

Naive Bayes

* Probabilistic based classification algo.

conditional probability (wiki = conditional probability example section).

$$P(A|B) = Pr(A=a|B=b) = \frac{P(A \cap B)}{P(B)}, P(B) \neq 0$$

* read eq's in english to get idea. $P(\text{dice}=3 | \text{dice is odd}) = \frac{1/6}{1/2} = \frac{1}{3}$

Independent Events & Mutually Exclusive Events

A & B are said to be independent if

$$\begin{cases} P(A|B) = P(A) \\ P(B|A) = P(B) \end{cases}$$

A: getting value of 6 in die 1 throw.

B: getting a value of 3 in die 2 throw.

A & B are said to be mutually exclusive iff

$$P(A|B) = P(B|A) = 0.$$

A: getting value of 6 in dice

$$\text{when } P(A \cap B) = P(B \cap A) = 0$$

B: " " " 3 " "

Bayes' Theorem

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

posterior ← $P(A|B)$ likelihood ← $P(B|A)$ prior ← $P(A)$ evidence ← $P(B)$

wiki - Bayes' Theorem
Problem - A more complicated example.

Proof:- $P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A, B)}{P(B)}$

$A \cap B = B \cap A$ — set theory.

$$\text{So, } P(A|B) = \frac{P(B, A)}{P(B)}$$

$$P(B|A) = \frac{P(B, A)}{P(A)} \Rightarrow P(B, A) = P(B|A) \times P(A)$$

$$\text{So, } P(A|B) = \frac{P(B|A) P(A)}{P(B)} \quad \text{— Hence, Proved.}$$

Ex. 3 machines. M_1 : 20% o/p.
 M_2 : 30% o/p.
 M_3 : 50% o/p.

fraction of defective items for M_1 : 5%
 M_2 : 3%
 M_3 : 1%

Item chosen at random,
 found to be defective. Probability
 that it belongs to 3rd m/c??

$$\rightarrow P(M_3|D) = \frac{P(D|M_3)P(M_3)}{P(D)}$$

$$P(D) = \sum_{i=1}^3 P(D|M_i) \times P(M_i) = 0.2 \times 0.05 + 0.3 \times 0.03 + 0.5 \times 0.01 = 0.024$$

$$P(M_3|D) = \frac{(0.01) \times (0.5)}{0.024} = \frac{5}{24} \text{ Ans}$$

Naive Bayes Algorithm (Naïve - Naive Bayes classifier).

Consider a data pt with n -features & K -classes. (c_1, c_2, \dots, c_K)

* If $K=2$, then binary classification

$$x = \langle x_1, x_2, \dots, x_n \rangle$$

task: Given x (all n features), we have to predict its class label.
 i.e. $\boxed{P(c_k|x)} \forall k \in [1, K]$ - whichever is highest, we'll select that as the class label.

$$P(c_k|x) = \frac{P(x|c_k)P(c_k)}{P(x)} \rightarrow P(x \cap c_k) = P(x, c_k)$$

Denominator constant for all class. So, only numerator matters.

So, $P(c_k|x) \propto P(c_k, x)$ - whichever has highest val. of $P(c_k, x)$ will be ans.

$$P(c_k, x) = P(c_k, x_1, x_2, x_3, \dots, x_n)$$

$$P(c_k, x_1, x_2, \dots, x_n) = P(x_1, x_2, x_3, \dots, x_n, c_k)$$

$$= P(x_1|x_2, x_3, \dots, x_n, c_k) \cdot P(x_2|x_3, x_4, \dots, x_n, c_k) \cdot \dots \cdot P(x_{n-1}|x_n, c_k) \cdot P(x_n|c_k) \cdot P(c_k)$$

$$= P(x_1|x_2, x_3, \dots, x_n, c_k) \cdot P(x_2|x_3, x_4, \dots, x_n, c_k) \cdot \dots \cdot P(x_{n-1}|x_n, c_k) \cdot P(x_n|c_k) \cdot P(c_k)$$

$$= P(x_1|x_2, \dots, x_n, c_k) \cdot P(x_2|x_3, \dots, x_n, c_k) \cdot \dots \cdot P(x_{n-1}|x_n, c_k) \cdot P(x_n|c_k) \cdot P(c_k)$$

Assuming that each feature x_i is conditionally independent of every other feature

$$\{P(A|B, C) = P(A|C) \text{ if } A \text{ \& } B \text{ are independent}\}$$

$$P(c_k, x_1, x_2, \dots, x_n) = P(c_k)P(x_1|c_k)P(x_2|c_k) \dots P(x_n|c_k)$$

$$= P(c_k) \prod_{i=1}^n P(x_i|c_k)$$

$$\text{So, } P(c_k|x_1, x_2, \dots, x_n) \propto P(c_k) \prod_{i=1}^n P(x_i|c_k)$$

$$= \frac{1}{Z} P(c_k) \prod_{i=1}^n P(x_i|c_k)$$

$$Z = P(x)$$

shatterline.com for example

Naive Bayes on text data

↳ very widely used.

