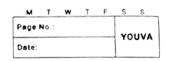
-> Laplace Additive Smoothing Son Cep and the state of t

М	T	W	T	F	S	s	
Page No.:							
Date:					YOUVA		

where w_1 w_2 w_3 w_2 w_3 w_4 w_5 w_1 w_2 w_3 w_4 w_5 w_6 w_1 w_2 w_3 w_4 w_6	
Eq: We have dataset Ceview Sentiment w_1 w_2 w_3 w_2 w_2 w_1 Unique words \rightarrow w_1 w_2 w_3 Sentiment Perfect 1 1 1 O(-12) Perfect 1 1 1 O(-12) Perfect 1 1 1 O(-12) Perfect 1 2 D (444) Perfect 1 1 1 O(-12) Perfect 1 O(-12) Perf	
Prize Sentiment W1 W2 W2 W1 W2 W2 W1 W2 W2 Prinque MDRds -> W1 W2 W3 Prinque 1 1 1 1 0 (-142) Prinque 2 1 0 2 (444) Prinque 3 1 2 0 (444) Prinque 3 1 2 0 (444) Prinque 3 1 2 0 (444) Prinque 401 New Acquire -> Ex: NIR got new Acquired -> Prinque 10 + NIR Acquired -> Prinque 10 +	
wi we we will we we we sentiment of the place Additive smoothing the surious is two or - we remisered. P(+ve review) P(+ve review) P(-ve review) P(-ve review) P(-ve review) P(-ve review) Reviewy P(-ve review) P(-ve review) P(-ve review) P(-ve review) P(-ve review) P(-ve review)	
wi we we we will we we sentiment of the service of taplace Additive smoothing for the service is to predict tokether the service is two or - we reviewed on the service of the service is to predict tokether the service is two or - we reviewed on the service is the service of the service is the service of the service of the service is the service of t	
Unique unords - wi wo was sentiment Applying Row - wi wo was sentiment runian 2 1 0 2 1 (4m) runian 3 1 2 0 1 (4m) existen 3 1 2 0 1 (4m) Existen 3 1 2 0 1 (4m) Existen 4 new roinian > { wi wo your wo you you wo y	
Unique widerds -> w1 w2 w3 Sentiment Applying BOW -> w1 w2 w3 Sentiment Scriew1 1 1 1 0 (-14)	†
· Unique wiords · Applying BOW · WI . W2 . W3 . Sentiment PRIVILED I I I O (-12) - LEUIEUR 2 I D 2	
Applying BOW > w1 w2 w3 Sentiment Surview 1 1 1 1 0 (-1/2) ruying 2 1 0 2 1 (4/4) ruying 3 1 2 0 1 (4/4) ruying 3 1 2 0 1 (4/4) Canderstanding rued of taplace Additive smoothing Ex: will got new register > { w1 w1 w2 } Our task is to predict whether the ruying is two or - we P(+ve reviewy) P(-ve runiewy) Apply Bowl on new qury Reviewy w1 w2 w3	
Applying BOW > w1 w2 w3 Sentiment Surview 1 1 1 1 0 (-1/2) ruyin 2 1 0 2 1 (4/4) ruyin 3 1 2 0 1 (4/4) ruyin 3 1 2 0 1 (4/4) **Charles 4 of taplace Additive smoothing Ex: Wile got new register > { w, w, w2 3° **Our tosk is to predict whether the ruyin is two or - we P(+ve reviewy) P(-ve runiewy) **Review 4 w1 w2 w3	
Applying BOW > w1 w2 w3 Sentiment Surview 1 1 1 1 0 (-1/2) ruying 2 1 0 2 1 (4/4) ruying 3 1 2 0 1 (4/4) ruying 3 1 2 0 1 (4/4) Canderstanding rued of taplace Additive smoothing Ex: will got new register > { w1 w1 w2 } Our task is to predict whether the ruying is two or - we P(+ve reviewy) P(-ve runiewy) Apply Bowl on new qury Reviewy w1 w2 w3	
Paylew 1 1 1 0 (-12) rayiow 2 1 0 2 1 (444) rayiow 3 1 2 D 1 (444) rayiow 3 1 2 D 1 (444) Canderstanding reed of taplace Additive smoothing Ex: WIL got new reiniew 2 & w, w, w 23° Our task is to predict tobether the rayious is two or - we P(+ve review 4) P(-ve rayiow 4).	
Lindenstanding ried of taplace Additive smoothing Ex: Will got new rejuiew > { w, w, w, 2} Our tosk is to predict whether the review is two or - we P(+ve reviewy) P(-ve reviewy) Reviewy w,	
Understanding reed of Laplace Additive smoothing Ex: will got new reiniew? { w,	
· Understanding reed of Laplace Additive smoothing • Ex: We got new xoixiew > { w, w, w2} • Own task is to predict Juhother the xurieve is two or - we P(+ve xeviewy) P(-ve xuriewy) • Apply Born on New qury Review y w1 w2 w3	
· Understanding red of Laplace Additive smoothing Ex: Will got new reiniew? { w,	
* Our tosk is to predict whether the review is two or - we P(+ve reviewy) P(-ve reviewy) * Apply Bont on New qury Fensewy wi we was	
P(+ve xeviewy) P(-ve xeviewy). P(+ve xeviewy) P(-ve xeviewy).	7
P(+ve reviewy) P(-ve reviewy). • Apply Bont on New gury Reviewy W1 W2 W3	
P(+ve review 4) P(-ve runiew 4). • Apply Bont on new query Review 4 w1 w2 w3	
P(+ve xeviewy) P(-ve xuriewy). • Apply Bont on new gray Reviewy wi w2 w3	
Review 4 w1 w2 w3	
Review 4 w1 w2 w3	
Reviewy w1 w2 w3	
Reviseury wi w2 w3	
3 , 0 , 0	
· P(+1/24) = P(W1=31+1/2) X, P(W2=0)+1/2) X P(W2=0)+1/2) XP	
)/_t)
$\frac{1}{2} \times \frac{1}{2} \times \frac{1}{2}$	'(न।
The state of the s)(+\ -
)(- 1\
P(-16 74) = P(w1=3)-16) XP(w2=0/-16) XP(w3=0/-16) XP(=	
$\frac{1}{3}$	
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	

