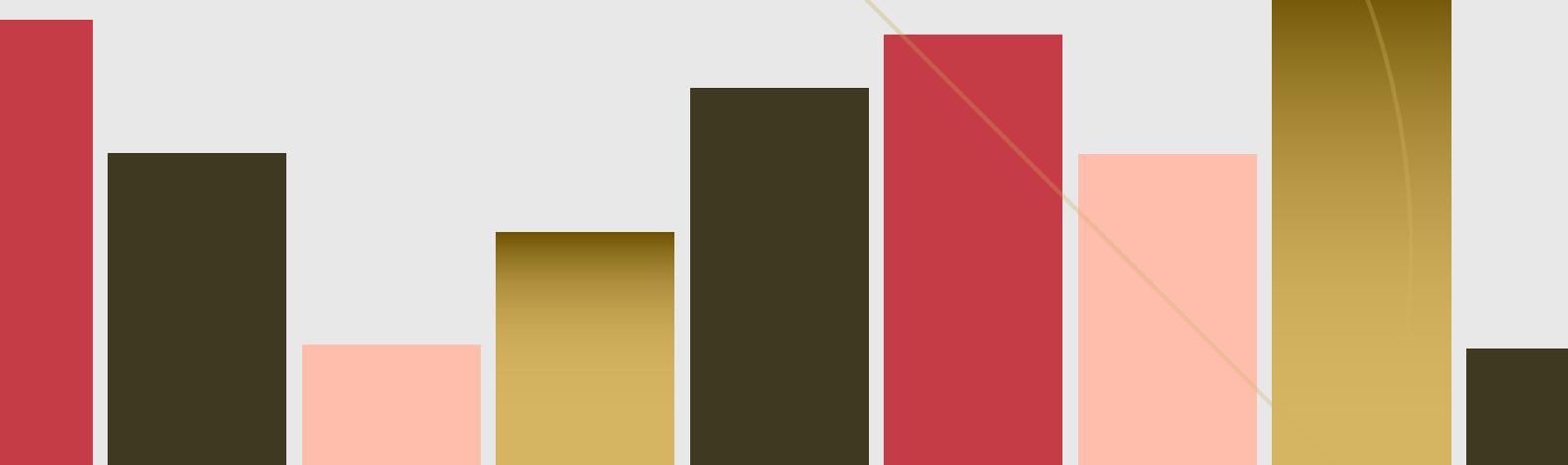


Natural Language
Processing

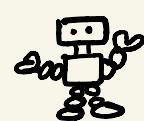
deeplearning.ai



Start date: 09/03



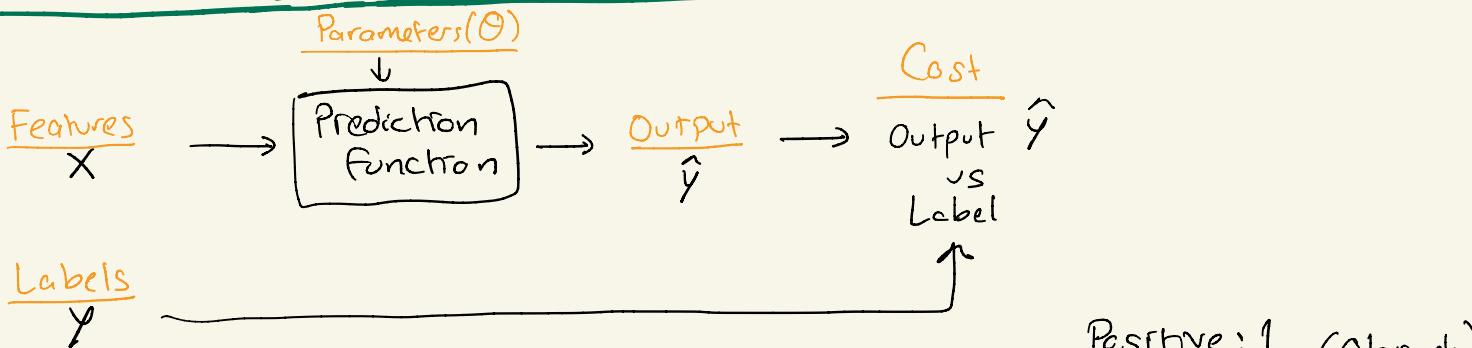
Natural Language Processing



Week 1

Sentiment Analysis with Logistic Regression

COURSE 1: NLP with
Classification and
Vector Spaces



- Logistic regression - sentiment analysis
 - ↳ 1) Process raw tweets
 - 2) Extract useful features
 - 3) Train logistic regression classifier (minimize the cost)

Positive: 1 (Obviously)
Negative: 0

- Represent a text as a vector (in dimension " V : vocabulary")

SLOW

- 1) Build a vocabulary (To encode any text as an array of numbers)
- 2) Vocabulary, $V \Rightarrow$ The list of **unique** words from all text (tweets in this case)
- 3) **Feature extraction:** Assigning 1 and 0 depending on the appearance of every word from the vocabulary above appears in the tweet
 - It appears: Assign 1 to that feature
 - If not: Assign 0 to that feature

- With this representation, a logistic regression model will have to learn $n+1$ parameters

Sparse representation: a small relative number of non-zero values

This representation of the tweet will have a number of features equal to the size of the entire vocabulary above.