SPAM filter: email $\begin{cases} \text{SPAM} \\ \text{not-SPAM} \end{cases}$

Text classification

NB widely used.

review $\begin{cases} \text{True} \\ \text{False} \end{cases}$

Task: $r(y=1|\text{text}_a), r(y=0|\text{text}_a)$

text \rightarrow stopwords stemming n-gram \rightarrow pre-processing \rightarrow binary bag of words, vector representation

text $\rightarrow \{w_1, w_2, w_3, \dots, w_n\}$ \rightarrow BoW features.

$$P(y=1|text) = P(y=1) * P(w_1|y=1) * P(w_2|y=1) * \dots * P(w_n|y=1).$$

$$P(y=1|text) = P(y=1) \prod_{i=1}^n P(w_i|y=1)$$

$$\text{Similarly } P(y=0|text) = P(y=0) \prod_{i=1}^n P(w_i|y=0).$$

$$P(y=1) = \frac{\# \text{ train pts with } y=1}{\text{total train pts.}} \quad P(y=0) = \frac{\# \text{ train pts with } y=0}{\text{total train pts.}}$$

$$P(w_i|y=1) = \frac{\# \text{ data pts which contain word } w_i \text{ \& } y=1}{\# \text{ data pts with } y=1}$$

Similarly for $y=0$.

NB is used as a baseline model (as its very simple) for text classification.

Laplace Smoothing

Q: what if the query text has a word that wasn't present in training set?

→ say text = $\{w_1, w_2, w_3, w'\}$ ^{new word.}

$$P(y=1|text) = P(y=1) P(w_1|y=1) P(w_2|y=1) P(w_3|y=1) P(w'|y=1).$$

If we completely ignore $P(w'|y=1)$ then it'll be same as assuming $P(w'|y=1)=1$ which is absurd.

$$\text{If we go by defn, } P(w'|y=1) = \frac{\# \text{ data pts where } w' \text{ occurs \& } y=1}{\# \text{ data pts with } y=1} = \frac{0}{n_1} = 0.$$

This will make the entire probability 0.

So, we do Laplace smoothing (additive smoothing).

$$P(w'|y=1) = \frac{0 + \alpha}{n_1 + \alpha K}$$

K = no of distinct value the feature w' can take.

α - can be any numeric value typically taken as 1.

Here since we are using binary bow representation $K=2$.

Let $n_1 = 100$.

$$P(w'|y=1) = \frac{0 + \alpha}{100 + 2\alpha}.$$

$$\text{Case 1: } \alpha = 1. \quad P(w'|y=1) = \frac{1}{102}$$

$$\text{Case 2: } \alpha = 10000 \quad P(w'|y=1) = \frac{10000}{20100} \approx \frac{1}{2}$$

Laplace smoothing when applied is applied to all the words, not to only the words not present in training.

$$P(w_i|y=1) = \frac{\# \text{ data pts with } w_i \text{ \& } y=1 + \alpha}{\# \text{ data pts with } y=1 + \alpha K}.$$

* As $\alpha \uparrow$, it moves the likelihood probabilities to uniform distⁿ.

So, if we have a likelihood with small num / denom, then we have less confidence in the ratios. So, if we have just large enough α , it'll smoothen its probab. value towards uniform distⁿ.

Log probabilities

Since all the probabilities are b/w 0 & 1, & also text features are having many features, continued product may cause underflow of data.

So, we take log of all the class priors & likelihood probabilities & add them together.

The class having max^m value is selected.

Since log is monotonically inc^t, its suitable for this purpose.

Bias variance trade-off

high bias — underfitting
high variance — overfitting

Only parameters — α

high α — smooths out likelihood probability towards uniform distrib.

case is $\alpha=0$ with small change in training data, ~~probabilities~~ likelihood probab. changes a lot, \Rightarrow high variance. \rightarrow overfitting

case is α large for all w_i , $P(w_i|y=1) \approx 1/2$ if α significantly high, so, underfitting.

So, how to find right α ?

α in naive bayes & k in KNN are hyper-parameters. These are calculated using cross validation.

