

# Natural Language Processing

[gpande3@gatech.edu](mailto:gpande3@gatech.edu)

## HW-2

### 1. A:

$$Q1: P(+|S) = ? \quad P(-|S) = ?$$

Using Conditional Independence of Naive Bayes:

$S = [ \text{"This", "Film", "was", "hilarious", "I", "didn't", "Yawn", "once", "Not", "single", "bland", "moment", "Every", "Minute", "laugh"} ]$

$$P(+|S) = P(+|) P(\text{hilarious}|+) P(\text{Yawn}|+) P(\text{bland}|+) P(\text{laugh}|+)$$

$$P(-|S) = P(-) P(\text{hilarious}|-) P(\text{Yawn}|-) P(\text{bland}|-) P(\text{laugh}|-)$$

$$P(+|) = \frac{3}{6} = 0.5$$

$$P(-) = \frac{3}{6} = 0.5$$

$$P(\text{hilarious}|+) = \frac{2}{13} \quad P(\text{hilarious}|-) = \frac{2}{16}$$

$$P(\text{Yawn}|+) = \frac{1}{13} \quad P(\text{Yawn}|-) = \frac{3}{16}$$

$$P(\text{bland}|+) = \frac{1}{13} \quad P(\text{bland}|-) = \frac{2}{16}$$

$$P(\text{laugh}|+) = \frac{5}{13} \quad P(\text{laugh}|-) = \frac{2}{16}$$

We will calculate log Probabilities, so

$$P(+|S) = P(+|) \sum_{i=1}^n \log P(\text{token}|+) = -3.757$$

$$P(-|S) = P(-) \sum_{i=1}^n \log P(\text{token}|-) = -3.737$$

$$P(-|S) > P(+|S)$$

Hence label assign to given sentence  $S$ :  $-$

B:

⑥ Recomputing Using Laplace add-1 smoothing:

$$P(+)=\frac{3}{6}=0.5$$

$$P(-)=\frac{3}{6}=0.5$$

$$P(\text{hilarious}|+)=\frac{2+1}{13+6}=\frac{3}{19}$$

$$P(\text{Yawn}|+)=\frac{1+1}{13+6}=\frac{2}{19}$$

$$P(\text{bland}|+)=\frac{1+1}{13+6}=\frac{2}{19}$$

$$P(\text{laugh}|+)=\frac{5+1}{13+6}=\frac{6}{19}$$

$$P(\text{hilarious}| -)=\frac{2+1}{16+6}=\frac{3}{22}$$

$$P(\text{Yawn}| -)=\frac{3+1}{16+6}=\frac{4}{22}$$

$$P(\text{bland}| -)=\frac{2+1}{16+6}=\frac{3}{22}$$

$$P(\text{laugh}| -)=\frac{2+1}{16+6}=\frac{3}{22}$$

$$\begin{aligned} P(+|s) &= \log(0.5) + \log(3/19) + \log(2/19) + \log(2/19) + \log(6/19) \\ &= -3.556 \end{aligned}$$

$$\begin{aligned} P(-|s) &= \log(0.5) + \log(3/22) + \log(4/22) + \log(3/22) + \log(3/22) \\ &= -3.637 \end{aligned}$$

$$P(-|s) < P(+|s)$$

s would be classified as "+" sentence.

C:

©: Additional feature that can be extracted from the text to improve classification is to make use of negation words like "didn't", "Not", as this words along with other words changes the meaning of the sentence. "didn't Yawn" = "Not Yawn" and "Not single" = "Every". These words can be useful in bigram or n-gram language modelling.

2:

(a). Likelihood equation for unregularized logistic Regression

$$L(w) = \prod_{i=1}^n (p(z^{(i)})^{y^{(i)}} (1 - p(z^{(i)}))^{1-y^{(i)}})$$

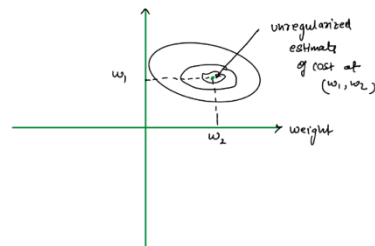
where  $p(z)$  is the conditional probability, i.e. sigmoid function.

$$p(z) = \frac{1}{1 + e^{-z}}$$

taking log of the likelihood function:

$$\begin{aligned} \lambda(w) &= \log(L(w)) \\ &= \sum_{i=1}^n y^{(i)} \log(p(z^{(i)})) + (1 - y^{(i)}) \log(1 - p(z^{(i)})) \end{aligned}$$

Now, the objective of this function is to minimize the cost by varying weights. for the unregularized cost, we could say that we find the global cost minimum for a particular value of weights, and these weights are unconstrained weights.

