

Recap

24 June 2023

09:32

Naive Bayes → working -



text data →

imdb →

Sentiment
analysis

Advanced Naive Bayes

Types of Naive Bayes

(1)

Log prob
(2)

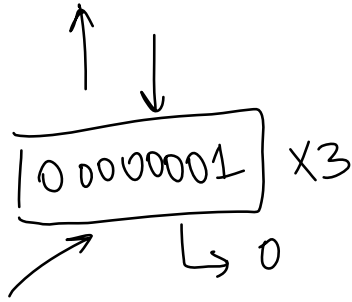
Laplace smoothing

(3)

[Numerical Stability]

23 June 2023 15:06

→ Underflow → represent decimals in memory



floating point
↓
binary 0 and 1

$$0.000001 < 0.00003$$

$$0 = 0$$

2500

100 wls

cgpa | iq

placement

{ 8.1, 81 } → y | N

$$P(y | 8.1, 81) = \frac{P(y)}{P(8.1 | y) P(81 | y)}$$

0 < x < 1

2500 prob

$$P(N | 8.1, 81) = 0$$

$$\log(P(y)) + \log(P(8.1 | y)) + \log(P(81 | y)) + 2500 \dots$$

$$0.1 \times 0.3 \times 0.4 \times 0.6$$

log probabilities = log-prob

$$\log(P(a) P(b) P(c) P(d) \dots)$$

$$\log(ab) = \log a + \log b$$

$$\log(0.5)$$

-ve number

$$y = -153$$

$$N = \lceil \frac{-135}{\dots} \rceil$$

N class ✓

$$\log(0.3 \times 0.5 \times 0.7)$$

$$\dots + \log(0.7)$$

$$\log(0.5) + \log(0.5) + \log(0.7)$$

$$\begin{array}{ccc} \text{"} & \text{"} & \text{"} \\ -1.2 & -0.7 & -0.3 \end{array} = \boxed{-2.2}$$

N

What is Underflow in Computing

24 June 2023 07:47

Underflow is a condition that can occur in computing when a number nears zero and the computer can no longer store it accurately in memory using floating-point representation. It happens when a calculated result is a smaller absolute value than the computer can actually represent.

Most computers use a form of representation called floating-point to represent real numbers. This representation has a certain precision limit, and it can only represent numbers between a certain minimum and maximum value. If a number is too close to zero (but not zero), it might be smaller than the smallest representable positive number in the machine's floating-point representation. When an operation on such small numbers is performed, the machine might round the result to zero, leading to a loss of precision.

Underflow can be a problem in certain domains, such as machine learning, where calculations often involve probabilities. Probabilities are positive numbers that can be very close to zero. When multiplying many small probabilities together, the result can underflow. One common way to avoid underflow in such scenarios is to perform calculations in the log domain, where addition and subtraction are used instead of multiplication and division, thereby maintaining higher numerical precision.

[Laplace Additive Smoothing] ←

23 June 2023 15:13

binary bow

review	sentiment
w_1, w_2, w_3	0
w_1, w_3, w_3	1
w_2, w_2, w_1	0

Bias Variance Trade-off

	w_1	w_2	w_3	sentiment
σ_1	1	1	1	0 -ve
σ_2	1	0	1	1 +ve
σ_3	1	1	0	0 -ve

1/3

$\sigma_4 \rightarrow w_1, w_1, w_1 \rightarrow +ve, -ve$

$$P(+ve | \sigma_4) = \frac{P(+ve)}{P(+ve) + P(-ve)} = \frac{1/3}{1/3 + 2/3} = \frac{1}{3}$$

$P(-ve | \sigma_4) = \frac{2}{3}$

$$0 = P(-ve | \sigma_4) = \frac{P(-ve)}{P(+ve) + P(-ve)} = \frac{2/3}{1/3 + 2/3} = \frac{2}{3}$$

$\log(0) = \text{undefined}$

epsilon

$\frac{1 + \alpha}{3 + n\alpha}$

$\alpha = 1$ = default

$n = 2$

$\alpha = 1$

vary \rightarrow depends

$\frac{1 + 1}{3 + 2(1)} = \frac{2}{5}$

Bias Variance Tradeoff

prob $\rightarrow 0 \rightarrow$

$\frac{0 + 1}{1 + 2(1)} = \frac{1}{3}$

0.1 \rightarrow 100000

Bias Variance Tradeoff

$P(-)$

hyper

$\frac{+ \alpha}{1 + n\alpha}$

$0 \rightarrow 0.00001$

$P(_)$

↳ Laplace smoo

$$= \frac{+\alpha}{+\eta\alpha} \rightarrow \text{flexibility} \rightarrow \text{control bias and variance}$$

high bias
 $\alpha \rightarrow$ low bias
 high varian $\alpha \rightarrow$ low variance