Feature Importance and Interpretability

NB words with high value of

$P(w_i|y=1)$ \rightarrow important words/features in determining that point belongs to +ve class.

So, feature importance can be directly obtained from models.

{ +ve feature: — find words (w_i) with highest value of $P(w_i|y=1)$.

{ -ve feature: — find words (w_i) with highest value of $P(w_i|y=0)$.

Interpretability

$\alpha_1 \rightarrow y_1 = 1$
{ w_1, w_2, \dots, w_{10} }

I am concluding $y_1 = 1$ because x_2 contains words w_3, w_6, w_{10} which have a high value of $P(w_i|y=1)$,
 $P(w_3|y=1)$ $P(w_6|y=1)$ $P(w_{10}|y=1)$

and also $y_2 \neq 0$ because w_8, w_{11}, w_{20} .
 $P(w_8|y=0)$, $P(w_{11}|y=0)$, $P(w_{20}|y=0)$ are ^{very} small.

Imbalanced data

$n_1 \rightarrow n_1$ — (+)ve class
 $n_2 \rightarrow n_2$ — (-)ve class.

$n_1 \gg n_2$ Or $n_2 \gg n_1$
 \downarrow
Imbalance cond.

Say $\frac{n_1}{n} = 0.9$ & $\frac{n_2}{n} = 0.1$. $\rightarrow 0.9$.

$$P(y=1|w_1, w_2, \dots, w_n) = \frac{P(y=1)}{P(y=0)} \prod_{i=1}^n \frac{P(w_i|y=1)}{P(w_i|y=0)}$$

$$P(y=0|w_1, w_2, \dots, w_n) = \frac{P(y=0)}{P(y=1)} \prod_{i=1}^n \frac{P(w_i|y=0)}{P(w_i|y=1)}$$

If the likelihoods are almost the same, then the majority class has an advantage due to high class prior.

Solⁿ — ① upsampling/downsampling $P(y=1) = P(y=0) = 1/2$

② Simply drop the probabilities.

③ modified NB — not often used.

Another problem with imbalanced data.

(i)ve $n_1 = 900$, $\rightarrow P(w_i|y=1) = \frac{\boxed{}}{900} \rightarrow (0-900)$

(-ve) $n_2 = 100$, $\rightarrow P(w_i|y=0) = \frac{\boxed{}}{100} \rightarrow (0-100)$

So, in minority class have small num & denoms.

when we apply Laplace smoothing, the impact of α is more on minority than on majority.

Ex: minority $\frac{2}{100}$ $\rightarrow 2\%$

majority $\frac{18}{900}$ $\rightarrow 2\%$

for the same α , diff behavior. which

$\alpha = 10$
minority $\frac{2+10}{100+20} \approx 10\%$

majority $\frac{18+10}{900+10} \approx 3.04\%$

less change.

① upsampling/downsampling
② hacks — α_1, α_2 — diff values

Outliers

Outlier at last time is taken care by Laplace smoothing.
During training time - if occurs very few times, then ignore.
(say less than 10 times)

Missing value

① text data - no case of missing values.

② categorical data - say i^{th} data pt has its j^{th} feature missing
then consider the missing as a new category NaN

③ numerical features → use mean/median to figure out.

Numerical features

* Big assumption - feature f_j is Gaussian distⁿ.

Say we want to calc. $P(y=1 | x_{i1}, x_{i2}, x_{i3}, \dots, x_{id})$.
 \downarrow i^{th} data pt
 $= P(y=1) \prod_{j=1}^d P(x_{ij} | y=1)$.
 \downarrow j^{th} feature.
(trivial)

for $P(x_{ij} | y=1) \rightarrow$ consider all pts having $y=1$.
then calc mean & var. of j^{th} feature
say μ & σ .

Then we assume feature f_j to be a Gaussian distⁿ of $N(\mu, \sigma)$.

Now, we find $P(x_{ij} | y=1)$ by using the pdf of the Gaussian curve
& value of x_{ij} .

This is also called Gaussian Naive Bayes.

- * NB can't use similarity/distance matrix. It needs features.
- * NB can easily handle large dimensions.

Best and worst case of NB

① conditional independence - of features - assumption.

If true - NB performs very well.

If false - NB deteriorates.

② text-classification → NB used as the baseline model for
spam filtering, review polarity.

③ categorical features → extensively used.

④ real-value features → seldom NB used.

⑤ Interpretability & feature importance is great.
We can reason as to why the model gave a certain o/p.

runtime complexity = low.

train-time complexity = low.

runtime space - low (prior & likelihood only reqd).

* can easily overfit unless Laplace smoothing done with proper
 α from cross validation.