BANA 273 Session 4

Assignment Project Exam Help

Ex https://eduassistpro.github.io/ Add WeChat edu\_assist\_pro

Prof. Vibs Abhishek
The Paul Merage School of Business
University of California, Irvine

### Agenda

- Classification using Exact Bayes & Naïve Bayes
- Reminders
  - Assignment Project Exam Help
  - Assignmen https://eduassistpro.github.io/
  - Project proposal (1 para) d ck Canvas for all due dates dd WeChat edu\_assist\_pro
  - Project guidelines posted to Canvas (Announcements page)



## Big Picture View of Course Progress

- Databases, Data Warehousing, SQL
- RFM & Pivot Tables
   Assignment Project Exam Help
   Classificatio
- - Bayesian ( https://eduassistpro.github.io/
  - Decision Tree (ID3) eChat edu\_assist\_pro
- Association Rules
  - Apriori
- Clustering
  - K Means



## A classic: Microsoft's Paperclip

Assignment Project Exam Help

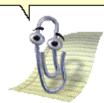
https://eduassistpro.github.io/

Add WeChat edu\_assist\_pro

It looks like you're writing a letter.

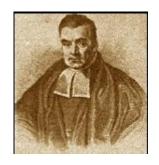
Would you like help?

- Yes, I need help
- just piss off and leave me alone!
- Don't show me this tip again





## **Exact Bayes**



**Thomas Bayes** 

#### For each record to be classified:

- 1. Find all other records pieculte and the predict https://eduassistpro.github.io/
- 2. Determine elong to and which class is more prevenue.
- 3. Assign that class to the new record

## Predict class attribute "Play" using Exact

Bayes

Outlook	Temp	Humidity	Windy	Play						
Sunny	Hot	High	False	No						
Sunny	Hot	High	True	No						
Overcast	Hot	High	False	Yes						
Rainy	Mild	Assignn	nealse p	rčije	t Ex	an	n Hel	n		
Rainy	Cool	Normal	False	Yes				P		
Rainy	Cool	Nor http	s://ed	luass	sisto	rO-	aithu	b.io/		
Overcast	Cool	Nor			arotp.		p.	Humidity	Windy	Play
Sunny	Mild	High Ado	d Hwe(	Chat	edu_	as	sist_	<b>PrO</b> ligh	False	?
Sunny	Cool	Normal	False	Yes						
Rainy	Mild	Normal	False	Yes	Outloo	ok	Temp.	Humidity	Windy	Play
Sunny	Mild	Normal	True	Yes	Sunn	У	Cool	High	True	?
Overcast	Mild	High	True	Yes						
Overcast	Hot	Normal	False	Yes						
Rainy	Mild	High	True	No						



#### **Notes**

- Bayesian classifier works best with categorical attribute Assignment Project Exam Help
  - Unlikely to merical variables
- Numerical a https://eduassistpro.github.jo/d and converted to https://eduassistpro.github.jo/d and and converted to https://eduassistpro.github.jo/d and converted to https://eduassist
- When the number of attributes is large (say 20), it becomes hard to find exact matches



## Exact Bayes – Cutoff Probability Method

- Establish a cutoff probability for the class of interest above which we consider that a record belongs Asthanniast Project Exam Help
- Find all the thttps://eduassistpro.giththe.inew record
- Determine the probability records belong to the class of interest
- If that probability is above the cutoff probability, assign the new record to the class of interest



### Example – Exact Bayes

	Sunny	Overcast	Rainy	Total
Play=Yes	2	3	2	7
Play=No	Assignme	ng Project I	zxam Heip	16
Total	5 https:	<del>//eduassist</del>	pro.aithub.	. 23 iO/

 $\begin{array}{c} Add \ WeChat \ edu\_assist\_pro \\ \hbox{P(Play=Yes \mid outlook=sunny)} \end{array} =$ 

P(Play=Yes | outlook=overcast) = 25%

 $P(Play=Yes \mid outlook=rainy) = 33\%$ 

Conclusion: No matter what the outlook, predict Play = No

**Cutoff probability method**: Specify cutoff probability p

If Probability(Play=Yes | outlook = ?) > p then predict Play = Yes

Suppose p = 37%

Under what outlook would we forecast play = Yes?



Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu\_assist\_pro Prof. Vibs

The Paul Merage School of Business University of California, Irvine

## **Conditional Probability**

• Rules of probability:  $P(A_1, A_p \mid B=1)$ 

https://eduassistpro.github.io/
This is correct only if the events A<sub>1</sub>
Add WeChat edu\_assist\_pro

• Let's start by assuming that they are, then:

P(Outlook=Sunny, Temp=High| Play=Yes) =
P(Outlook=Sunny| Play=Yes) \* P(Temp=High| Play=Yes)



## Apply Bayes' Rule

$$P(B \mid A) = \frac{P(A \mid B)P(B)}{P(A)}$$

B = the eveats spanned Project Exam Help

A = the event " p = High" https://eduassistpro.github.io/

P(Play = Yes | Outlook Withhy renocultions ist\_pro

$$= \frac{P(Outlook = sunny, Temp = High | Play = Yes) \cdot P(Play = Yes)}{P(Outlook = sunny, Temp = High)}$$

$$= \frac{P(Outlook = sunny | Play = Yes) \cdot P(Temp = High| Play = Yes) \cdot P(Play = Yes)}{P(Outlook = sunny, Temp = High)}$$



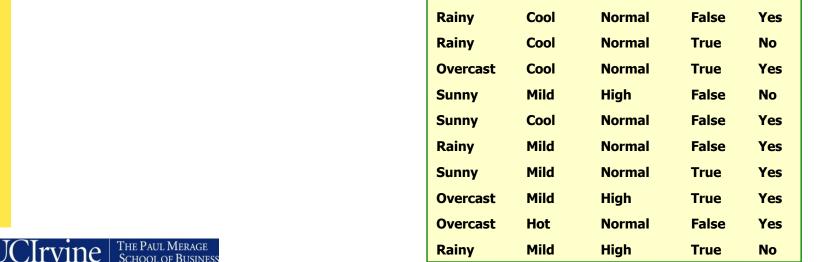
## Meaning of conditional independence

- P(outlook=sunny,Temp=High | Yes) with P(outlook=sunny,Yes); \* P(Femp=High | Yes)
- This means https://eduassistpro.github.io/onditional independence hetweedian edu\_assistempo
- If the conditional dependence is not extreme, it will work reasonably well



## Probabilities for weather data

				NO CC			)			GC	λία	•	
Ou	tlook		Tempe	rature		Hur	nidity		V	/indy		Pl	ay
	Yes	No		Yes	No		Yes	No		Yes	No	Yes	No
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	.2/9	m <mark>ent</mark>	High	3/9	4/5	False Help	6/9	2/5	9/	5/
Overcast	4/9	0/5	Mild	818911	ment	Projec Norma	6/9X		Help	3/9	3/5	14	14
Rainy	3/9	2/5	Cool	1.44	11		• 4		tate tar	/			
				htt	os://€	duass	ıstpr	o.g	h Fal	<del>U</del> /	ay		
				<b>Δ</b> d	d W	Sunny Chatte	2 du			ie No	0		
				Au	id VV	Rainy	Mild	_ Higl					
						Rainy	Cool	Nor					



## **Terminology**

- Frequency Chart also called contingency table (on previous Assiglement Project Exam Help
- Probability https://eduassistpro.github.io/ xcel – Pivot
- Create the c Add WeChat edu\_assist\_pro Table
- How to open ARFF file in Excel?
  - Launch Excel, Open File, Delimited, comma delimited
- Can also use SQL to compute entries in table.



## Probabilities for weather data

Oı	utlook		Temper	ature	ure Humidity		V	Vindy		P	lay		
	Yes	No		Yes	No		Yes	No		Yes	No	Yes	No
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	.2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/	5/
Overcast	4/9	0/5	Mild	818311	menu	Project	<b>6</b> /9	<u>am</u>	False Help	3/9	3/5	14	14
Rainy	3/9	2/5	Cool	htti	00://6	duocci	otor		thub				

Titips.//eduassistpro.gitriub.io/

• A new day: A dd We Chat edu\_assist property of the cool of the c



Ou	ıtlook		Tempe	rature		Hu	midity		V	Vindy		PI	ay
	Yes	No		Yes	No		Yes	No		Yes	No	Yes	No
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/	5/
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5	14	14
Rainy	3/9	2/5	Cool	3/9	1/5	Desis	4 17		T T = 1				

Assignment Project Exam Help

Outlook	Tem	Evidonos F
Sunny	Coohttps://eduassistpro	.github.to/luence E

Pr [yes|E]=Pr [Outlook=SulfnyWes]xhat redu\_assist\_propes]