$\alpha = \text{small} \rightarrow \boxed{\alpha = 0} \text{ min } \boxed{\text{min} = 0} \quad 0.063 \quad 0.07$

$\frac{100}{Y} = f_1 | f_2 | f_3 \dots | Y = \frac{500}{500} \leftarrow \text{sample}$

$$P(Y|X) = P(Y) P(f_1|Y) P(f_2|Y) P(f_3|Y) \dots$$

500 rows
 $N=100$

$\frac{0}{500} \quad \frac{1}{500} \quad \frac{3}{500} \quad \frac{12}{500}$

high variance
 $n=2$

$\alpha = 0, 0.01, 0.007$

overfitting \leftarrow high variance

$\alpha = \text{very high}$
 $= 1000$

$$\frac{1+1000}{500+2000} = \frac{1001}{2500} = \frac{10}{25} = \boxed{\frac{2}{5}}$$

$= 10000$

$$\frac{1+10000}{500+20000} = \frac{10000}{20500} \approx \boxed{\frac{1}{2}}$$

$\alpha = 100000$

$\frac{100000}{2000500}$

$\alpha \uparrow \quad \frac{1}{2}$

$\frac{+\alpha}{+\eta\alpha} = \frac{1}{2}$

$\alpha = 1$

$$P(Y|X) = \frac{P(Y) P(f_1|Y) P(f_2|Y) \dots P(f_n|Y)}{P(Y) P(f_1|N) P(f_2|N) \dots P(f_n|N)}$$

$P(N|X) = \frac{P(N)}{1}$

same

$$P(N|X) = \frac{P(N)}{P(X|N) + \frac{P(N)}{2} + \dots}$$

data $\rightarrow Y$
 $\rightarrow N$
underfitting

$X \rightarrow$ query

$\alpha =$ hyperparameter
term

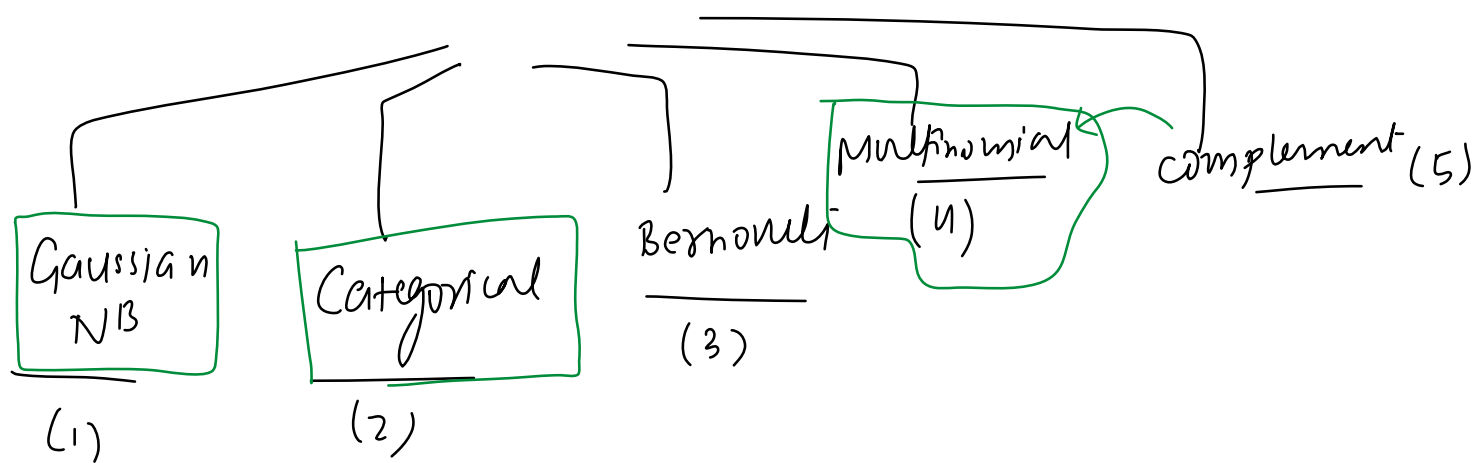
$\alpha \uparrow$ high
 lead to high bias
 or
 underfitting

$\alpha \downarrow$ low = 0
 high variance
 or
 overfitting

Types of Naïve Bayes

24 June 2023

07:48



Categorical Naïve Bayes

Categorical Naïve Bayes is a variant of the Naïve Bayes algorithm designed specifically to handle categorical data.