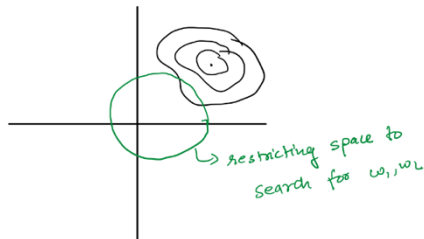


Now adding regularization term (L2)

L2:

$$\lambda(w) = \left[ \sum_{i=1}^n -y^{(i)} \log(p(z^{(i)})) - (1 - y^{(i)}) \log(1 - p(z^{(i)})) \right] + \lambda \|w\|_2^2$$

Addition of regularization means we are increasing the cost by the euclidean norm of the weight vector. So, we penalize the cost if we go too far on our weights to overcome overfitting. we restrict our search space to the regularized estimate.



Thus the weight will reduce in regularized logistic regression.

$$\|\theta^*\|_2^2 \leq \|\hat{\theta}\|_2^2$$

3: a:

$$\begin{aligned} \textcircled{3}. \textcircled{a} P(w) &= P(w_1, w_2, \dots, w_n) \\ &= P(w_1) P(w_2 | w_1) P(w_3 | w_1, w_2) \dots P(w_n | w_1, w_2, \dots, w_{n-1}) \quad - \textcircled{1} \\ &= \prod_{i=1}^n P(w_i | w_{i-1}) \end{aligned}$$

To make n-gram model work, we make an assumption

$$P(w_m | w_{m-1}, \dots, w_1) \approx P(w_m | w_{m-1}, \dots, w_{m-n+1}) \quad - \textcircled{2}$$

meaning we always only look for previous  $n-1$  tokens

So, putting value of eq<sup>n</sup> ② in eq<sup>n</sup> ①.

$$\begin{aligned} P(w) &= P(w_1) P(w_2 | w_1) P(w_3 | w_1, w_2) \dots P(w_m | w_1, \dots, w_{m-n+1}) \\ &= \prod_{m=1}^M P(w_m | w_{m-1}, \dots, w_{m-n+1}) \end{aligned}$$

for Bi-gram,  $n=2$

$$\begin{aligned} P(w) &= \prod_{m=1}^M P(w_m | w_{m-1}, \dots, w_{m-n+1}) \\ &= \prod_{m=1}^M P(w_m | w_{m-1}) \end{aligned}$$

B:

⑥: S: [BOS] I like cheese made at home [EOS]

$$\text{Count}([BOS]) = 4$$

$$\text{Count}([EOS]) = 4$$

$$\text{Count}(I) = 3$$

$$\text{Count}(\text{like}) = 2$$

$$\text{Count}(\text{cheese}) = 4$$

$$\text{Count}(\text{made}) = 3$$

$$\text{Count}(\text{home}) = 3$$

$$\text{Count}(is) = 2$$

$$\text{Count}(\text{at}) = 2$$

$$\text{Count}(\text{tasty}) = 1$$

$$\text{Count}(\text{salty}) = 1$$

$$\text{Count}(\text{that}) = 1$$

$$P([BOS] \mid \text{I like cheese made at home [EOS]}) =$$

$$P(I \mid [BOS]) \times P(\text{like} \mid I) \times P(\text{cheese} \mid \text{like}) \times P(\text{made} \mid \text{cheese}) \times P(\text{home} \mid \text{made})$$

$$= \frac{3}{4} \times \frac{2}{3} \times \frac{1}{2} \times \frac{1}{4} \times \frac{1}{3} \times \frac{1}{3}$$

$$= \frac{1}{12 \times 12}$$

$$= \frac{1}{144}$$

C:

- ①: Perplexity: It is the measure of the uncertainty.  
 Lower perplexity corresponds to higher likelihood  
 higher perplexity corresponds to lower likelihood.

$$\text{Perplexity}(w) = 2^{-\frac{L(w)}{M}}$$

In our case of bigram model.

$$L(w) = \log_2 \frac{1}{144}$$

$$\text{Perplexity}(w) = 2^{-\frac{1}{M} \log(1/144)}$$

$$M = 8 \quad [\because \text{total } 8 \text{ tokens in the given sentence}]$$

$$= 2^{\frac{1}{8} \log(144)}$$

$$= (144)^{1/8} \quad [\because 2^{\log_2 6} = 6^{\log_2 2} = 6]$$

D:

- ① Laplace smoothing.