- 4) Generate counts → features into the logistic regression classifier

↳ the number of times that a word appears in $+/-$ class

- Get the positive frequency in any word in the vocabulary, count the times as it appears in the positive tweets.
- ↳ Same thing applies to negative frequency.

freqs: dictionary mapping from (word, class) to frequency

$$X_m = [1, \sum_w \text{freqs}(w, 1), \sum_w \text{freqs}(w, 0)]$$

⇒ Represents a tweet as a vector of three

↓ features of tweet m ↓ bras
 sum of positive frequencies sum of negative frequencies

Preprocessing

- Stemming
- &
- Stop words

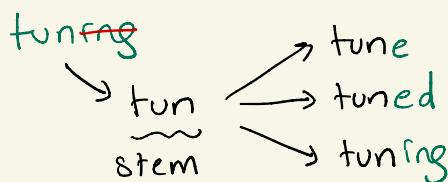
To process a tweet:

Preprocessing steps in NLP

- Tokenizing the string
- Lowercasing
- Removing stop words & punctuation
- Stemming

stop words

- 1) Remove all the words that don't add meaning to the tweets
*Stop words and punctuation marks
- 2) Every word that also appears in stop words and punc. marks lists should be eliminated. → Sometimes punctuations shouldn't be eliminated, evaluate case-by-case basis
- 3) **Stemming**: transforming any word to its base stem
the set of characters that are used to construct the word and its derivatives



- 4) To reduce vocabulary even further, lowercase every one of the words.

Generating matrix

After extracting features using a frequencies dictionary mapping

⇒ matrix X : m rows 3 columns : Every row would contain the features for each one of the tweets

$$\begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \\ \vdots & \vdots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} \end{bmatrix}$$

$\text{freqs} = \text{build_freqs}(\text{tweets}, \text{labels}) \Rightarrow \text{build frequencies dictionary}$

$X = \text{np.zeros}((m, 3)) \Rightarrow \text{initialize matrix } X$

for i in range(m):

$p\text{-tweet} = \text{process_tweet}(\text{tweets}[i]) \Rightarrow$ stemming, deleting stop words, URLs, handles, and lowercasing

$X[i, :] = \text{extract_features}(p\text{-tweet}, \text{freqs})$

• **NLTK package** ⇒

(natural language toolkit)

(nltk.org)

• tokenize and tag

• stemming

• parsing → display a parse tree

• semantic reasoning

• identifying named entities

```

import nltk
from nltk.corpus import twitter_samples # sample Twitter dataset from NLTK
import matplotlib.pyplot as plt
import random
nltk.download('twitter_samples') # pseudo-random number generator
all_positive_tweets = twitter_samples.strings('positive_tweets.json')
all_negative_tweets = twitter_samples.strings('negative_tweets.json')
print('Number of positive tweets:', len(all_positive_tweets))
print('Number of negative tweets:', len(all_negative_tweets))
print('\nThe type of all-positive-tweets is:', type(all_positive_tweets))
print('The type of a tweet entry is:', type(all_negative_tweets[0]))
# data is stored in a list, individual tweets are stored as strings.

```

Pie chart for visualizing the data

```

frg = plt.figure(figsize=(5,5))
labels = 'Positives', 'Negatives'
sizes = [len(all_positive_tweets), len(all_negative_tweets)] # sizes for each slide
plt.pie(sizes, labels=labels, autopct='%.1f%%', shadow=True,
        startangle=90)
plt.axis('equal')
plt.show()
print('1033[92m' + all_positive_tweets[random.randint(0,5000)]) # green
print('1033[91m' + all_negative_tweets[random.randint(0,5000)]) # red

```

nltk.download('stopwords')

```

import re # library for regular expression operations
import string # for string operations
from nltk.corpus import stopwords # module for stop words
from nltk.stem import PorterStemmer # module for stemming
from nltk.tokenize import TweetTokenizer # module for tokenizing strings

```

Removing hashtags, retweet marker, hyperlinks

```

tweet = all_positive_tweets[2277] # choosing a sample
tweet2 = re.sub(r'^RT[\s]+', '', tweet) # remove 'RT'
tweet2 = re.sub(r'https?:\/\/.*[\r\n]*', '', tweet2)
# remove hyperlinks
tweet2 = re.sub(r'#[^\s]', '', tweet2)
# remove only hashtag sign
print(tweet2)

```



Tokenizing the string • tokenize module - NLTK

tokenize: split the strings into individual words without blanks or tabs
substrings

tokenizer = TweetTokenizer (preserve_case=False, strip_handles=True, reduce_len=True)

↳ # instantiate tokenizer class

- By default, it is set to True.
- If set to False, the tokenizer will downcase everything.

tweet_tokens = tokenizer.tokenize(tweet2)

```
print()
print('Tokenized string:')
print(tweet_tokens)
```

Remove stop words and punctuations

stopwords_english = stopwords.words('english') # import the english stop words 1st from NLTK

print('Stop words\n')
print(stopwords_english)
print('\n Punctuation\n')
print(string.punctuation)

Clean up

```
tweets_clean = []
for word in tweet_tokens: # go through every word in the list
    if (word not in stopwords_english and # removing stop words
        word not in string.punctuation): # removing punctuations
        tweets_clean.append(word)
print('removed: ') # let's see what our cleaned tweet looks like
print(tweets_clean)
```

Stemming → NLTK has different modules for stemming. Martin Porter's algorithm

stemmer = PorterStemmer() # instantiate stemming class

tweets_stem = [] # create an empty list to store the stems

for word in tweets_clean:
 stem_word = stemmer.stem(word) # stemming word
 tweets_stem.append(stem_word) # append to the list
print("stemmed words:")
print('tweets_stem')

Porter Stemmer module uses

Building and Visualizing Word Frequencies

```
import nltk  
from nltk.corpus import twitter_samples  
import matplotlib.pyplot as plt  
import numpy as np  
nltk.download('stopwords')  
from utils import process_tweet, build_freqs
```

```
all_positive_tweets = twitter_samples.strings('positive_tweets.json')  
all_negative_tweets = twitter_samples.strings('negative_tweets.json')  
tweets = all_positive_tweets + all_negative_tweets #concatenate  
print("Number of tweets:", len(tweets))
```

* labels = np.append(np.ones((len(all_positive_tweets))), np.zeros((len(all_negative_tweets))))
↳ #numpy array representing labels of the tweets

word frequency dictionary:

```
def build_freqs(tweets, ys):
```

↳ an mx1 array with the sentiment label of each tweet

```
yslist = np.squeeze(ys).tolist()
```

↳ #converting np array into list because zip

function needs an iterable

squeeze is necessary otherwise the list ends up with one element.

if ys is already a list, no need to do this.

```
freqs = {} # start with empty dictionary
```

```
for y, tweet in zip(yslist, tweets):
```

```
    for word in process_tweet(tweet):
```

```
        pair = (word, y)
```

alternative:

```
freqs[pair] = freqs.get(pair, 0) + 1
```

{ if pair in freqs:

```
    freqs[pair] += 1
```

else:

```
    freqs[pair] = 1
```

```
return freqs
```

Each key is a 2-element tuple containing a (word, y) pair:

- word: an element in a processed tweet
- y: an integer representing the corpus $\Rightarrow 1$ for positive, 0 for negative
- value associated with a key is the number of times that word appears in the specified corpus

('shame', 0.0): 19

2 helper functions:

process_tweet()

↳ clean, tokenize, removes stopwords, stemming

build_freqs()

↳ Counts how often a word in the corpus was associated with positive label '1' and negative label '0'.

Reminder of < dictionaries > in Python

Dictionaries:

- Mutable { collection }
- Indexed
- Stores items as key-value pairs and uses hash tables underneath to allow lookups.

• Declared using curly brackets { }

↳ dictionary = { 'key1': 1, 'key2': 2 }

Adding a new entry:

- With square brackets []

↳ dictionary['key3'] = -5

Accessing values and lookup keys:

- 2 methods:

1) square brackets

print(dictionary['key1'])

* Works if the lookup key is in the dictionary

2) get()

* Allows to set a default value if the key does not exist

print(dictionary.get('key7', -1))

```

freqs = build_freqs(tweets, labels)
print(f'type(freqs) = {type(freqs)}') # data type
print(f'len(freqs) = {len(freqs)}') # length of the dictionary

```

Table of word counts

```

keys = ['-', '+', '_'] # select words to appear in the report.
data = [] # list representing our table of word counts
for word in keys:
    pos = 0 # initialize + and - counts [word], <pos-count>, <neg-count>
    neg = 0
    if (word, 1) in freqs: # retrieve number of positive counts
        pos = freqs[(word, 1)]
    if (word, 0) in freqs: # retrieve number of negative counts
        negative = freqs[(word, 0)]
    data.append([word, pos, neg]) # append the word counts to the table

```

each element has a sublist:

data

Visualizing with scatter plot

wide discrepancies between positive and negative values \Rightarrow logarithmic scale rather than raw counts

```

fig, ax = plt.subplots(figsize=(8, 8))
x = np.log([x[1]+1 for x in data]) # adding 1 is to avoid log[0]
y = np.log([x[2]+1 for x in data])
ax.scatter(x, y)
plt.xlabel("Log positive count")
plt.ylabel("Log negative count")
for i in range(0, len(data)):
    ax.annotate(data[:, 0], (x[i], y[i]), fontsize=12)
    ax.plot([0, 9], [0, 9], color='red')
plt.show()

```

Logistic Regression



- Sigmoid function which outputs a probability between zero and one

$$h(x^{(i)}, \theta) = \frac{1}{1 + e^{-\theta^T x^{(i)}}}$$

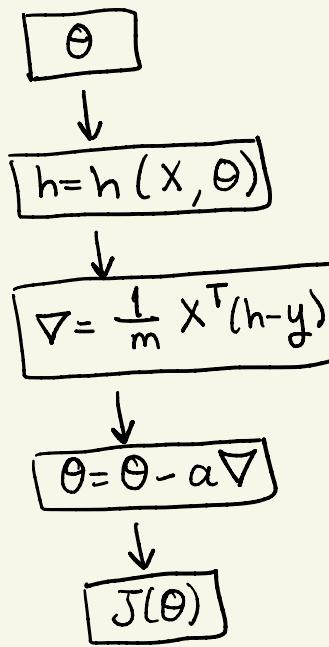
- For classification:

Threshold = 0.5

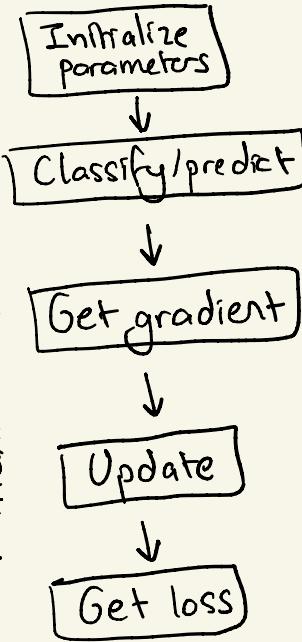
↳ Dot product between θ^T and x equals to zero at this value

$$\theta^T x^{(i)}$$

Training LR



Stop-parameter
α
to
• Goes until
or maximum number of iterations.



Gradient descent

code:
z = input (scalar/array)
 $h = 1 / (1 + np.exp(-z))$

Testing LR

$$\text{pred} = h(x_{\text{val}}, \theta) \geq 0.5 \Rightarrow \begin{bmatrix} 0.3 \\ 0.8 \\ 0.5 \\ \vdots \\ \text{him} \end{bmatrix} \geq 0.5 = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ \text{predm} \end{bmatrix}$$

$$\sum_{i=1}^m \frac{(\text{pred}^{(i)} == y_{\text{val}}^{(i)})}{m} \Rightarrow \begin{bmatrix} 0 \\ 1 \\ \vdots \\ \text{predm} \end{bmatrix} == \begin{bmatrix} 0 \\ 0 \\ \vdots \\ y_{\text{valm}} \end{bmatrix} \Rightarrow \begin{bmatrix} 1 \\ 0 \\ \vdots \\ \text{predm} == y_{\text{valm}} \end{bmatrix}$$

↳ accuracy: an estimate of the times that logistic regression will correctly work on unseen data

Cost Function for Logistic Regression

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log h(x^{(i)}, \theta) + (1-y^{(i)}) \log (1-h(x^{(i)}, \theta))]$$

number of the training examples y⁽ⁱ⁾ h(x⁽ⁱ⁾, θ) result vector
 0 any 0
 1 0.99 ~0
 1 ~0 -inf

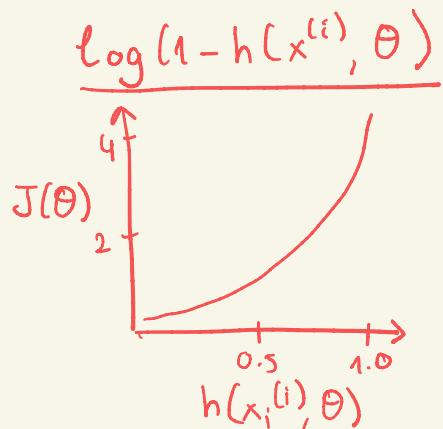
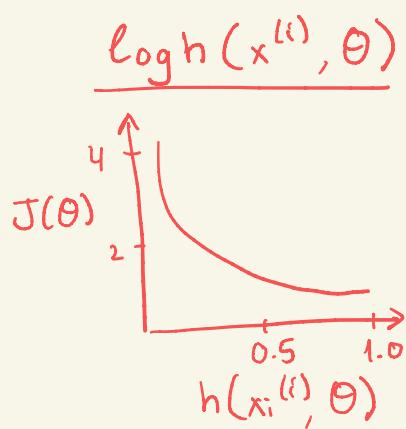
- ensures that cost will always be a positive number

★ this term is relevant if the label is 1. When prediction is close to the label value, loss is small. If they disagree, the overall cost goes up.

$$(1-y^{(i)}) \log (1-h(x^{(i)}, \theta)) \rightarrow$$

y ⁽ⁱ⁾	h(x ⁽ⁱ⁾ , θ)	result vector
1	any	0
1	0.01	~0
0	~1	-inf

★ this term is relevant when the label is 0.



• Visualizing tweets and logistic regression model •

```
import nltk
from os import getcwd
import pandas as pd
from nltk.corpus import twitter_samples
import matplotlib.pyplot as plt
import numpy as np

# functions defined before
from utils import process_tweet, build_freqs
all_positive_tweets = twitter_samples.strings('positive_tweets.json')
all_negative_tweets = twitter_samples.strings('negative_tweets.json')

tweets = all_positive_tweets + all_negative_tweets
labels = np.append(np.ones((len(all_positive_tweets), 1)),
                   np.zeros((len(all_negative_tweets), 1)), axis=0)

train_pos = all_positive_tweets[:4000]    # splitting dataset into
train_neg = all_negative_tweets[:4000]      # two: 1) training
                                             2) testing

train_x = train_pos + train_neg
print("Number of tweets: ", len(train_x))

data = pd.read_csv('logistic_features.csv'); # pre-prepared
                                             # CSV containing
                                             # features
data.head(10) # printing first 10 columns

X = data[['bias', 'positive', 'negative']].values # get numerical values
y = data['sentiment'].values;
print(X.shape)

theta = [7e-08, 0.005239, -0.005517] # calling a pretrained
                                         # LR model

# split our feature space into 2 parts
# We have a 3D feature space → bias, pos. value, neg. value
# If we ignore bias, we can plot each tweet in a cartesian
# plane.

fig, ax = plt.subplots(figsize=(8, 8))
colors = ['red', 'green']
ax.scatter(X[:, 1], X[:, 2], c=[colors[int(k)] for k in y], s=0.1)
plt.xlabel('Positive')
plt.ylabel('Negative')
```



plot the model alongside the data (just for visualization)

draw a gray line to show the cutoff between the pos. and neg. regions.

gray line: $z = \theta * x = 0$

$$\hookrightarrow z = \theta * x = 0$$

$$x = [1, \text{pos}, \text{neg}]$$

$$z(\theta, x) = \theta_0 + \theta_1 * \text{pos} + \theta_2 * \text{neg} = 0$$

$$\text{neg} = (-\theta_0 - \theta_1 * \text{pos}) / \theta_2$$

there will be red and green lines that point in the direction of the corresponding sentiment.

- they are calculated using perpendicular (b) line to the separation line
- lines should point in the same direction as the derivative of the Logit function. $\rightarrow \text{direction} = \text{pos} * \theta_2 / \theta_1$

$$\begin{aligned} \# f(\text{pos}, \text{neg}, w) &= w_0 + w_1 * \text{pos} + w_2 * \text{neg} = 0 \\ \# s(\text{pos}, w) &= (-w_0 - w_1 * \text{pos}) / w_2 \end{aligned} \quad \left. \begin{array}{l} \text{equation for} \\ \text{the separation} \\ \text{plane} \end{array} \right\}$$

def neg(theta, pos):

$$\text{return } (-\theta[0] - \text{pos} * \theta[1]) / \theta[2]$$

$$\# df(\text{pos}, w) = \text{pos} * w_2 / w_1 \quad \rightarrow \text{equation for the directions of the sentiment change}$$

def direction(theta, pos):

$$\text{return } \text{pos} * \theta[2] / \theta[1]$$

fig, ax = plt.subplots(figsize=(8, 8))

colors = ['red', 'green']

ax.scatter(x[:, 1], x[:, 2], c=[colors[int(k)] for k in Y], s=0.1)

plt.xlabel("Positive")

plt.ylabel("Negative")

max_pos = np.max(x[:, 1]) # represent LR model in this chart

offset = 500

ax.plot([0, max_pos], [neg(theta, 0), neg(theta, max_pos)], color='gray')

ax.arrow(offset, neg(theta, offset), offset, direction(theta, offset), head_width=500, head_length=500, fc='g', ec='g')

ax.arrow(offset, neg(theta, offset), -offset, -direction(theta, offset), head_width=500, head_length=500, fc='r', ec='r')

plt.show()

PROGRAMMING ASSIGNMENT - Week 1

```
import nltk
from os import getcwd
nltk.download('twitter_samples')
nltk.download('stopwords')
import numpy as np
import pandas as pd
from nltk.corpus import twitter_samples
from utils import process_tweet, build_freqs
all_positive_tweets = twitter_samples.strings('positive_tweets.json')
all_negative_tweets = twitter_samples.strings('negative_tweets.json')
test_pos = all_positive_tweets[4000:]
train_pos = all_positive_tweets[:4000]
test_neg = all_negative_tweets[4000:]
train_neg = all_negative_tweets[:4000]
train_x = train_pos + train_neg
test_x = test_pos + test_neg
train_y = np.append(np.ones((len(train_pos), 1)), np.zeros((len(train_neg), 1)),
axis=0)
test_y = np.append(np.ones((len(test_pos), 1)), np.zeros((len(test_neg), 1)),
axis=0)
print("train_y.shape = " + str(train_y.shape))
print("test_y.shape = " + str(test_y.shape))
freqs = build_freqs(train_x, train_y)  $\Rightarrow$  creating frequency dictionary
    ↳ for y, tweet in zip(ys, tweets):
        for word in process_tweet(tweet):
            pair = (word, y)
            if pair in freqs:
                freqs[pair] += 1
            else:
                freqs[pair] = 1
print("type(freqs) = " + str(type(freqs)))
print("len(freqs) = " + str(len(freqs.keys())))

```

Part 1: Logistic regression

Part 1.1: Sigmoid

$$h(z) = \frac{1}{1 + \exp^{-z}}$$

def sigmoid(z): $\Rightarrow z$: logits

$$h = 1 / (1 + np.exp(-z))$$

return h

Part 1.2: Cost function and gradient

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h(z(\theta)^{(i)})) + (1-y^{(i)}) \log(1-h(z(\theta)^{(i)}))$$

m: training examples

$y^{(i)}$: actual label of the i-th training example

$h(z(\theta)^{(i)})$: model's prediction for the i-th training example

$$\nabla_{\theta_j} J(\theta) = \underbrace{\frac{1}{m} \sum_{i=1}^m}_{\text{the gradient of}} (h^{(i)} - y^{(i)}) x_j \quad \begin{array}{l} \text{weight vector: } \theta \\ \text{Index across all } \\ 'm' \text{ training } \\ \text{examples} \end{array}$$

the gradient of
the cost function

$$\theta_j = \theta_j - \alpha \times \nabla_{\theta_j} J(\theta)$$

index of the
weight θ_j : j

feature associated
w/ weight θ_j : x_j

def gradientDescent(x, y, theta, alpha, num_iters):

"Input:

x: matrix of features $\Rightarrow m, n+1$

y: corresponding labels of matrix x $\Rightarrow m, 1$

theta: weight vector of dimension $\Rightarrow n+1, 1$

alpha: learning rate

num_iters: number of iterations you want to train your model for

Output:

J: the final cost

theta: final weight vector

m = x.shape[0]

for i in range(0, num_iters):

z = np.dot(x, theta)

h = 1 / (1 + np.exp(-z))

J = -(float(1)/m) * (np.dot(y.transpose(), np.log(h)) +

np.dot((1-y).transpose(), np.log(1-h)))

$$J = -\frac{1}{m} \times (y^T \cdot \log(h) + (1-y)^T \cdot \log(1-h))$$

$$\theta = \theta - \frac{\alpha}{m} \times (x^T \cdot (h-y))$$

```
theta = theta - ((alpha/m) * np.dot(np.transpose(x), (h-y)))
```

```
J = float(J)
```

```
return J, theta
```

Part 2: Extracting features

```
def extract_features(tweet, freqs):
```

```
    word_L = process_tweet(tweet)
```

```
    x = np.zeros((1,3))
```

```
    x[0,0] = 1
```

```
    for word in word_L:
```

```
        x[0,1] += freqs.get((word, 1.0), 0)
```

```
        x[0,2] += freqs.get((word, 0.0), 0)
```

```
    assert (x.shape == (1,3))
```

```
    return x
```

Part 3: Tranning your model

```
X = np.zeros((len(tram_x), 3))
```

```
for i in range(len(tram_x)):
```

```
    X[i, :] = extract_features(tram_x[i], freqs)
```

stack the features for
all training examples
into a matrix X

```
y = train_y
```

```
J, theta = gradient Descent(X, y, np.zeros((3,1)), 1e-9, 1500)
```

Part 4: Test your logistic regression

```
def predict_tweet(tweet, freqs, theta):
```

```
"y pred = sigmoid(x · θ)"
```

$\tilde{(3,1)}$ vector of weights

```
x = extract_features(tweet, freqs)
```

```
y - pred = sigmoid(np.dot(x, theta))
```

```
return y - pred
```

Part 5: Check performance

```
def test_logistic_regression(test_x, test_y, freqs, theta):
```

" Inputs:

test_x: a list of tweets

test_y: (m, 1) vector w/ the corresponding labels

freqs: a dictionary w/ the frequency of each pair

theta: weight vector of dimension (3,1)

Output:

accuracy

"

y-hat = []

for tweet in test_x:

y-pred = predict_tweet(tweet, freqs, theta)

if y-pred > 0.5:

y-hat.append(1.0)

else:

y-hat.append(0)

accuracy = (y-hat == np.squeeze(test_y)).sum() / len(test_x)

return accuracy

Week 2

Sentiment Analysis with Naive Bayes

Bayes' Rule

Conditional probabilities

probability of B ,
given A happened

looking at the elements of
set A , the chance that
one also belongs to set B

$$P(\text{positive} | \text{"happy"}) = \frac{P(\text{positive} \cap \text{"happy"})}{P(\text{"happy"})}$$

$$P(\text{"happy"} | \text{positive}) = \frac{P(\text{"happy"} \cap \text{positive})}{P(\text{positive})}$$

$$\rightarrow P(\text{positive} | \text{"happy"}) = P(\text{"happy"} | \text{positive}) \times \frac{P(\text{positive})}{P(\text{"happy"})}$$

$$\rightarrow P(X|Y) = P(Y|X) \times \frac{P(X)}{P(Y)} \quad | \text{Bayes' rule}$$

Naive Bayes for Sentiment Analysis

naive: model makes the assumption that the features used for classification are all independent.

Positive tweets corpus

- I am happy because
- I am learning NLP
- I am happy, not sad.

Negative tweets corpus

- I am sad, I am not learning NLP
- I am sad, not happy

word	pos	neg
I	3	3
am	3	3
happy	2	1
because	1	0
learning	1	1
NLP	1	1
sad	1	2
not	1	2
Total	13	13

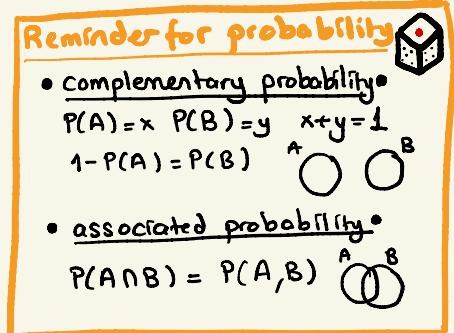
conditional probability
→

word	pos	neg
I	0.24	0.24
am	0.24	0.24
happy	0.15	0.08
because	0.08	0
learning	0.08	0.08
NLP	0.08	0.08
sad	0.08	0.15
not	0.08	0.15
Sum	1	1

$$P(I|Pos) = 3/13$$
$$P(I|Neg) = 3/13$$

$$\prod_{i=1}^m \frac{P(w_i|Pos)}{P(w_i|Neg)}$$

→ If > 1 :
positive sentence



Laplacian smoothing:

- a technique to avoid the probabilities being zero

$$P(w_i | \text{class}) = \frac{\text{freq}(w_i, \text{class})}{N_{\text{class}}}$$

- $\text{class} \in \{\text{Positive, Negative}\}$
- $N_{\text{class}} = \text{frequency of all words in class}$

- $V = \text{number of unique words in vocabulary}$

- $\Rightarrow 1$: avoids the probability being zero
 ↳ it adds a new term to all the frequencies that is not correctly normalized by N_{class}

additive
smoothing

- ↳ to account for this, add a new term in the denominator V .

- ↳ all probabilities will sum to 1

Reasons for using Laplacian:

- To avoid $P(w_i | \text{class}) = 0$

Log likelihood:

- conditional probability of each word:

$$\text{ratio}(w_i) = \frac{P(w_i | \text{Pos})}{P(w_i | \text{Neg})}$$

↳ ratios are essential in Naïve Bayes' for binary classification

$$\prod_{i=1}^m \frac{P(w_i | \text{pos})}{P(w_i | \text{neg})} \Rightarrow \text{likelihood}$$

$$\frac{P(\text{pos})}{P(\text{neg})}$$

prior ratio:

- taking a ratio between the negative and positive tweets
- important for unbalanced datasets - otherwise gen. 1

$$\frac{P(\text{pos})}{P(\text{neg})} \cdot \prod_{i=1}^m \frac{P(w_i | \text{pos})}{P(w_i | \text{neg})} \Rightarrow \text{Full Naïve Bayes' formula}$$

Multiplication of many numbers between zero and one

→ risk of numerical underflow

to solve:

$$\log\left(\frac{P(\text{pos})}{P(\text{neg})} \prod_{i=1}^m \frac{P(w_i | \text{pos})}{P(w_i | \text{neg})}\right)$$

$$1) \log(a * b) = \log(a) + \log(b)$$

$$2) \log \frac{P(\text{pos})}{P(\text{neg})} + \sum_{i=1}^m \log \frac{P(w_i | \text{pos})}{P(w_i | \text{neg})}$$

log prior

$$\log \text{likelihood}$$

$$\Rightarrow \sum_{i=1}^m \lambda(w_i)$$

$$\lambda(w) = \log \frac{P(w_i | \text{pos})}{P(w_i | \text{neg})}$$

Lambda

Trapping Naïve Bayes

✓ No gradient descent - Just counting the frequencies of the words in a corpus

! Steps of a supervised ML project (NLP):

Step 1: Collect and annotate corpus

Step 2: Preprocess

- Lowercase
- Remove punctuations, urls, nones
- Remove stop words
- Stemming
- Tokenize sentences

Step 3: 1) Word count: $\text{freq}(w, \text{class})$

$$2) P(w_i | \text{class}) \rightarrow \frac{\text{freq}(w_i, \text{class}) + 1}{N_{\text{class}} + V_{\text{class}}}$$

Step 4: Get lambda

$$\lambda(w) = \log \frac{P(w | \text{pos})}{P(w | \text{neg})}$$

Step 5: Get the log prior

$$\text{Logprior} = \frac{D_{\text{pos}} (\text{number of + tweets})}{D_{\text{neg}} (\text{number of - tweets})}$$

★ If dataset is balanced, $D_{\text{pos}} = D_{\text{neg}}$ and $\text{logprior} = 0$

Visualizing Naive Bayes

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```

```
from utils import confidence_ellipse
```

- calculate the likelihoods for each tweet

- plot the samples:

```
fig, ax = plt.subplots(figsize=(8,8))
colors = ['red', 'green']
```

```
ax.scatter(data.positive, data.negative, c=[colors[int(k)] for k in data.sentiment],
           s=0.1, marker='*')
```

```
plt.xlim(-250, 0)
plt.ylim(-250, 0)
plt.xlabel("Positive")
plt.ylabel("Negative")
```

- confidence ellipse:

* a way to visualize a 2D random variable

```
fig, ax = plt.subplots(figsize=(8,8))
```

```
colors = ['red', 'green']
```

```
ax.scatter(data.positive, data.negative, c=[colors[int(k)] for k in data.sentiment],
           s=0.1, marker='*')
```

```
plt.xlim(-200, 40)
```

```
plt.ylim(-200, 40)
```

```
plt.xlabel("Positive")
```

```
plt.ylabel("Negative")
```

```
data_pos = data[data.sentiment == 1]
```

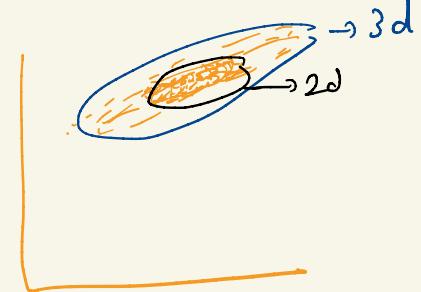
```
data_neg = data[data.sentiment == 0]
```

↳ precalculated
likelihoods

2 sd { Confidence_ellipse(data_pos.positive, data_pos.negative, ax, n_std=2, edgecolor='black', label=r'\$2/\sigma_{\mu}\$')
confidence_ellipse(data_neg.positive, data_neg.negative, ax, n_std=2, edgecolor='orange')

3 d { Confidence_ellipse(data_pos.positive, data_pos.negative, ax, n_std=3, edgecolor='black', linestyle=':', label=r'\$3/\sigma_{\mu}\$')
confidence_ellipse(data_neg.positive, data_neg.negative, ax, n_std=3, edgecolor='orange', linestyle=':')

```
ax.legend()
plt.show()
```



Testing Naive Bayes

Tweet t: I passed the NLP interview

↳ sum of the lambda score of each word + log prior > 0 = positive sentiment

$X_{val}, Y_{val}, \lambda, \log \text{prior}$

↳ score = predict($X_{val}, \lambda, \log \text{prior}$)

pred = score > 0

$$\frac{1}{m} \sum_{i=1}^m (\text{pred}_i == Y_{val_i}) \Rightarrow \text{accuracy}$$

Applications of Naive Bayes

Reminder:

$$P(\text{pos} | \text{tweet}) \approx P(\text{pos}) P(\text{tweet} | \text{pos})$$

$$P(\text{neg} | \text{tweet}) \approx P(\text{neg}) P(\text{tweet} | \text{neg})$$

$$\frac{P(\text{pos} | \text{tweet})}{P(\text{neg} | \text{tweet})} = \frac{P(\text{pos})}{P(\text{neg})} \prod_{i=1}^m \frac{P(w_i | \text{pos})}{P(w_i | \text{neg})}$$

- author identification
- spam detection
- information retrieval
- word disambiguation

- One of the earliest use of Naive Bayes was filtering relevant and irrelevant documents in a database. \Rightarrow We only need to calculate likelihood of whether the doc is irr. or rel.

Information retrieval:

$$P(\text{document}_k | \text{query}) \propto \prod_{i=0}^{|\text{query}|} P(\text{query}_i | \text{document}_k)$$

retrieve document if $P(\text{document}_k | \text{query}) > \text{threshold}$

Naive Bayes Assumptions

1) Independence

- Assume independence between the predictors or features associated with each class

↳ Words in a text are independent of one another

↳ Under/over estimates the conditional probabilities of individual words

2) Relative frequency in corpus

- Naive Bayes relies on the distribution of the training dataset

Error Analysis

- Causes of errors :
- 1) Semantic meaning loss in the pre-processing step
Due to ↳ Removing punctuation and stop words
 - 2) Word order affecting the meaning of the sentence
 - 3) Quirks of human language
↳ Adversarial attacks

Week 2: Assignment

def count_tweets (result, tweets, ys):

" Input:
 result: a dictionary that will be used to map each pair to its frequency
 tweets: a list of tweets
 ys: a list corresponding to the sentiment of each tweet (either 0 or 1)

Output:
 result: a dictionary mapping each pair to its frequency

"

pair = (word, y)

if pair in result:

 result [pair] += 1

else:

 result [pair] = 1

output: {('happi', 1): 1, ('trick', 0): 1}

return result

freqs = count_tweets ({}, train-x, train-y) \Rightarrow building freqs dictionary

def train-naive-bayes (freqs, train-x, train-y):

log likelihood = {}

logprior = 0

vocab = set ([pair[0] for pair in freqs.keys()])

V = len(vocab)

N_pos = N_neg = 0

for pair in freqs.keys():

 if pair[1] > 0:

 N_pos += 1

 else:

 N_neg += 1



```
<indent> D = len(train-y)
```

```
D_pos = len(list(filter(lambda x: x > 0, train-y)))
```

```
D_neg = len(list(filter(lambda x: x <= 0, train-y)))
```

```
logprior = np.log(D_pos) - np.log(D_neg)
```

```
for word in vocab:
```

```
    freq_pos = lookup(freqs, word, 1)
```

```
    freq_neg = lookup(freqs, word, 0)
```

```
p_w_pos = (freq_pos + 1) / (N_pos + V)
```

```
p_w_neg = (freq_neg + 1) / (N_neg + V)
```

```
loglikelihood[word] = np.log(p_w_pos / p_w_neg)
```

```
return logprior, loglikelihood
```

```
logprior, loglikelihood = train_naive_bayes(freqs, train_x, train_y)
```

```
def naive_bayes_predict(tweet, logprior, loglikelihood):
```

```
word_1 = process_tweet(tweet)
```

```
p = 0
```

```
p += logprior
```

```
for word in word_1:
```

```
    if word in loglikelihood:
```

```
        p += loglikelihood[word]
```

```
return
```

```
def test_naive_bayes(test_x, test_y, logprior, loglikelihood):
```

```
accuracy = 0
```

```
y_hats = []
```

```
for tweet in test_x:
```

```
    if naive_bayes_predict(tweet, logprior, loglikelihood) > 0:
```

```
        y_hat_i = 1
```

```
    else:
```

```
        y_hat_i = 0
```

```
error = np.mean(np.absolute(y_hats - test_y))
```

```
accuracy = 1 - error
```

```
return
```

$$p = \logprior + \sum_i^N (\loglikelihood_i)$$

def get_ratio(freqs, word): → defining the level of positiveness or negativeness
 pos-neg-ratio = {'positive': 0, 'negative': 0, 'ratio': 0.0}
 pos-neg-ratio['positive'] = lookup(freqs, word, 1)
 pos-neg-ratio['negative'] = lookup(freqs, word, 0)
 pos-neg-ratio['ratio'] = ((pos-neg-ratio['positive']+1)/(pos-neg-ratio['negative']+1))

def get_words_by_threshold(freqs, label, threshold):
 word_list = {}
 for key in freqs.keys():
 word, _ = key
 pos-neg-ratio = get_ratio(freqs, word)
 if label == 1 and pos-neg-ratio >= threshold:
 word_list[word] = pos-neg-ratio
 elif label == 0 and pos-neg-ratio <= threshold:
 word_list[word] = pos-neg-ratio
 return word_list

0 for neg, 1 for pos threshold value
 ↑ ↑

get_words_by_threshold(freqs, label=↑, threshold=↑)

Week 3

Vector Space Models