Data - all features are categorical

$x_q \in \{\text{sunny, Hot, High, False}\}$

y/N likelihood

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot ✓	High	False	No ✓
Sunny	Hot ✓	High	True	No ✓
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No ✓
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No ✓
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No ✓

↓

↓

$$P(Y | X_q) = \frac{P(Y)}{9/14} \quad \text{likelihood}$$

$$P(N | X_q) = \frac{P(N)}{5/14} \quad \text{likelihood}$$

Laplace additive smoothing

$$P(\text{out} = \text{sunny} | N) = \frac{3 + \alpha}{5 + \alpha \cdot 4} = \frac{3 + 1}{5 + 1(3)}$$

$$P(\text{temp} = \text{Hot} | N) = \frac{2 + 1}{5 + 1(3)} = \frac{3}{8}$$

Question on Categorical NB

24 June 2023

14:40

Let's say I have 4 features in my dataset 2 of them are categorical like gender and is married and 2 are numerical like age and height. Can I apply Categorical Naïve Bayes?

1. **Transform numerical features into categorical ones:** You can discretize numerical features by binning them into different categories. For example, you could create an "age group" feature that bins age into categories like "0-18", "19-35", "36-50", "51+". This allows you to treat the numerical feature as categorical, so you can use Categorical Naive Bayes. However, you should be aware that this may lead to loss of information, as binning reduces the granularity of the data.
2. **Use a mixed Naive Bayes model:** These models can handle both numerical and categorical data by making different assumptions for different types of features. For instance, numerical features could be modelled using a Gaussian distribution while categorical features could be modelled using a multinomial or categorical distribution. You might need to look for a different library or implement it yourself.
3. **Use another type of model:** Some machine learning models can naturally handle mixed data types. Decision trees and their ensemble variants (like random forests and gradient boosted trees) are capable of handling both numerical and categorical features without requiring any explicit feature transformation.

Bernoulli Naïve Bayes

23 June 2023 15:12

BernoulliNB implements the naive Bayes training and classification algorithms for data that is distributed according to multivariate Bernoulli distributions; i.e., there may be multiple features but each one is assumed to be a binary-valued (Bernoulli, Boolean) variable.

Therefore, this class requires samples to be represented as binary-valued feature vectors; if handed any other kind of data, a BernoulliNB instance may binarize its input (depending on the binarize parameter).

Table 13.1: Data for parameter estimation examples.

	docID	words in document	in $c = \text{China}$? —
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Chinese Chinese Chinese Tokyo Japan	?

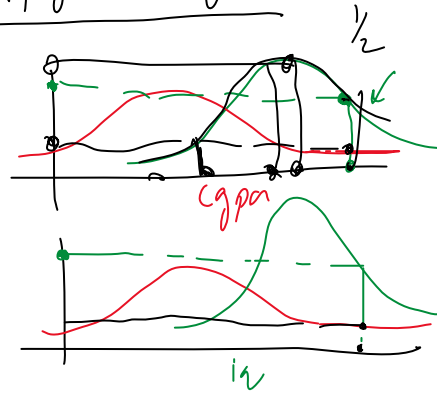
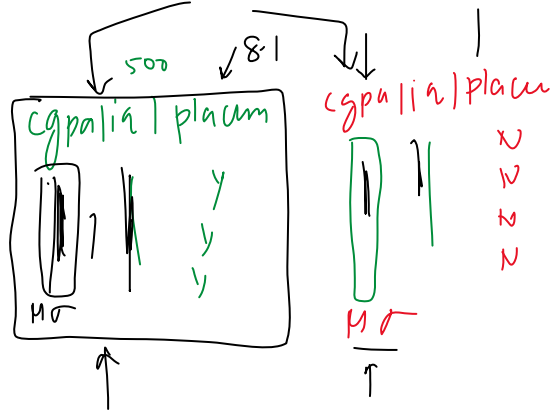
Data → all features are numerical ←