×Pr [Humidity=High|yes]×Pr [Windy=True|yes]×Pr [yes]/Pr [E]

$$= \frac{\frac{2}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{9}{14}}{\Pr[E]} = \frac{0.0053}{\Pr[E]}$$



Ou	tlook		Tempe	rature		Hur	midity		V	Vindy		PI	ay
	Yes	No		Yes	No		Yes	No		Yes	No	Yes	No
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/	5/
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5	14	14
Rainy	3/9	2/5	Cool	3/9	1/5	Dunia	4 IV		T T = 1				

Assignment Project Exam Help

Outlook	Tem	Evidonos F
Sunny	Coohttps://eduassistpro	.github.to/luence E

$$= \frac{\frac{3}{5} \cdot \frac{1}{5} \cdot \frac{4}{5} \cdot \frac{3}{5} \cdot \frac{5}{14}}{\Pr[E]} = \frac{0.0206}{\Pr[E]}$$

#### Normalize...

• Pr[Yes | E] + Pr [No | E] = 1  
- Play can be either "Yes" or "No" 
$$\frac{0.0053}{Pr[E]} + \frac{0.0206}{Pr[E]} = 1$$
Assignment Project Exam Help

$$Pr[Yes \mid E] = \frac{0.00}{Pr[E] \text{ ttp}} / \frac{0}{Pr[E] \text{ ttp}} + \frac{1}{2}$$
Programme trope to Extend Trope to Ex

$$Pr[Yes | E] = \frac{0.0053 \text{ WeCharosdu_assist_pro}}{0.0053 + 0.0206}$$

$$\Pr[No \mid E] = \frac{0.0206}{0.0053 + 0.0206} = 0.795$$



## Example of Naïve Bayes Classifier

_						•
Name	Give Birth	Can Fly	Live in Water	Have Legs	Class	A: attributes
human	yes	no	no	yes	mammals	A. allibutes
python	no	no	no	no	non-mammals	Mumammala
salmon	no	no	yes	no	non-mammals	M: mammals
whale	yes	no	yes	no	mammals	N: non mammala
frog	no	no	sometimes	yes	non-mammals	N: non-mammals
komodo	no	no	no	yes	non-mammals	
bat	yes	yes 🔥 😋	monm	yes t Di	mammals+ T	xam Help
pigeon	no	yes 🔼	<b>Agnm</b>	yes I I	non mammals	zxam Heip
cat	yes	no	no	yes	mammals	
leopard shark	yes	no	ye		_	
turtle	no	no	sonttne	s://edi	uassist	pro.github.io/
penguin	no	no	so	<b>.</b> ,,, o a .		prorgitiraeno
porcupine	yes	no	no	yes	mammal	
eel	no	no	yes 11	mex/	non-ma	_assist_pro
salamander	no	no	yes sometimed	yes C	Hennacut	<u>ı_assisi_pro</u>
gila monster	no	no	no	yes	non-ma	
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	?



## Degenerate Probabilities (Pr[Outlook=Overcast|No)=0

- Could be a "true" representation of the real-world
  - Of course, one does not have to worry in that case
  - Rare Assignment Project Exam Help
- The training ugh
  - Is it EVER https://eduassistpro.githuh.jo/when "Play=no"? Add WeChat edu\_assist\_pro
  - If the answer is yes, a large uld have captured that fact
    - What does one do when data set is not big enough?
- We treat degeneracy seriously and try to remove it
  - Laplace approach



## The "zero-frequency problem"

Why does degeneracy matter? (e.g. "Humidity = high" for class "yes")

Pr [Humidity=High|yes]=0
Assignment Project Fxam Help

- Probability https://eduassistpro.github.io/
- (No matter how likely the o are!)
  Add WeChat edu\_assist\_pro
  Remedy: add 1 to the cou attrib value-class combination (Laplace estimator)
- Result: probabilities will never be zero! (also: stabilizes probability estimates)



