarthy

More Definitions

- Events A and B are independent (or unrelated) if

$$P(A \wedge B) = P(A) * P(B) = P(B \wedge A)$$

What about $P(X=3 \wedge X=7)$?
Events can be independent and disjoint if one of them is null or never Happens!
“promoted” \wedge “audited”?

$$P(A_1 \text{ and } A_2, \dots, \text{ and } A_k) = \prod_{i=1}^k P(A_i)$$

- How do we know whether two events are disjoint, independent?

- Must be inferred from the problem, domain, context!
- 2 cards are drawn from a deck **without replacement**! Independent or dependent? (HW)
- What about **with replacement**? (HW)

- For example, when two dices are rolled, outcome on each dice is independent of each other!

$P(X=4 \wedge Y=5) = 6/36 * 6/36 = 1/36$
where as events “good grade” \wedge “completing the project” are **not independent** events (related)
Contrast with $P(X=1 \wedge S=6)$!

1,1	1,2	1,3	1,4	1,5	1,6
2,1	2,2	2,3	2,4	2,5	2,6
3,1	3,2	3,3	3,4	3,5	3,6
4,1	4,2	4,3	4,4	4,5	4,6
5,1	5,2	5,3	5,4	5,5	5,6
6,1	6,2	6,3	6,4	6,5	6,6



2/11/2021



© Sharma, Chakravarthy

More Definitions

- If Events A and B are **not independent**, the previous expression generalizes to

$P(A \wedge B) = P(A) * P(B|A)$ instead of $P(A) * P(B)$. Suppose we use this for

$$P(X=4 \wedge Y=5) = P(X=4) * P(Y=5 | X=4) = 6/36 * 1/6 = 1/36 \text{ (same as earlier, why?)}$$

- Check $P(X=1 \wedge S=6)$
- Suppose 20% of all UTA students are in the COE, and 15% of engineering majors are female. (there are non-COE females! At UTA)
- $P(\text{COE student} \wedge \text{female}) = P(\text{COE student}) * P(\text{female} | \text{COE student})$
 $0.2 * 0.15 = 0.03$ or 3% of UTA students are “female engineers”
- When A and B are **actually independent**
 $P(B|A)$ is the same as $P(B)$ and reverts to $P(A \wedge B) = P(A) * P(B)$

- $P(B|A)$ is also known as conditional probability!
- Probability of B conditionally occurring on A
- $P(B|A) = P(B)$ if A and B are independent!



2/11/2021



© Sharma, Chakravarthy

Bayes theorem

- We can express $P(A \wedge B) = P(A) * P(B|A)$ as

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

Similarly,

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

giving rise to Bayes theorem

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

Bayes theorem provides a way of
Going from $P(X|Y)$ to $P(Y|X)$ or from
A sample labeled population to a
unseen data outcome prediction!
This is what classification is

Significance: if I know individual
probabilities of A and B, and $P(B|A)$,
I can compute the Probability of $P(A|B)$!



2/11/2021



© Sharma Chakravarthy

Example of Bayes Theorem

- Given:

- A doctor knows that meningitis causes stiff neck 50% of the time $P(S|M)$ is 0.5 **conditional probability**
- Prior probability of any patient having meningitis is 1/50,000 $P(M)$ **sample or population probability**
- Prior probability of any patient having stiff neck is 1/20 $P(S)$

- If a patient has stiff neck, what's the probability he/she has meningitis? **Posterior probability**
 $P(M|S)$

$$P(M|S) = \frac{P(S|M)P(M)}{P(S)} = \frac{0.5 \times 1/50000}{1/20} = 0.0002$$



2/11/2021



© Sharma Chakravarthy

Bayesian Classifiers

- How can we apply this to classification using a training data set
- Consider each attribute and class label as random variables
- Given a record with attributes (A_1, A_2, \dots, A_n)
 - Goal is to predict class C
 - Specifically, we want to find the value of C that maximizes $P(C|A_1, A_2, \dots, A_n)$
- Can we estimate $P(C|A_1, A_2, \dots, A_n)$ directly from data?

15



Bayesian Classifiers

- Approach:
 - Compute the posterior probability $P(C|A_1, A_2, \dots, A_n)$ for all values of C using the Bayes theorem
- $$P(C|A_1, A_2, \dots, A_n) = \frac{P(A_1, A_2, \dots, A_n|C)P(C)}{P(A_1, A_2, \dots, A_n)}$$
- Choose value of C that maximizes $P(C|A_1, A_2, \dots, A_n)$ **posterior probability!**
 - Equivalent to choosing value of C that maximizes numerator $P(A_1, A_2, \dots, A_n|C)P(C)$
- How to estimate $P(A_1, A_2, \dots, A_n|C)$? **How do we know $P(C)$?**

16



Naïve Bayes Classifier

$$P(C | A_1 A_2 \dots A_n) = \frac{P(A_1 A_2 \dots A_n | C) P(C)}{P(A_1 A_2 \dots A_n)}$$

- In the above expression, we do not know how to compute probability $P(A_1 A_2 \dots A_n | C)$ because of multiple A's!
 - This is where the assumption of **independence of A_i events** comes into use!
 - Also why it is called Naïve!
 - Assuming independence among attributes A_i when class is given:
 - $P(A_1, A_2, \dots, A_n) = P(A_1) * P(A_2) * \dots * P(A_n)$
 - $P(A_1, A_2, \dots, A_n | C)$ can be written as $P(A_1 | C) * P(A_2 | C) * \dots * P(A_n | C)$
- Can estimate $P(A_i | C)$ for all A_i and C_j .
- New point is classified to C_i if $P(C_i) \prod P(A_i | C_i)$ is maximal.

17



Simple one attribute example

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Overcast	Yes
Sunny	No
Overcast	Yes
Rainy	No

P(A|no) P(A|yes)

Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
Grand Total	5	9

Weather	No	Yes	P(A)
Overcast		4	=4/14 0.29
Rainy	3	2	=5/14 0.36
Sunny	2	3	=5/14 0.36
All	5	9	
	=5/14 0.36	=9/14 0.64	P(C)

Compute $P(\text{Yes}|\text{Sunny}) = P(\text{Sunny}|\text{yes}) * P(\text{yes})/P(\text{Sunny}) = 3/9 * 0.64/0.36 = 0.6$

Compute $P(\text{yes}|\text{overcast}) = P(\text{overcast}|\text{yes}) * P(\text{yes})/P(\text{overcast}) = 4/9 * 0.64/0.29 = 0.98$

$P(\text{No}|\text{Sunny}) = 0.4$

Compute $P(\text{No}|\text{overcast})$

- Step 1: convert data into frequency table (can be done for each attribute)
- Step 2: create probabilities (likelihood) table for the sample data
- Step 3: use Naïve Bayes equation to compute the posterior probability for each class
- Can be easily extended to multi-attribute and multi-class!
- For numerical attributes, normal distribution is assumed!



Simple one attribute example

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No

Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
Grand Total	5	9

Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
All	5	9
	=5/14 0.36	=9/14 0.64

Compute $P(\text{Yes}|\text{Sunny}) = P(\text{Sunny}|\text{yes}) * P(\text{yes})/P(\text{Sunny}) = 3/9 * 0.64/0.36 = 0.6$

Compute $P(\text{yes}|\text{overcast}) = P(\text{overcast}|\text{yes}) * P(\text{yes})/P(\text{overcast}) = 4/9 * 0.64/0.29 = 0.98$

$P(\text{No}|\text{Sunny}) = 0.4$

- To compute $P(\text{No}|\text{overcast})$, we need to compute $P(\text{overcast}|\text{No})$ what is the value!
- In this sample, there were no no-play days when it was overcast!
- Due to this, the probability of the computation $P(\text{No}|\text{overcast})$ becomes 0.
- This is corrected by adding a small value (usually 1) as correction called '**Laplace Correction**' (may be a parameter)



Naïve Bayes Classifier

- If one of the conditional probability is zero, then the entire expression becomes zero
- Most algorithms use one of the following or take a parameter for this correction
- Probability estimation:

$$\text{Original: } P(A_i | C) = \frac{N_{ic}}{N_c}$$

$$\text{Laplace: } P(A_i | C) = \frac{N_{ic} + 1}{N_c + c}$$

$$\text{m - estimate: } P(A_i | C) = \frac{N_{ic} + mp}{N_c + m}$$

c: number of classes

p: prior probability

m: parameter

HW: What is $P(\text{No}|\text{overcast})$ With Laplace Correction!

- M is a constant; p is a weight relative to observed data



Multiple attribute example

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

21



- Class: $P(C) = N_C/N$
 - e.g., $P(\text{No}) = 7/10$, $P(\text{Yes}) = 3/10$
- For discrete attributes:

$$P(A_i | C_k) = |A_{ik}| / N_k$$
 - where $|A_{ik}|$ is number of instances having attribute A_i and belongs to class C_k
 - Examples:
 - $P(\text{Status}=\text{Married} | \text{No}) = 4/7$
 - $P(\text{Refund}=\text{Yes} | \text{Yes}) = 0$
 - $P(\text{Refund}=\text{Yes} | \text{No}) = 3/7$
 - $P(\text{Refund}=\text{No} | \text{No}) = 4/7$

What is Gaussian Naïve Bayes?

- So far we have seen computations when the A's are **categorical**. But how to compute the probabilities when A is a **continuous variable**?
- If we assume that A follows a particular distribution, then you can plug in the probability density function of that distribution to compute the probability of likelihoods.
- If you assume the A's follow a Normal (aka Gaussian) Distribution, which is fairly common, we substitute the corresponding probability density of a Normal distribution and call it the **Gaussian Naïve Bayes**. You need the **mean** and **variance** of the A to compute this formula.



2/11/2021



© Sharma Chakravarthy

How to Estimate Probabilities from Data?

