

School of Engineering and Applied Science (SEAS), Ahmedabad University

B.Tech (CSE Semester VI):
Machine Learning (CSE 523)

Project Submission #2: Linear Regression

Submission Deadline: March 04, 2020 (11:59 PM)

- Group No.: 8
- Project Domain: Natural Language Processing
- Project Title: Sentimental Analysis on movie reviews
- Name of the group members:
 1. Aditya Shah (1741007)
 2. Dhruvil Shah (1741024)
 3. Nilay Patel (1741038)
 4. Varun Patel (1741080)

Contents

1	Implementation code	2
2	URL links	4
2.1	Main folder	4
2.2	Codes	4
2.3	Data set	4
3	Inference	5
3.1	What we have done and why is it important?	5
3.2	How we have implemented?	5
3.3	Analysis and Results	5

1 Implementation code

```
1 import numpy as np
2 from sklearn.datasets import load_files
3 reviews_train = load_files("aclImdb/train/") #loading training file
4 text_train, y_train = reviews_train.data, reviews_train.target
5
6 print("Number of documents in train data: {}".format(len(text_train))) #printing
   document in train folder
7 print("Samples per class (train): {}".format(np.bincount(y_train))) #No.of
   samples
8
9 reviews_test = load_files("aclImdb/test/") #loading testing file
10 text_test, y_test = reviews_test.data, reviews_test.target
11
12 print("Number of documents in test data: {}".format(len(text_test))) #
   printing document in test folder
13 print("Samples per class (test): {}".format(np.bincount(y_test)))
14 from sklearn.feature_extraction.text import CountVectorizer
15 vect = CountVectorizer(min_df=5, ngram_range=(2, 2)) #to tokenize and to
   build new words
16 X_train = vect.fit(text_train).transform(text_train) # fitting the model and
   transform into vector
17 X_test = vect.transform(text_test) #transforming into vector
18
19 print("Vocabulary size: {}".format(len(vect.vocabulary_)))
20 print("X_train:\n{}".format(repr(X_train)))
21 print("X_test: \n{}".format(repr(X_test)))
22
23 feature_names = vect.get_feature_names()
24 print("Number of features: {}".format(len(feature_names))) #No.of features
25
26 from sklearn.model_selection import GridSearchCV
27 from sklearn.linear_model import LogisticRegression
28 param_grid = {'C': [0.001, 0.01, 0.1, 1, 10]} #Taking parameter with discrete
   value
29 grid = GridSearchCV(LogisticRegression(), param_grid, cv=5) #Searching by Cross
   validation = 5
30 grid.fit(X_train, y_train)
31
32 print("Best cross-validation score: {:.2f}".format(grid.best_score_))
33 print("Best parameters: ", grid.best_params_)
34 print("Best estimator: ", grid.best_estimator_)
35
36 import matplotlib.pyplot as plt
37 import mglearn
38 mglearn.tools.visualize_coefficients(grid.best_estimator_.coef_, feature_names,
   n_top_features=25) #to create figures
39 plt.show()
40
41 lr = grid.best_estimator_
42 lr.fit(X_train, y_train) #fitting the model
43 lr.predict(X_test) #predicting value
44 print("Score: {:.2f}".format(lr.score(X_test, y_test)))
45
46 pos = ["I've seen this story before but my kids haven't. Boy with troubled past
   joins military, faces his past, falls in love and becomes a man. "
47        "The mentor this time is played perfectly by Kevin Costner; An ordinary man
   with common everyday problems who lives an extraordinary "
48        "conviction, to save lives. After losing his team he takes a teaching
```

```

49 position training the next generation of heroes. The young troubled "
50 "recruit is played by Kutcher. While his scenes with the local love
51 interest are a tad stiff and don't generate enough heat to melt butter, "
52 "he compliments Costner well. I never really understood Sela Ward as the
53 neglected wife and felt she should of wanted Costner to quit out of "
54 "concern for his safety as opposed to her selfish needs. But her presence
55 on screen is a pleasure. The two unaccredited stars of this movie "
56 "are the Coast Guard and the Sea. Both powerful forces which should not be
57 taken for granted in real life or this movie. The movie has some "
58 "slow spots and could have used the wasted 15 minutes to strengthen the
59 character relationships. But it still works. The rescue scenes are "
60 "intense and well filmed and edited to provide maximum impact. This movie
61 earns the audience applause. And the applause of my two sons."]
62 print("Pos prediction: {}". format(lr.predict(vect.transform(pos))))
63
64 neg = ["David Bryce\'s comments nearby are exceptionally well written and
65 informative as almost say everything "
66 "I feel about DARLING LILI. This massive musical is so peculiar and over
67 blown, over produced and must have "
68 "caused ruptures at Paramount in 1970. It cost 22 million dollars! That is
69 simply irresponsible. DARLING LILI "
70 "must have been greenlit from a board meeting that said \"hey we got that
71 Pink Panther guy and that Sound Of Music gal... "
72 "lets get this too\" and handed over a blank cheque. The result is a hybrid
73 of GIGI, ZEPPELIN, HALF A SIXPENCE, some MGM 40s "
74 "song and dance numbers of a style (daisies and boaters!) so hopelessly old
75 fashioned as to be like musical porridge, and MATA HARI "
76 "dramatics. The production is colossal, lush, breathtaking to view, but the
77 rest: the ridiculous romance, Julie looking befuddled, Hudson "
78 "already dead, the mistimed comedy, and the astoundingly boring songs
79 deaden this spectacular film into being irritating. LILI is"
80 " like a twee 1940s mega musical with some vulgar bits to spice it up. STAR
81 ! released the year before sadly crashed and now is being "
82 "finally appreciated for the excellent film is genuinely is... and Andrews
83 looks sublime, mature, especially in the last half hour....."
84 "but LILI is POPPINS and DOLLY frilly and I believe really killed off the
85 mega musical binge of the 60s..... "
86 "and made Andrews look like Poppins again... which I believe was not
87 Edwards intention. Paramount must have collectively fainted "
88 "when they saw this: and with another $20 million festering in CATCH 22,
89 and $12 million in ON A CLEAR DAY and $25 million in PAINT YOUR WAGON...."
90 "they had a financial abyss of CLEOPATRA proportions with $77 million tied
91 into 4 films with very uncertain futures. Maybe they should have asked seer "
92 "Daisy Gamble from ON A CLEAR DAY .....LILI was very popular on immediate
93 first release in Australia and ran in 70mm cinemas for months but it failed "
94 "once out in the subs and the sticks and only ever surfaced after that on
95 one night stands with ON A CLEAR DAY as a Sunday night double. Thank "
96 "god Paramount had their simple $1million (yes, ONE MILLION DOLLAR) film
97 LOVE STORY and that $4 million dollar gangster pic THE GODFATHER "
98 "also ready to recover all the $77 million in just the next two years....
99 for just $5m.... incredible!"]
100 print("Neg prediction: {}". format(lr.predict(vect.transform(neg))))