Page No.: YOUVA
Both probability lucame zero, sero cans
paodict the sentiment
This is the problem in some cases the particular feature in new query
- may for a particular class cg ((w1=3 +w)
- that class, this loads to 0 of
- probability whole product will be
- O can me some it using log prob
- Ans: No, lurause log(0) is undéfined
This problem is solved by Laplace.
· It will prevent the particular term
to not lu come zerco
uith very small no 12-> 0.0000]
lut is senvally not good as
v
· Second soln -1" Claplace Additive smothing
- For LAS, you add something in nume.
$A \cup A \cup$
Ex: hu have Probabity -) [] the



L'is generally!

so even if we have =) (0/12 2 will grevent is from lu coming serio

Denominatos

· (itd =) what is this m?

on actually vory, It depends on which depends on which

· Examply: P(+14 x4)

= P(w1=3 | tre) X P(w2=0|tw) X P(w3=0|tw) X P(+w)

x 1 x 1 x ~1

Applying LAS (d=1, n=2)

1 nd 2 + nd 2 + nd 2 + nd 3 + 3nd

Binefit: Pauvents from total pub = 0

- But why did sue use complex thing instead
- The onsure is in Bias Trade off.
- Except for Gaussian Main Bayes, all

	M T W T F	5 5 VOUNA
	Dete:	AOUA
Bias variance tradeoff		
· Why +d?		
		6
· le rouse of Bias variance	1	. 6
entrol is high bias we will	con a work	
let's take sconario	ر ا	. 5
D & is small		
	1 .	
· let's say $\chi = 0$ (con't be -we) .	
- data. 11, 12, 13	1 .	•
		•
YUS - YUS	, 5	9
		•
1000 raws V		
er sous are labelled yes	2 50	
• P(Va. IV) P(V)		
· P(Yes X) = P(Y) XP(f) Y) X P(f2)	Y) e.	•
Let's say this hab	,	•
1. way small = 1	,	. •
500		tia
this indicates for walnu december &	anda siti a l	
	and the	1
Coly Con all a land of h	y.	
· El (an also be zero -) 0	luken	
500		

	M T W T F S S
	Page No.: YOUVA
-	Date:
	data is changed a little bid
	· Let's say une changed the date again
	$now \Rightarrow 3$
	· We can see, If training data is changed
•	little bit there Is too much vouiation
•	in 0/P, as when one turn is zero,
	whole grob was getting zero
u. 🖜	
	This indicate high variance, that
	means if a value is too small
4	there is a chance of oineyithing
•	(high mornance)
-	
_	(2) & walue is money large
_	
	L=1000,
	· ond particular feature probability was
	is facility and the second of
	50.0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
4	• Applying LAS = 1 + 1000 (Assume n = 2)
A	500,72(1000)
6	New Prob = 2 That feeture 5
A	of that feeture s
<u> </u>	· X = 10,000 ' 1400 A ' 1100 A
<u> </u>	1 + 10000 = 1
	500 +20000 2
0	
2	As you in wearing I am the
2	probabilities are reaching to half
	running

М	T	W	T	F	s	S
Page No.:					YOUVA	
Dete:					'	

Let's check again

tend to reach half

· same gois for P(N(X) =1 all of el's prob also reaches half

· Since both of their all prop are half.

The product of P(X|X) value is (JXJXJZJ-)
the " " P(N|X) " will also be (JXJXJ)

· This means beth P(YIX) and P(NIX)

The thing that will make P(Y|X) and P(X|X) and P(X|X) different will depend on deata rewhich class is more the or No $P(X|X) = P(Y) \times P(H|X) \times P(H|X) \cdot \dots \cdot \dots$

this Value will be more than P(N)

· and noue POI / s class will be selected

· What does it mean?

point, un will assign that class whicheur is more, your result

	Page No.: Page No.: VOUVA
•	which indicates underfitting (Migh Bias)
•	Condusion
•	
	· high alpha(L) => high Bias / under fitting
•	· vous alpha(2) = Nigh variance our fitting
•	
•	· & can be tuned, which means
0	
•	underfitting between ourlitting and
•	
-0	· Reasons for using laplace Additing Smoothly
-	
	- Prementing Prop to lacome zur
	I can reduce underfitting or overfitting
	using L
	sites of the control of the second of the se
	n value will be understand when
	un will study different types of
	Maire Bayes 1 00: