

1. Goal: implement a classifier that can distinguish between Economic + Onion words
Articles have been pre-processed.

V = Set of all words found in the articles (vocabulary)

$|V|$ = size of vocabulary

Article^(S) = $\langle X_1, \dots, X_n \rangle$ where $X_i \in \Sigma^{0,15}$

1 if i^{th} word appears in dictionary.

0 otherwise

Label space = {1, 2, 3}

Article^(S) = $\begin{cases} 1 & \text{if economic} \\ 2 & \text{if onion} \end{cases}$

Applying Naive Bayes:

2 assumptions

1) data is i.i.d

2) for each pair of features X_i, X_j

w/ $i \neq j$ we assume conditional independence given

label Y (this is the Naive Bayes assumption)

Naive Bayes > classification rule:

$$\hat{y} = \arg \max_y P(Y=y | X=x)$$

Can be rewritten using Bayes Rule

$$\hat{y} = \arg \max_y \frac{P(X=x | Y=y) P(Y=y)}{P(X=x)} \quad (\text{Bayes Rule})$$

$$\hat{y} = \arg \max_y P(Y=y | X=x) P(X=x)$$

$$= \arg \max_y P(X=X_1, \dots, X_n | Y=y) P(Y=y)$$

$$= \arg \max_y \left(\prod_{i=1}^n P(X_i=X_i | Y=y) \right) P(Y=y) \quad (\text{conditional independence})$$

A note about #params needed

- without conditional indep assumptions

$$2(2 - 1) \text{ params}$$

w/ conditional independence assumptions

$$\hookrightarrow 2n$$

a) Compute the log Product (x) where x is a vector of numbers in log space. and returns the product of those numbers in log space

\hookrightarrow The product of numbers who are already in

log space is simply the sum of all

~~vector components~~

$$\text{i.e. } \prod_i^n p_i = \sum_{i=1}^n \log p_i$$

Training NB.

b). Compute $\text{func}[D] = \text{NB_X}[\text{func}(X_{\text{Train}}, Y_{\text{Train}})]$

Output D : $2 \times V$ matrix where for any word

index $w \in \{1, \dots, V\}$ and class index $y \in \{1, 2, 3\}$

the entry $D(y, w)$ is the map estimate prior

$P(X_w=1 | Y=y)$ w/ a Beta (2, 1) distribution

$$\hat{D}_{\text{map}} = \frac{\#(\text{Observed } i \text{ s } 1 | Y=y) + \text{halfcounted } i's}{\#(Y=y) + \text{halfcounted } 1's + \text{halfcounted } 0's}$$

j 1). Separate the X_{Train} examples by classes

$$\downarrow 2) P(X_w=1 | Y=1) = \text{sum}(X_{\text{Train}, 1, \text{examples}}) + 1 / \text{size of } X_{\text{Train}, 1, \text{examples}} + 1$$

$$P(X_w=1 | Y=2) = \text{sum}(X_{\text{Train}, 2, \text{examples}}) + 1 / \text{size of } X_{\text{Train}, 2, \text{examples}} + 1$$

c) The rule for $p(y=1)$ is the empirical count of labels, $\#1 / \text{size of vector } y$.

$$\frac{\#1}{\#Y} = 1/3$$

$$1/3$$

d) wrote the function $\text{EqHAT} = \text{NB_Classify}(D, p, y)$

X is a $m \times V$ matrix containing m feature vectors

y = $m \times 1$ vector of predicted class labels

$y^{(t)}$ = the predicted label of feature vector x

The NB optimal classifier is

$$\hat{y} = \arg \max_{y \in Y} P(Y|D)$$

$$= \arg \max_{y \in Y} P(D|Y) P(Y) \quad (\text{Bayes Rule})$$

$$= \arg \max_{y \in Y} P(X_1, \dots, X_n | Y=y) P(Y=y)$$

$$= \arg \max_{y \in Y} \prod_{i=1}^n P(X_i | Y=y) P(Y=y)$$

1) We've previously calculated $P(X_i | Y=y)$

2) we apply PMF Bernoulli ($X_i | b$) where $b = P(X_i | Y=y)$

3) we take the Log Prod of the log of the result of

above and the log of the prior

4). $y^{(t)}$ = the label w/ maximum probability

e) Classification error: $\frac{\# \text{ of labels that agree}}{\# \text{ of total samples}}$

g) The test error is higher than the train error.

This is because we trained the classifier on the

training data. Testing the classifier using the

training data effectively amounts to having

memorized the data.

b). Retraining the classifier using the smaller training

set + attempting to reclassify the test samples

produces a lower successful classification

rate than using the larger training set. This is

intuitive if we look at the params of the classifier

$$P(X_i | Y=y) \text{ and } P(Y=y)$$

Certainly, larger data sets are more likely to

produce accurate estimates of the parameters. Also

since we are using the MLE to estimate the prior

we are likely to get inaccurate measurements. As

as a result, the prior has the potential to skew

the classifier towards the incorrect classifier.

i). Find the n most popular words per label :

Looking at the results for the $n=5$ most popular

words for each label we see

Economic most popular words are more economic

In nature (e.g. reform, servant, member)

This is appropriate for articles that mostly deal

with topics related to the political economy.