```

2 URL links

2.1 Main folder

https://drive.google.com/open?id=1cfI15BJaI1iIsV5_aYnx3d6ygzzYa4pB

2.2 Codes

https://drive.google.com/open?id=18ea_fqozyvh5Itk6U9yIK_mxsIvK5beY

2.3 Data set

https://drive.google.com/open?id=1Rbs4SUiWbtkxpBzS7jrJ5_ep28tMCF9P

3 Inference

3.1 What we have done and why is it important?

We have done separation of positive and negative reviews by giving 1 and 0 and then we have first taken one word and then see accuracy on the training and unseen test data and then we have taken word pair and perform same operation. We have implemented tfidf model, so that large number of words and documents can be understood by computer, so we vectorize and transform it into numericals. This is important because large number of documents can be easily formed into small numbers so that vectorization and transformation can be easily done. We have applied logistic regression because we are giving input data and by grid estimator and parameter of cross validation and fitting model we are predicting magnitude whether it is negative or positive.

3.2 How we have implemented?

We have implemented this above method by giving path of test and train and then reading files and see whether document or text is positive or negative. If positive then append 1 and if negative then append -1 and then we vectorize the whole and transforming whole. Then we split data by stopwords and then taking one word or features calculating accuracy of system and then same with word pair. Then we implement tfidf to divide into more smaller into numerical so that computer can understand. The same we apply for testing data. After that as we have 1 and 0, we decided to implement logistic regression so that we can train our model with parameter by taking 5 discrete numbers as Cross validation points and then finding best grid estimator and then calculating magnitude and plotting graph and by best grid estimator we fit model with training or testing set and then predicting particular value and printing score prediction.

3.3 Analysis and Results

$$0 \leq R_\theta(x) \leq 1$$

$$h_\theta(x) = g((\theta)^T \cdot x)$$

$$g(x) = \frac{1}{1+e^{-x}}; g(x) \text{ is Logistic function}$$

$$h_\theta(x) = \frac{1}{1+e^{-\theta^T \cdot x}}$$

$$\text{Trainingset} = (x_1, y_1), \dots, (x_m, y_m)$$

$$m \text{ examples } x \in [x_0, x_1, \dots, x_n]$$

$$x_0 = 1, y = 0 \text{ or } 1$$

$$h_\theta(x) = \frac{1}{1+e^{-\theta^T \cdot x}}$$

Cost function is given by:

$$J_\theta = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_\theta x_i - y_i)^2$$

$$= \frac{1}{m} \sum_{i=1}^m \text{cost}(h_\theta x_i, y_i)$$

where,

$$\text{cost}(h_\theta(x), y) = \frac{1}{2} (h_\theta x - y)^2$$

$$\text{cost}(h_\theta(x), y) = -\log(h_\theta(x)) \text{ if } y=1$$

$$\text{cost}(h_\theta(x), y) = -\log(1 - h_\theta(x)) \text{ if } y=0$$

$$\text{cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1 - y) \log(1 - h_\theta(x))$$