We have given  $K = 0.1$

# of Unique tokens in the given vocab = 12

$$\text{Hence} \quad P(w_1, w_2) = \frac{\text{Count}(w_2, w_1) + 0.1}{\text{Count}(w_1) + 12 \times 0.1}$$

$$P(w) = \frac{3+0.1}{4+1.2} \times \frac{2+0.1}{3+1.2} \times \frac{1+0.1}{2+1.2} \times \frac{1+0.1}{4+1.2} \times \frac{1+0.1}{3+1.2} + \frac{1+0.1}{3+1.2}$$

$$= \frac{4.1}{5.2} \times \frac{3.1}{4.2} \times \frac{2.1}{3.2} \times \frac{2.1}{5.2} \times \frac{2.1}{4.2} \times \frac{2.1}{4.2}$$

$$= \frac{4.1 \times 3.1 \times 2.1^4}{5.2^2 \times 4.2^3 \times 3.2}$$

$$= 0.000977$$

E:

⑤: After adding [UNK] our vocab will become:

[BOS] i made cheese at home [EOS]

[BOS] i like home made cheese [EOS]

[BOS] cheese made at home is [UNK] [EOS]

[BOS] i like cheese [UNK] is [UNK] [EOS]

S: [BOS] i like pepperjack cheese [EOS]

adding [UNK] to the unknown vocab words.

S: [BOS] I like [UNK] cheese [EOS]

Now,  $\text{count}([BOS]) = 4$

$\text{count}(\text{home}) = 3$

$\text{count}([EOS]) = 4$

$\text{count}(\text{at}) = 2$

$\text{count}([UNK]) = 3$

$\text{count}(\text{'s}) = 2$

$\text{count}(\text{cheese}) = 4$

$\text{count}(\text{I}) = 3$

$\text{count}(\text{like}) = 2$

$\text{count}(\text{made}) = 3$

$$P(S) = P(I|BOS) \times P(\text{like}|I) \times P([UNK]|\text{like}) \times P(\text{cheese}|[UNK]) \\ \times P([EOS]|\text{cheese})$$

$$= \frac{3+0.1}{4+1} \times \frac{2+0.1}{3+1} \times \frac{0.1}{2+1} \times \frac{1+0.1}{3+1} \times \frac{1+0.1}{4+1}$$

$$= \frac{3.1}{5} \times \frac{2.1}{4} \times \frac{0.1}{3} \times \frac{1.1}{4} \times \frac{1.1}{5}$$

$$= \frac{3.1 \times 2.1 \times 1.1^2 \times 0.1}{5^2 \times 4^2 \times 3}$$

$$= 0.000656$$



#### 4:a:

④ ①% Most hateful words in Naive Bayes

[ "goid", "non-white", "dumb", "ape", "liberal", "scum", "asian", "nude", "non", "jeos"]

least Hateful Words

[ "-", "hanks", "html", "sf", "irishcentral", "facebook", "email", "10.00", "radio", "match"]

from the result, we can see that most hateful words are the common hate speech words and the model is able to identify it correctly. Also similarly for the least hateful words we can see that the trends follows the least hateful words in the vocab.

#### B:

⑥: For logistic Regression with regularization:

best tuned hyperparameters as

Learning Rate = 0.9 | Epochs = 2000

Regularization	Train Accuracy	Test Accuracy
0.0001	99.35	69.01
0.001	98.82	69.27
0.01	89.34	67.19
0.1	76.76	65.10
1	67.95	56.25
10	66.38	59.90

for unregularized version or

$\lambda = 0$

Train Accuracy = 99

Test Accuracy = 69

We can clearly observe from the above results and regularization parameter that when our lambda / regularization parameter reaches or close to zero, then it tends towards the overfitting solution which can be seen by the test accuracy approaching  $\approx 99.9\%$ , while the test accuracy does not improve much.

Another observation is if the lambda or regularization parameter is increased and reaches 1 or above model tends to underfit due to decrease in the train/test accuracy.

**C:** Implemented Bigram and remove all stop words using nltk library. Also I removed all the punctuations and numbers from the vocabulary, but since the size of the dataset was too small, the accuracy wasn't improved rather it decreased to 90 percent for train dataset and 50 percent for test dataset.