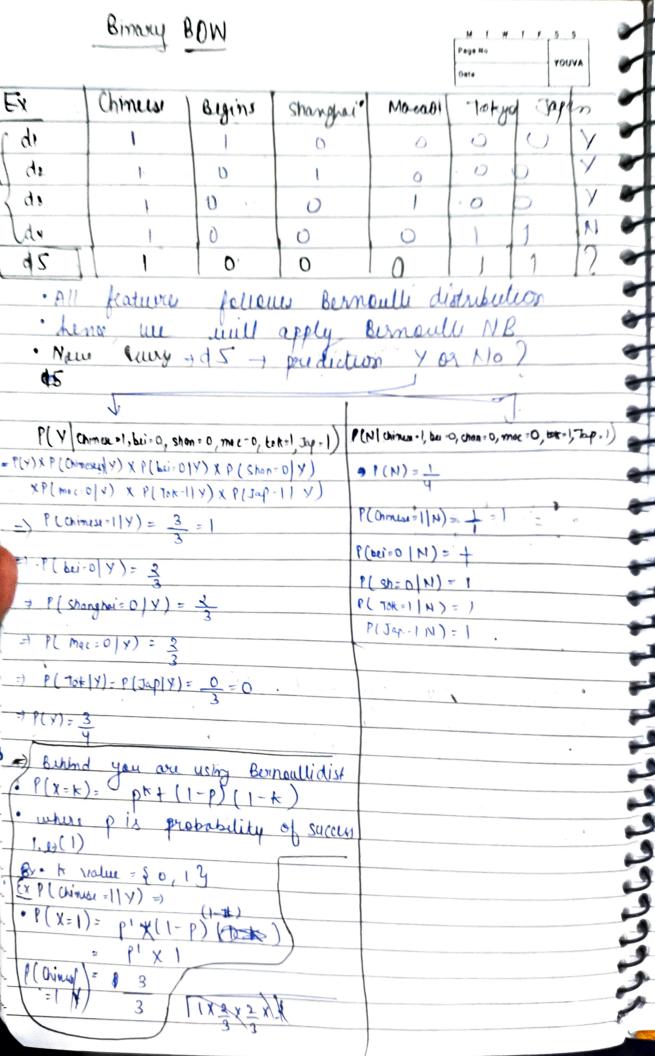
| LOGIC E | 3FHIND | TYPES | OF | Data: | YOUVA |
|--|-------------------------|--------------|---------|-------------|-----------|
| (1) Multinoulli | / Categorica | 1 MB | | | |
| -> Used u | then all ox con so foll | input fo | edeve | es are | abuteon |
| eg Data → | | . , | | | . 4 |
| | Gender | Pclass | | Surwind | |
| | n E | Pl | , | 1. | |
| 1 | | P2 | | 0 | • |
| | | P3 - | . 1 | D | |
| | | | - | | |
| • | 100 | | | | · |
| even though | - Gender | 18 leerno | lli | but as | |
| alle alle a | pplying | ategorical 1 | VB | we are | |
| even though due are a alsummin distribution | 7 0 11 7 | o be | otego | rical [mul | rinoulli' |
| | | | ••• | | |
| · Acc to -1 | Vaire Baul's | , | · | * | |
| ~~~~ | | • | , , , | | |
| new g | ury = & N | 1. F3.3- | + kid | 'ct -10,1 | |
| , | 0 | | | | , |
| PIIIm | | | 1 | | Q. |
| P(1 m, p | 3) | P | (0)r | M1 P3) | |
| = P(1) x P(M/1) X F | P/ P3/1) | = 12(0) | X P/m | 10) X P (P3 | 1-1 |
| P(1) = 392 | | | | | 10) |
| 89) | | [≥ Y(0) | = 49 | 11 | |
| $=) P(M 1) = \frac{2}{3}$ | 1, 1, 1 | =) P(m) | 70)= | | <u>)</u> |
| | 1 (2 | | | | |
|) f(p3/1). = :1 | | 7 7 | Y3 [0): | 8 11.5 | 1 |
| 2 a . V | 1 | • | | 1 7 7 1 | |
| 392 X 2 X 1= | 0.6 | × 490 | | X 5 = 0. | 9 |
| 091 3 5 | | 091 | 15 | 0. | |

=>

| Page No Date: YOUVA |
|---|
| · Sina P(0 m, P3) is nightst. O will be prodicteon ie person will not survive |
| · that's how categorical / multinolli works · and LAS will be applied as same |
| te the temporally del n 18 to the month of categories in column |
| 2) Bernoulli Naine Bayes |
| · works well on data whom lack feature follows Bernoulli distribution is bimory features (2 categories) |
| 1 0 |
| |
| · where can we find this type of data! · > Using BOW (binary) |
| -) In this, it only tells whether the word is there or not instead of - |
| • Ex -1 |



| W | I | NF. | 7 | £ | 5 | 1 |
|-------|-----|-----|---|---|----|-----|
| Parpe | No. | | | | | |
| | | - | - | | 40 | AVE |

| | | | | | 1 | |
|--------|-----|--------|------|----|------|--|
| That's | how | Burnos | ılli | Û. | used | |

And Because in Burnoulli, you find push of Hoa

- · Breams in book Bernaulti you not only find the prob of word present but also prob
- . In Burnoulli + P(Y|X) indicates word is present =) P(N|X) " " is not "

Categorial

In categorical, you only account of for people of word present

Ex classification - Dog, (at, house

- · No of prob will be P(Dog|x), P(Cod|x)
 P(porse | boox).
- net a dogs, prob of not cat, as horse

LAS

· In some way ... also apply LAS

Multinomial Naive Bayes when · seatures arc. discrete EX Shanghai Chinuse Bujing Mocous Toty Jup 0 d3 0 0 6 0 0 0 0 for · Prodict

| Page No.: | |
|-----------|-----|
| | |
| Date: | JVA |

· Tuo puob will be calculated 1) P(Y | Chimes=3 | bei=0 | sha=0 | dap=0 | mac=0 | tok=1, Tap=1) = 2001 · Note • In other types, ex Bernoulli NB, une considered all features to follow Bernoulli distribution Inp 1 80, similarly In Multinomial, what is the sning that follows multinomial distribution? Ans you will towat each column as a multimoully distribution and each score multinomial distribution Touat crucy new sentence as bag which · and whenever new query comes, effectively you are asking tell what is the propability of 3 chinese words, I takyo word, I Japan word Bimida to Ou get placed I student apt out and 6 not placed. Similarly we are asking prob of 3 Chimise words, I tokyo and I tag

of the 5 words taken out

| Page No.: YOUYA | |
|--|---------------|
| Dete: | |
| Can be rewritten as: | - |
| chinese chinese chinese tokego Japan P(chinese Y) X P(chinese Y) X P(chinese Y) X P(totyo Y) X P/R | |
| P(Chimise Y) X P(chimise / Y) X P(chimise / Y) X P/ Tokyo (Y) X P/R | 14) |
| .0 | |
| Durat is this P (Chimese Y)) | 9 |
| No. 1 March 1 March 2 | 3 |
| Ans Out of total Yes words, how many are climse | 9 |
| Chimse (Y) | 9 |
| 0 | • |
| | 0 |
| · A STATE OF THE S | 9 |
| =) What le D/tokunly 9 | 9 |
| Correct Is P (lorgo Y-) / | |
| P) (Japly) is also same | - |
| | - |
| In short $P(Y X) =$ | - |
| | - |
| $= \left(\frac{5}{3}\right)^{3}\left(\frac{0}{0}\right)\left(\frac{0}{0}\right) \times P(x) = P(x x)$ | - |
| (8) (8) | • |
| | 8 |
| There are also prob of beiging and shary | 9 |
| but they are not prisent in tech | 0 |
| auta their roked to pour is o | 0 |
| => P(bei(y)0 = 1 Hence no need to | 0 |
| mention | 0 |
| land the way of the property of the same | 0 |
| · Refer to 04, you will understand | 9 |
| | 4 |
| Or en multinomial, un de also | 9 |
| multiplying it with permutation | () |
| multierging it with pernutation but in this we are not | |
| euhy? | |
| | A. |

ulty we are not multiplying it with !-5 / → no of total words in ywy

3/ 1/ 1/ → 3 rategorhs = chines = 3, 10/1, Jap 1 Ans: Because, we are only asking for that lest combination thin le. chinese chinese chinese tokyo Japan Note: · Since, nows represent multinomial that's cuty in laplace additue smoothing + n2

not no of features

not no of categories in features Out of core Maine Bayes can apply partial fix
useful suhen detaset is longe
Druide large data into chunks
Because dert data is usually large Training will continue from at that chunk part

Other algos also provide this method