$\{8.1, 81\} \rightarrow Y / N \ 8.1 \ 81$

cgpa | iq | placement
↓ ↓
nume number
1000 students
500 Y 500 N

$$P(Y | \text{cgpa}=8.1, \text{iq}=81) = \frac{1}{2} \rightarrow \text{prior} \quad P(\text{cgpa}=8.1 | Y) \quad P(\text{iq}=81 | Y)$$

$$P(N | \text{cgpa}=8.1, \text{iq}=81) = \frac{P(N)}{1/2} \frac{P(\text{cgpa}=8.1 | N)}{1/2} \frac{P(\text{iq}=81 | N)}{1/2}$$



$$\text{cgpa} = 3.00012$$

Question on Gaussian NB

24 June 2023 16:17

Why Laplace Additive Smoothing not applied on Gaussian Naïve Bayes?

Multinomial Naive Bayes

23 June 2023 15:08

textual data Multinomial naive bayes → discrete

Multinomial Naive Bayes is a variant of the Naive Bayes algorithm that is particularly suited for classification tasks involving discrete features, such as text classification where features correspond to word counts or frequencies within the documents.

$$P(w|c) = (T_{c,w} + 1) / (\text{text}_c + B)$$

Here, $T_{c,w}$ is the count of word w in class c , text_c is the total count of words in class c , and B is the size of the vocabulary.

	f_1	f_2	f_3	f_4	y
→	2	1	4	3	0
	1	2	10	11	1
	0	5	6	11	1
	9	1	2	3	0

Table 13.1: Data for parameter estimation examples.

	docID	words in document	in $c = \text{China?}$
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Chinese Chinese Chinese Tokyo Japan	?

BOW → binary
non-bin (count)

movies review | sentiment
→ great-mov
fine —
— —
1
0
1
→ BOW discrete
great faith epic
→ $\frac{2(1)}{0}$ 0 0
0 3(1) 2(1)
1 1 1
(count) bow → Multinomial Naive Bayes
fractions → Tf-idf

Table 13.1: Data for parameter estimation examples.

	docID	words in document	in $c = \text{China?}$
training set	1	Chinese Beijing Chinese	yes
	2	Chinese Chinese Shanghai	yes
	3	Chinese Macao	yes
	4	Tokyo Japan Chinese	no
test set	5	Chinese Chinese Chinese Tokyo Japan	?

$$P(w|c) = (T_{c,w} + 1) / (\text{text}_c + B)$$

Here, $T_{c,w}$ is the count of word w in class c , text_c is the total count of words in class c , and B is the size of the vocabulary.

$$B = n \alpha$$

$$B = n \alpha \quad \alpha = 1$$

$$B = n$$

$$\frac{3}{4} \quad \frac{1}{8} \quad \frac{1}{8}$$

	1	2	3	4	5	
	chinese	Beijing	Shanghai	Macao	Tokyo	Japan
d1	2	1	0	0	0	Y
d2	2	0	1	0	0	Y
	2	0	1	0	0	Y

$$P(\text{china} | Y) = \frac{1}{8}$$

a1	0	0	0	0	0	0	Y
d2	2	0	1	0	0	0	Y
d3	1	0	0	1	0	0	Y
d4	1	0	0	0	1	1	N
d5	3	0	0	0	1	1	?

$$p(\text{beign} | Y) = 1$$

$$p(\text{chinese} | Y) = \frac{1}{8}$$

$$\frac{5}{8} \rightarrow \frac{0}{8} \frac{0}{8}$$

$$\frac{0 + \alpha}{8 + n\alpha} \quad \alpha = 1$$

$$\frac{0 + 1}{8 + 6 \times 1} = \frac{1}{14}$$

size of the vocab

$$p(Y | \text{chinese} = 3, \text{be} = 0, \text{sha} = 0, \text{ma} = 0, \text{tor} = 1, \text{jap} = 1) = p(\text{chinese} | Y)^3 p(\text{be} | Y)^0$$

$$p(N | \text{chinese} = 3, \text{be} = 0, \text{sha} = 0, \text{ma} = 0, \text{tor} = 1, \text{jap} = 1)$$

$$p(Y) \Rightarrow \frac{3}{14} \quad p(\text{chinese} | Y)^3 p(\text{beign} | Y)^0 \dots p(\text{Japan} | Y)^1$$

$$p(\text{chinese} | N) = \frac{1+1}{3+6(1)} = \frac{2}{9}$$

$$p(N) = \frac{1}{14}$$

multinomial distribution

Complement Naïve Bayes

23 June 2023 15:12

Out of Core Naïve Bayes

23 June 2023 15:13