- For continuous attributes:
 - Discretize** the **range** into bins and thus transform the attribute into an ordinal attribute.
 - Probability density estimation:**
 - Assume attribute follows a normal distribution
 - Use data to estimate parameters of distribution (e.g., mean and standard deviation)
 - Once probability distribution is known, can use it to estimate the conditional probability $P(A_i | C)$

23



How to Estimate Probabilities from Data?

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Normal distribution:

$$P(A_i | C_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(A_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$

- One for each (A_i, C_j) pair

- For (Income, Class=No):

- If Class=No
 - sample mean = $770K/7 = 110K$
 - sample variance = 2975
 - Sample std = $\sqrt{2975} = 54.54$

$$P(\text{Income} = 120 | \text{No}) = \frac{1}{\sqrt{2\pi(54.54)^2}} e^{-\frac{(120-110)^2}{2(2975)}} = 0.0072$$

24



Example of Naïve Bayes Classifier

- Given a Test Record: $X = (\text{Refund}=\text{No}, \text{Married}, \text{Income} = 120K)$

$$P(X|\text{Class}=\text{No}) = P(\text{Refund}=\text{No}|\text{Class}=\text{No}) \\ \times P(\text{Married}|\text{Class}=\text{No}) \\ \times P(\text{Income}=120K|\text{Class}=\text{No}) \\ = 4/7 \times 4/7 \times 0.0072 = 0.0024$$

$$P(X|\text{Class}=\text{Yes}) = P(\text{Refund}=\text{No}|\text{Class}=\text{Yes}) \\ \times P(\text{Married}|\text{Class}=\text{Yes}) \\ \times P(\text{Income}=120K|\text{Class}=\text{Yes}) \\ = 1 \times 0 \times 1.2 \times 10^{-9} = 0$$

➤ Since $P(X|\text{No})P(\text{No}) > P(X|\text{Yes})P(\text{Yes})$

➤ Therefore $P(\text{No}|X) > P(\text{Yes}|X)$
=> Class = No

naïve Bayes Classifier:

$P(\text{Refund}=\text{Yes}|\text{No}) = 3/7$
 $P(\text{Refund}=\text{No}|\text{No}) = 4/7$
 $P(\text{Refund}=\text{Yes}|\text{Yes}) = 0$
 $P(\text{Refund}=\text{No}|\text{Yes}) = 1$
 $P(\text{Marital Status}=\text{Single}|\text{No}) = 2/7$
 $P(\text{Marital Status}=\text{Divorced}|\text{No}) = 1/7$
 $P(\text{Marital Status}=\text{Married}|\text{No}) = 4/7$
 $P(\text{Marital Status}=\text{Single}|\text{Yes}) = 2/3$
 $P(\text{Marital Status}=\text{Divorced}|\text{Yes}) = 1/3$
 $P(\text{Marital Status}=\text{Married}|\text{Yes}) = 0$

For taxable income:
 If class=No: sample mean=110
 sample variance=2975
 If class=Yes: sample mean=90
 sample variance=25

25



Example of Naïve Bayes Classifier

Name	Give Birth	Can Fly	Live in Water	Have Legs	Class
human	yes	no	no	yes	mammals
python	no	no	no	no	non-mammals
salmon	no	no	yes	no	non-mammals
whale	yes	no	yes	no	mammals
frog	no	no	sometimes	yes	non-mammals
komodo	no	no	no	yes	non-mammals
bat	yes	yes	no	yes	mammals
pigeon	no	yes	no	yes	non-mammals
cat	yes	no	no	yes	mammals
leopard shark	yes	no	yes	no	non-mammals
turtle	no	no	sometimes	yes	non-mammals
penguin	no	no	sometimes	yes	non-mammals
porcupine	yes	no	no	yes	mammals
eel	no	no	yes	no	non-mammals
salamander	no	no	sometimes	yes	non-mammals
gila monster	no	no	no	yes	non-mammals
platypus	no	no	no	yes	mammals
owl	no	yes	no	yes	non-mammals
dolphin	yes	no	yes	no	mammals
eagle	no	yes	no	yes	non-mammals

Give Birth	Can Fly	Live in Water	Have Legs	Class
yes	no	yes	no	?

A: attributes

M: mammals

N: non-mammals

$$P(A|M) = \frac{6}{7} \times \frac{6}{7} \times \frac{2}{7} \times \frac{2}{7} = 0.06$$

$$P(A|N) = \frac{1}{13} \times \frac{10}{13} \times \frac{3}{13} \times \frac{4}{13} = 0.0042$$

$$P(A|M)P(M) = 0.06 \times \frac{7}{20} = 0.021$$

$$P(A|N)P(N) = 0.0042 \times \frac{13}{20} = 0.0027$$

$P(A|M)P(M) > P(A|N)P(N)$
=> Mammals



Naïve Bayes (Summary)

- Robust to isolated noise points
- Handle missing values by ignoring the instance during probability estimate calculations
- Robust to irrelevant attributes
- Independence assumption **may not hold** for some attributes
 - Use other techniques such as Bayesian Belief Networks (BBN)



Thank You !!!



For more information visit:

<http://itlab.uta.edu>



13 December 2018



28



BDA 2018 (Warangal)