Ou	Outlook Temperature		rature		Humidity		V	Vindy		Pl	ay		
	Yes	No		Yes	No		Yes	No		Yes	No	Yes	No
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/	5/
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5	14	14
Rainy	3/9	2/5	Cool	3/9	1/5	D .			TT 1				
- -	<ul> <li>Pretend that we add 3 rows of data containing only columns Outlook and Play:         <ul> <li>All 3 rows have pla</li> <li>1 row with Outlook</li> <li>Rainy. See resulting change in conditiona degenerate probability:</li> </ul> </li> </ul> <li>Pretend that we add 3 rows of data containing only columns Outlook and Play:         <ul> <li>All 3 rows have pla</li> <li>1 row with Outlook</li> <li>2 reast and 3<sup>rd</sup> with Outlook</li> <li>3 rows have pla</li> <li>4 reast and 3<sup>rd</sup> with Outlook</li> <li>5 reast and 3<sup>rd</sup> with Outlook</li> </ul> </li>												
Out	look		Tempe	rature		Hu	midity		V	Vindy		Pla	ау
	Yes	No		Yes	No		Yes	No		Yes	No	Yes	No
Sunny	2	4	Hot	2	2	High	3	4	False	6	2	9	8
Overcast	4	1	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	3	Cool	3	1								
Sunny	2/9	4/8	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/	8/
Overcast	4/9	1/8	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5	17	17
Rainy	3/9	3/8	Cool	3/9	1/5								

## Modified probability estimates

• In some cases, the number of rows to be added may need to be identificant and number of rows to be added setting we a

• Example: att https://eduassistpro.github.io/ s Play=No

Add WeChat edu\_assist\_pro

$$=(3+\mu/3)/\mu = (0+\mu/3)/\mu = (2+\mu/3)/\mu$$

Sunny Overcast Rainy



# Testing for Independence OPTIONAL (Information Theoretic Testing)

- Let A and B be two random variables
- Let D(A,B) = (H(A) + H(B) H(A,B))/H(A,B)
  - If A and Brigaindente Project Exam Help
    - H(A,B) =
    - D(A,B) = https://eduassistpro.github.io/
  - If A and B are linearly related edu\_assist\_prolated)
    - H(A,B) = H(A) = H(B)
    - D(A,B) = 1; this is the maximum
- If D() value is close to zero, assume independence
  - No need for looking up of statistical tables
  - Easy to implement



## Piecing it all together

- We want to estimate  $P(Y=1 | X_1,...,X_p)$
- But we don't shayon can under the profile  $X_1, \dots,$
- - True if we can assume (conditional) independence between  $X_1$ , ...,  $X_D$  within each class



## Piecing it all together

$$P(Y = 1 | X_1,...,X_p) = \frac{P(X_1,...,X_p | Y = 1)P(Y = 1)}{P(X_1,...,X_p)}$$
Assignment Project Exam Help

$$\approx \frac{P(X_1 \mid Y = 1) \cdot P(\text{thos}: \text{Pedpassistproble in training set})}{\text{Add} P(\text{VeChat edu}_assist_pro}$$

Proportion of rows with that predictor combination in the training set

Use the cutoff to determine classification of this observation. Default: cutoff = 0.5 (classify to group that is most likely)



### Advantages and Disadvantages

- The good
  - Simple
  - Can han Als sign amont broject des x am Help
  - High perform
  - Pretty robust https://eduassistpro.github.io/
- The bad
  - Need to categoried wind hat redu\_assist\_pro
  - Predictors with "rare" categories ility (Use Laplace fix)
  - No insight about importance/role of each predictor



## What is the probability of Play=Yes | Humidity=Normal and what would you predict for Play?

	Humidity High Humidity Normal Total
Play=Yes	5 7 12
Play=No	Assignment Project Exam Help
Total	https://eduassistpro.github.io/

A: 5/12, Predict Alaly We Chat edu\_assist\_pro

B: 7/19, Predict Play = Yes

C: 5/12, Predict Play = No

D: 7/19, Predict Play = No

E: None of the above



## Naive Bayes works better with categorical data because

- A: It takes less time to compute probabilities for categorical signment Project Exam Help
- B: It cannot contempts where different values for cannot contempt where different values is the contempt where different values are cannot contempt where different values is the contempt which is the contempt where different values is the contempt which is
- C: It needs the predicte Chart edu\_assish\_progome rows to compute accurate conditional probabilities
- D: Numeric data slows down the computation too much
- E: None of the above



## Data Preprocessing using Weka

- Follow steps on the following page:
- http://facesignmentaPuloject/iExamsHelplasses/ect5
  84/WEKA/p
  https://eduassistpro.github.io/
- File conversion and open edu\_assist\_in different applications
  - Excel, WordPad/TextEdit, Weka
  - CSV (text), XLSX (binary), ARFF (text)



#### Weka

• Run Naïve Bayes Classifier on cleaned and binned version of Abank clater of Exam Help

https://eduassistpro.github.io/

Add WeChat edu\_assist\_pro



#### **Next Session**

Testing and Validation

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu\_assist\_pro