If $y=1$,

$$\text{cost}(h_\theta(x), y) = -\log(h_\theta(x))$$

If $y=0$,

$$\text{cost}(h_\theta(x), y) = -\log(1 - h_\theta(x))$$

$$J_\theta = -\frac{1}{m} [\sum_{i=1}^m y_i \log(h_\theta(x_i)) + (1 - y_i) \log(1 - h_\theta(x_i))]$$

To fit parameter θ ,

Minimize J_θ w.r.t θ and get θ

To make a prediction given new x :

Output: $h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$

$p(y=1 | x, \theta)$

To minimize cost function,

Apply Gradient descent Algorithm and Repeat until convergence

$\theta_j = \theta_j - \alpha \frac{\Delta J}{\Delta \theta_j}$

Simultaneously update all θ_j

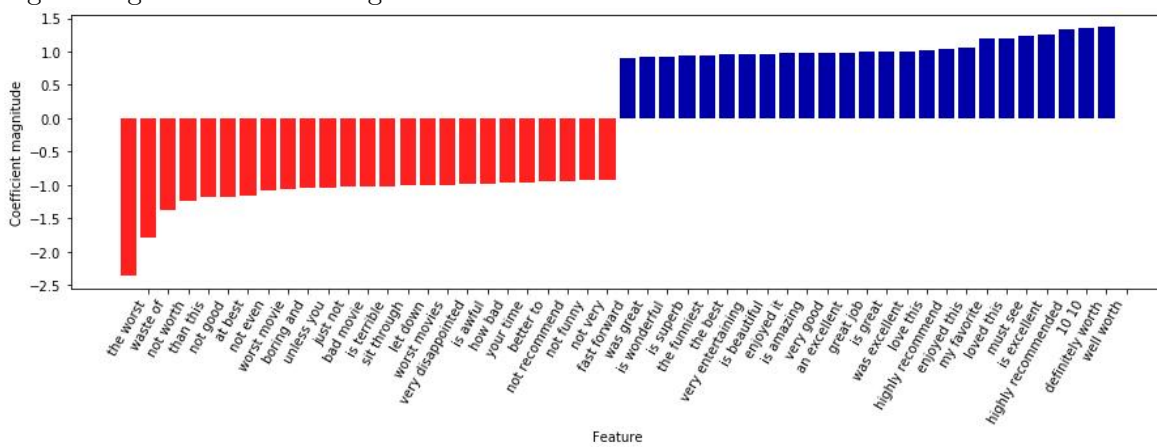
$\frac{\Delta J}{\Delta \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_\theta x_i - y_i) x_j^{(i)}$

Repeat

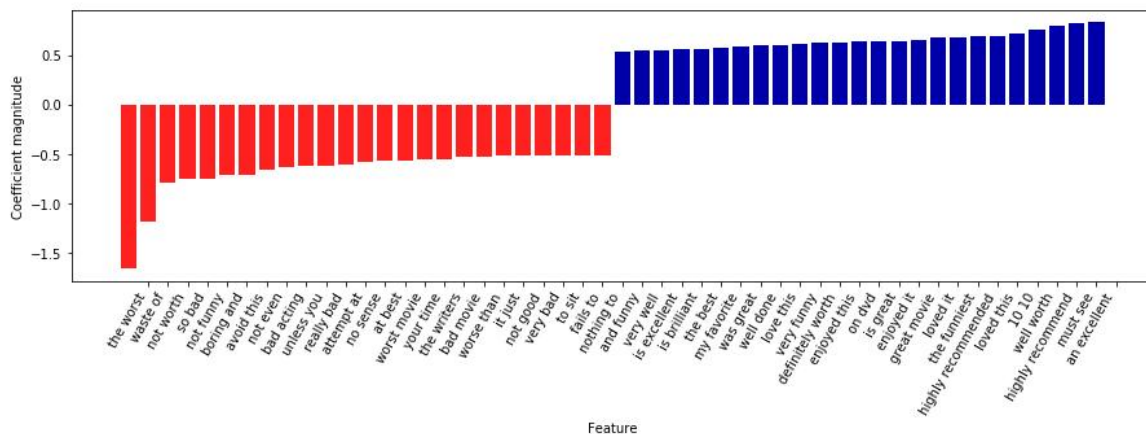
$\theta_j = \theta_j - \alpha \sum_{i=1}^m (h_\theta x_i - y_i) x_j^{(i)}$

$\theta = [\theta_0, \dots, \theta_n]$

Logistic regression for training data



Logistic regression for testing data



Parameters are used to find the best regression model which generalizes well i.e. works well with training as well as testing data sets by minimizing the mean squared error.