

Contents

| | | |
|----------|---|-----------|
| 1 | Sentiment Analysis with Logistic Regression | 2 |
| 1.1 | Intro to Supervised ML and Sentiment Analysis | 2 |
| 1.1.1 | Representation of Text | 2 |
| 1.1.2 | Feature Extraction with Frequencies | 4 |
| 1.1.3 | Preprocessing | 4 |
| 1.2 | Logistic Regression | 6 |
| 1.2.1 | Learning Parameters | 6 |
| 1.2.2 | Assessing model generalization | 7 |
| 1.2.3 | Deriving Gradient Descent for Logistic Regression | 7 |
| 2 | Sentiment Analysis with Naive Bayes | 10 |
| 2.1 | Probability and Bayes Rule | 10 |
| 2.2 | Naive Bayes | 12 |
| 2.2.1 | Assumptions of Naive Bayes | 15 |
| 2.2.2 | Error Analysis | 15 |
| 3 | Vector Space Models | 16 |
| 3.1 | Motivating Vector Space Models | 16 |
| 3.2 | Word by Word and Word by Doc | 17 |
| 3.3 | Distance and Similarity Metrics | 18 |
| 3.4 | Principal Components Analysis | 20 |
| 4 | Machine Translation | 21 |
| 4.1 | Transforming Word Vectors | 21 |

1 Sentiment Analysis with Logistic Regression

1.1 Intro to Supervised ML and Sentiment Analysis

In supervised ML, we have input features X and a set of labels Y . To get the most accurate predictions, we try to minimize our *error rates* or *cost function* as much as possible: to do this, we'll run our prediction function which takes in parameters θ to map your input features to output labels \hat{Y} . The best mapping is achieved when the difference between the expected values Y and the predicted values \hat{Y} is minimized, which the cost function does by comparing how closely your output \hat{Y} is to your label Y . You can then update your parameters and repeat the whole process until your cost is minimized.



Figure 1: Overview of supervised machine learning.

How about the supervised ML classification task of sentiment analysis? Suppose we're given a tweet that says, "I'm happy because I'm learning NLP": and the objective in the task is to predict whether a tweet has a positive or negative sentiment. We'll do this by starting with a training set where tweets with a positive label have a label of unit value, and tweets with a negative sentiment have a label of zero. To get started building a logistic regression classifier that's capable of predicting sentiments of an arbitrary tweet, we first need to process the raw tweets in our training data set and extract useful features. Then, we will train our logistic regression classifier while minimizing the cost. Finally, we'll be able to make predictions.

1.1.1 Representation of Text

How to represent text as a vector In order to represent text as a vector, we need to first build a vocabulary. We define the vocabulary V as the *set* of unique words from your input data (e.g. your listing of tweets). To get this listing, we quite literally need to comb through all words from all input data and save every new word that appears in our search. To represent a tweet as a vector, we can use a one-hot encoding with our vocabulary: i.e. each tweet will be represented with a length $|V|$ vector where elements are binary-valued - a one indicates the word is in the tweet and a zero indicates the absence of a word in a tweet. We call this a *sparse* representation because the number of non-zero entries is relatively small when compared with the number of zero entries. Realize that if we are running a logistic regression, we would require learning $|V| + 1$ parameters which can be problematic for large vocabularies. If not prohibitive, it would make training models take excessive time and making predictions would be expensive.

Negative and positive frequencies Let's discuss how to generate counts which can be used as features in our logistic regression classifier. Specifically, given a word, we wish to keep track of the number of times that it shows up as the positive class. Given another word, we wish to track how many times that word shows up in the negative class. Using both these counts, we can then extract features and use those features in our logistic regression classifier. Suppose we have the following corpus of tweets:

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

Then our vocabulary is given by

| Vocabulary |
|------------|
| I |
| am |
| happy |
| because |
| learning |
| NLP |
| sad |
| not |

For this particular example of sentiment analysis, we only have two sentiments (i.e. two classes): one class is associated with a positive sentiment and the other with a negative sentiment. So, taking your corpus, you'd have a set of two tweets that belong to the positive class, and two tweets which belong to the negative class. Let's calculate the positive frequencies by examining the first two tweets:

| Vocabulary | PosFreq(1) |
|------------|------------|
| I | 3 |
| am | 3 |
| happy | 2 |
| because | 1 |
| learning | 1 |
| NLP | 1 |
| sad | 0 |
| not | 0 |

The same logic applies to getting negative frequencies. We can calculate these by examining our last two training examples.

| Vocabulary | NegFreq(0) |
|------------|------------|
| I | 3 |
| am | 3 |
| happy | 0 |
| because | 0 |
| learning | 1 |
| NLP | 1 |
| sad | 2 |
| not | 1 |

So, we can now have an entire table for our corpus, where for each entry in V we associate with it a scalar value $\text{PosFreq}(1)$ and another scalar value $\text{NegFreq}(0)$. In practice, we use a Python dictionary `freqs` mapping from `(word, class)` \rightsquigarrow frequency.

1.1.2 Feature Extraction with Frequencies

Whereas we previously learned to encode a tweet as a vector of length $|V|$, we will now use our frequency counts to represent each tweet as a vector of length equal to one plus the number of classes in our set of labels. This gives us a much faster speed for our logistic regression classifier. How can we do this, exactly? We represent each tweet as follows:

$$\underbrace{X_m}_{\text{Features of tweet } m} = \left[\underbrace{1}_{\text{Bias}}, \underbrace{\sum_w \text{freqs}(w, 1)}_{\text{Sum Pos. Frequencies}}, \underbrace{\sum_w \text{freqs}(w, 0)}_{\text{Sum Neg. Frequencies}} \right]$$

I.e. the first feature is a bias unit equal to unit value, the second is the sum of positive frequencies for every unique word on tweet m , and the third is the sum of negative frequencies for every unique word on the tweet. So, to extract the features for this *representation*, we only have to sum frequencies of words, which is straightforward. Let's look at an example: "I am sad, I am not learning NLP". The only words in our vocabulary that don't appear in this sentence are "happy" and "because": if we sum up the $\text{PosFreq}(1)$ associated with the remaining words in our vocabulary, i.e. the words that appear in this tweet, we get a scalar value of eight. We do the same for the negative frequencies, and we get a scalar value of eleven. So, we represent "I am sad, I am not learning NLP" $\rightsquigarrow [1, 8, 11]$.

1.1.3 Preprocessing

There are two major concepts here: stemming and "stop words". We'll learn how to apply these preprocessing steps to our data.

Stop words Stop words are defined as those which don't add significant meaning to the tweets; we *might* also choose to remove punctuation (if we decide it doesn't provide information in our context). In practice, this means comparing our tweet against two sets: one with stop words (in English) and another with punctuation.

| Stop Words | Punctuation |
|------------|-------------|
| and | , |
| is | . |
| are | : |
| at | ! |
| has | " |
| for | ' |
| a | |

In practice the list of stop words and punctuation marks are much larger, but for pedagogical purposes these will serve well. We might start out with a tweet like

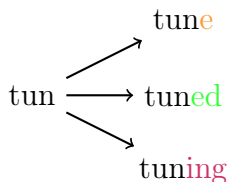
@YMourri and @AndrewYNg are tuning a GREAT AI Model at <https://deeplearning.ai>!!!

We then preprocess by stripping stop words "and", "are", "a" "at", and "a". The only punctuation that appears in this tweet that's also in our list is the exclamation point(s). We might further decide that tweets should have handles and URLs removed, because these don't add value for the specific task of sentiment analysis. In the end, we end up with a data point that looks like

tuning GREAT AI model

It's clearly a positive tweet, and a sufficiently good model should be able to classify it. Now that the tweet contains the minimum necessary information, we can perform *stemming* for every word.

Stemming Stemming in NLP is simply transforming any word to its base stem, which you could define as the set of characters that are used to construct the words and its derivatives. Let's look at the first word in the example: its stem is "tun", since



If we were to perform stemming on our entire corpus, the words "tune", "tuned", and "tuning" all get reduced to the stem "tun". So, your vocabulary would be significantly reduced in performing this process. You can further reduce the size of the vocabulary without losing valuable information by *lower-casing* every word, e.g. the words "GREAT", "Great", and "great" all get treated as the same word. Perhaps our final preprocessed tweet looks like

[tun, great, ai, model]

In summary, for our example of sentiment analysis on tweets, we might preprocess as follows:

1. Eliminate handles and URLs
2. Tokenize the string into words
3. Remove stop words like "and, is, a, on, etc."
4. Stemming - or convert every word to its stem. E.g. dancer, dancing, danced, becomes "danc". You can use Porter Stemmer to take care of this.
5. Convert all words to lowercase.

As an applied example:

I am Happy Because I am learning NLP @deeplearning $\xrightarrow{\text{Preprocessing}}$ [happy, learn, nlp] $\xrightarrow{\text{Feature Extraction}}$ [1, 4, 2]

where 1 is our bias term, 4 is the sum of positive frequencies, and 2 is the sum of negative frequencies. In practice, we are given a set of m raw tweets, and so we have to process them one-by-one to process them into an $m \times 3$ matrix, where each row describes the features for a given tweet.

$$\begin{bmatrix} 1 & X_1^{(1)} & X_2^{(1)} \\ 1 & X_1^{(2)} & X_2^{(2)} \\ \vdots & \vdots & \vdots \\ 1 & X_1^{(m)} & X_2^{(m)} \end{bmatrix}$$

The process is simple: (i) build the frequencies dictionary, (ii) initialize the matrix X to match the number of tweets, (iii) go through your sets of tweets and carefully preprocess by deleting stop words, stemming, deleting URLs/handles, and lowercasing, and finally (iv) extract the features by summing up the positive and negative frequencies of each of the tweets.

```

freqs = build_freqs(tweets, labels)      # Build frequencies dictionary.
X = np.zeros((m, 3))                    # Initialize matrix X.
for i in range(m):                      # For every tweet:
    p_tweet = process_tweet(tweets[i])   # Process tweet.
    X[i,:] = extract_features(p_tweet, freqs) # Extract features.

```

1.2 Logistic Regression

Previously, we've learned how to extract features, which we will now use to predict whether a tweet has a positive or negative sentiment. Logistic regression makes use of a sigmoid (or standard logistic) function which outputs a probability between zero and one. What's the recap from supervised machine learning? Recall figure 1.1: in the case of logistic regression our prediction function is going to be the standard logistic function:

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

where i denotes the observation number. Note that as $\theta^T x^{(i)}$ gets closer and closer to $-\infty$, the denominator of the sigmoid expression blows up and as a result the output values gets closer to zero. Conversely, as the inner product $\theta^T x^{(i)}$ gets closer to ∞ , the denominator of the sigmoid function approaches unit value and the resulting sigmoid expression evaluates to something near one. For classification, a threshold is needed, and it is natural to set it at $\frac{1}{2}$. For the logistic function, this threshold occurs when the inner product $\theta^T x^{(i)} = 0$. If the inner product is greater than (or equal to) zero, we classify as positive, else negative.

1.2.1 Learning Parameters

How to learn θ ? To train a logistic regression classifier, we need to iterate until we find a set of parameters θ that minimizes our cost function. Suppose we have a loss that depends only on the parameters θ_1, θ_2 : you might have a cost function that looks like follows, on the left, with the evaluation of the cost function plotted on the right as a function of the number of training iterations:



We might first initialize our parameters θ , then update our parameters in the direction of the *gradient of the cost function*. After a sufficient number of training steps, we will have updated θ to their optimal values where we are achieving near optimal cost. Let's quickly review this process of gradient descent for logistic regression:



1.2.2 Assessing model generalization

To analyze model fit, we need the following: $(X_{\text{val}}, Y_{\text{val}}, \theta)$, where we have *validation* data that was set aside during training, and a learned θ parameter vector. We will compute, for each example in X_{val} , the value of $h(\theta, x^{(i)})$ and compare it with our threshold value to make a prediction. In particular, our simple prediction function is given by

$$\hat{Y}_{\text{val}} = h(X_{\text{val}}, \theta) \geq \frac{1}{2}.$$

In particular, we will have a vector $h = [h_1 \ h_2 \ \dots \ h_m]$ where e.g. h_i could equal some float in $[0, 1]$, which we then convert into a binary label vector by applying our threshold. After building our predictions vector \hat{Y}_{val} , we can compare the predictions with the actual values and evaluate our test-set *accuracy*:

$$\text{Accuracy} = \sum_{i=1}^m \frac{(\text{pred}^{(i)} == \hat{Y}_{\text{val}}^{(i)})}{m}.$$

This metric gives an estimate of the number of times of logistic regression model will work correctly on unseen data.

1.2.3 Deriving Gradient Descent for Logistic Regression

Motivating where cost function comes from Let's examine the equation for the cost function for logistic regression:

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

The deep learning notes derive this equation in detail in the introduction. Let us briefly recap.

$$\Pr(y|x^{(i)}, \theta) = h(x^{(i)}, \theta)^{y^{(i)}} (1 - h(x^{(i)}, \theta))^{(1-y^{(i)})}.$$

We wish to maximize our function $h(\cdot, \theta)$ over the parameter space θ : when $y = 0$ we want $(1 - h(x^{(i)}, \theta))$ to be zero, and therefore $h(x^{(i)}, \theta)$ close to one. When $y = 1$, we want $h(x^{(i)}, \theta) = 1$. To model our entire dataset and not just one observation, we make an assumption of independence to arrive at a joint likelihood:

$$L(\theta) = \prod_{i=1}^m h(x^{(i)}, \theta)^{y^{(i)}} (1 - h(x^{(i)}, \theta))^{(1-y^{(i)})}.$$

Realize that if we “mess up” one prediction, we have the potential to “mess up” the entire cost function, which is what we want: we want a model that captures the entire dataset, where all datapoints are related. One issue: what happens when m grows? Then $L(\theta) \rightsquigarrow 0$, because the expressions $h(x^{(i)}, \theta)$ and correspondingly $(1 - h(x^{(i)}, \theta))$ are bounded between $(0, 1)$.

Optimization Using properties of logarithms (that they are monotone and maximizing a function under a monotone transformation doesn't change the optimum, and that they turn multiplication into addition), i.e.

$$\log(a * b * c) = \log a + \log b + \log c \quad \text{and} \quad \log a^b = b \log a.$$

We may now rewrite our optimization problem:

$$\begin{aligned}
\max_{h(x^{(i)}, \theta)} \log L(\theta) &= \log \prod_{i=1}^m h(x^{(i)}, \theta)^{y^{(i)}} (1 - h(x^{(i)}, \theta))^{1-y^{(i)}} \\
&= \sum_{i=1}^m \log h(x^{(i)}, \theta)^{y^{(i)}} (1 - h(x^{(i)}, \theta))^{1-y^{(i)}} \\
&= \sum_{i=1}^m \log h(x^{(i)}, \theta)^{y^{(i)}} + \log (1 - h(x^{(i)}, \theta))^{1-y^{(i)}} \\
&= \sum_{i=1}^m y^{(i)} \log h(x^{(i)}, \theta) + (1 - y^{(i)}) \log (1 - h(x^{(i)}, \theta))
\end{aligned}$$

We can then rescale by $\frac{1}{m}$ to get *average* cost. Recall we are maximizing over $h(x^{(i)}, \theta)$ in the equation above, and maximizing an equation is the same as minimizing its negative. Therefore,

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

A vectorized implementation is given by

$$\begin{aligned}
h &= g(X\theta) \\
J(\theta) &= \frac{1}{m} \cdot \left(-y^T \log(h) - (1 - y)^T \log(1 - h) \right)
\end{aligned}$$

Intuition for loss function of logistic regression Now, let's just go over some intuition here. Consider the term on the left-hand side of the parenthesized expression: this is the relevant term in your cost function when your label is 1. The term on the right is relevant when the label is zero. In general, this loss function simply says: the closer the prediction is to the observed label, the smaller the loss incurred. We can plot the cost as a function of our the prediction value for a single training example.



Deriving logistic regression gradient The general form of logistic regression is given by

Algorithm 1: General form of gradient descent

```

while not converged, and for all j do
  |  $\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$ 
end

```

We can work out the derivative using partial calculus to fill in the expression further:

Algorithm 2: Gradient descent for logistic regression

```

while not converged, and for all j do
  |  $\theta_j \leftarrow \theta_j - \frac{\alpha}{m} \sum_{i=1}^m (h(x^{(i)}, \theta) - y^{(i)}) x_j^{(i)}$ 
end

```

A vectorized implementation is given by

$$\theta := \theta - \frac{\alpha}{m} X^T (H(X, \theta) - Y).$$

Partial derivative of $J(\theta)$ It'll be helpful to first calculate the derivative of the sigmoid function.

$$\begin{aligned}
 h(x)' &= \left(\frac{1}{1 + e^{-x}} \right)' = \frac{-(1 + e^{-x})'}{(1 + e^{-x})^2} = \frac{-1' - (e^{-x})'}{(1 + e^{-x})^2} = \frac{0 - (-x)'(e^{-x})}{(1 + e^{-x})^2} = \frac{e^{-x}}{(1 + e^{-x})^2} \\
 &= \left(\frac{1}{1 + e^{-x}} \right) \left(\frac{e^{-x}}{1 + e^{-x}} \right) = h(x) \left(\frac{1 + e^{-x} - 1 + e^{-x}}{1 + e^{-x}} \right) = h(x) \left(\frac{1 + e^{-x}}{1 + e^{-x}} - \frac{1}{1 + e^{-x}} \right) = h(x)(1 - h(x)).
 \end{aligned}$$

The above was all for a computation of the derivative of the sigmoid function. But what about the derivative of $h(x^{(i)}, \theta) = \frac{1}{1 + e^{-\theta^T x^{(i)}}}$ with respect to θ_j ? Using the chain rule, because of the inner product $\theta^T x^{(i)}$, and applying toward θ_j , we see that the derivative would be

$$h(x^{(i)}, \theta) (1 - h(x^{(i)}, \theta)) x_j^{(i)}.$$

Now, we can compute the partial derivative of our loss function with respect to θ_j :

$$\begin{aligned}
\frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{-1}{m} \sum_{i=1}^m [y^{(i)} \log(h(x^{(i)}, \theta)) + (1 - y^{(i)}) \log(1 - h(x^{(i)}, \theta))] \\
&= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \frac{\partial}{\partial \theta_j} \log(h(x^{(i)}, \theta)) + (1 - y^{(i)}) \frac{\partial}{\partial \theta_j} \log(1 - h(x^{(i)}, \theta)) \right] \\
&= -\frac{1}{m} \sum_{i=1}^m \left[\frac{y^{(i)} \frac{\partial}{\partial \theta_j} h(x^{(i)}, \theta)}{h(x^{(i)}, \theta)} + \frac{(1 - y^{(i)}) \frac{\partial}{\partial \theta_j} (1 - h(x^{(i)}, \theta))}{1 - h(x^{(i)}, \theta)} \right] \\
&= -\frac{1}{m} \sum_{i=1}^m \left[\frac{y^{(i)} \frac{\partial}{\partial \theta_j} h(x^{(i)}, \theta)}{h(x^{(i)}, \theta)} + \frac{(1 - y^{(i)}) \frac{\partial}{\partial \theta_j} (1 - h(x^{(i)}, \theta))}{1 - h(x^{(i)}, \theta)} \right] \\
&= -\frac{1}{m} \sum_{i=1}^m \left[\frac{y^{(i)} h(x^{(i)}, \theta) (1 - h(x^{(i)}, \theta)) \frac{\partial}{\partial \theta_j} \theta^T x^{(i)}}{h(x^{(i)}, \theta)} + \frac{-(1 - y^{(i)}) h(x^{(i)}, \theta) (1 - h(x^{(i)}, \theta)) \frac{\partial}{\partial \theta_j} \theta^T x^{(i)}}{1 - h(x^{(i)}, \theta)} \right] \\
&= -\frac{1}{m} \sum_{i=1}^m \left[\frac{y^{(i)} h(x^{(i)}, \theta) (1 - h(x^{(i)}, \theta)) \frac{\partial}{\partial \theta_j} \theta^T x^{(i)}}{h(x^{(i)}, \theta)} - \frac{(1 - y^{(i)}) h(x^{(i)}, \theta) (1 - h(x^{(i)}, \theta)) \frac{\partial}{\partial \theta_j} \theta^T x^{(i)}}{1 - h(x^{(i)}, \theta)} \right] \\
&= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} (1 - h(x^{(i)}, \theta)) x_j^{(i)} - (1 - y^{(i)}) h(x^{(i)}, \theta) x_j^{(i)}] \\
&= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} (1 - h(x^{(i)}, \theta)) - (1 - y^{(i)}) h(x^{(i)}, \theta)] x_j^{(i)} \\
&= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} - y^{(i)} h(x^{(i)}, \theta) - h(x^{(i)}, \theta) + y^{(i)} h(x^{(i)}, \theta)] x_j^{(i)} \\
&= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} - h(x^{(i)}, \theta)] x_j^{(i)} \\
&= \frac{1}{m} \sum_{i=1}^m [h(x^{(i)}, \theta) - y^{(i)}] x_j^{(i)}
\end{aligned}$$

The vectorized version is simply given by

$$\nabla J(\theta) = \frac{1}{m} \cdot X^T (H(X, \theta) - Y).$$

2 Sentiment Analysis with Naive Bayes

2.1 Probability and Bayes Rule

Imagine you have an extensive corpus of tweets that can be categorized as either positive or negative, but not both.

2.2 Naive Bayes

Naive Bayes is often a “very good, quick, and dirty baseline” for many text classification tasks; it’s an example of supervised machine learning and as such shares many similarities with logistic regression. It’s called Naive because it makes the assumption that the features you’re using for classification are all independent, which in reality is *rarely* the case. As per usual, we start with two corpora: one for the positive tweets and one for the negative tweets:

| | | |
|--------------------------------------|--------------------|---------|
| Positive tweets | | |
| I am happy because I am learning NLP | word | Pos Neg |
| I am happy, not sad | I | 3 3 |
| | am | 3 3 |
| | happy | 2 1 |
| | because | 1 0 |
| | learning | 1 1 |
| | NLP | 1 1 |
| | sad | 1 2 |
| | not | 1 2 |
| Negative tweets | N_{class} | 13 13 |
| I am sad, I am not learning NLP | | |
| I am sad, not happy | | |

The above word frequencies table is the backbone input to our naive bayes algorithm: it allows us to compute conditional probabilities. E.g. $\Pr(I|\text{Pos}) = \frac{3}{13}$. We can do this for each word in our vocabulary, i.e. compute the conditional probability of it appearing in each class. Notice that if you sum over the probabilities for a particular class, you get 1.

| word | Pos | Neg |
|----------|----------------|----------------|
| I | $\frac{3}{13}$ | $\frac{3}{13}$ |
| am | $\frac{3}{13}$ | $\frac{3}{13}$ |
| happy | $\frac{2}{13}$ | $\frac{1}{13}$ |
| because | $\frac{1}{13}$ | 0 |
| learning | $\frac{1}{13}$ | $\frac{1}{13}$ |
| NLP | $\frac{1}{13}$ | $\frac{1}{13}$ |
| sad | $\frac{1}{13}$ | $\frac{2}{13}$ |
| not | $\frac{1}{13}$ | $\frac{2}{13}$ |
| Sum | 1 | 1 |

Let’s inspect some of the entries: notice that for a few words in the vocabulary, their conditional probabilities of appearing in either class are (nearly) identical: words that are equally probable don’t add anything to the sentiment. On the other hand, words like **happy**, or **sad**, **not** are “power” words which tend to express one sentiment or another. These words carry a lot of weight in determining your tweet sentiments. As a separate note, examine the word **because**: it only appears in the positive corpus, and so its conditional probability for the negative class is zero: when this happens we have no way of comparing between the two corpora which will become a problem for subsequent calculations. We’ll see how we can “smooth” our probability function.

Suppose we get a new tweet, “I am happy today; I am learning.” and we want to classify its sentiment. We use the following expression:

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

So, for our tweet example, we have (word by word, and skipping “today” because it doesn’t appear in our vocabulary):

$$\frac{\frac{3}{13}}{\frac{3}{13}} \times \frac{\frac{3}{13}}{\frac{3}{13}} \times \frac{\frac{2}{13}}{\frac{1}{13}} \times \frac{\frac{3}{13}}{\frac{3}{13}} \times \frac{\frac{3}{13}}{\frac{3}{13}} \times \frac{\frac{1}{13}}{\frac{1}{13}} = \frac{2}{13} > 1.$$

Because the ratio is greater than unit value, we conclude that overall the sentiment of the tweet is positive.

Laplacian smoothing This is a technique we use to avoid probabilities being identically zero. Typically, the expression used to calculate the conditional probability of a word, given the class, is

$$\Pr(w_i|\text{class}) = \frac{\text{freq}(w_i, \text{class})}{N_{\text{class}}} \quad \text{class} \in \{\text{Positive}, \text{Negative}\}$$

where N_{class} = frequency of all words in class. Laplacian smoothing does the following; supposing $|V|$ is the number of unique words in the vocabulary

$$\Pr(w_i|\text{class}) = \frac{\text{freq}(w_i, \text{class}) + 1}{N_{\text{class}} + |V|}. \quad (5)$$

By adding a one to our numerator, we ensure the expression is non-zero. However, this is not correctly normalized by N_{class} , and so we add a new term to the denominator $|V|$; this ensures the probabilities all sum to one. E.g. in our example table in section 2.2 describing positive and negative word frequencies in our corpora of tweets, we can use this to compute

$$\Pr(I|\text{Pos}) = \frac{3 + 1}{13 + 8}.$$

We can apply Laplacian smoothing to every entry in our table and end up with a new table of conditional probabilities where the column-sums are unit valued. Notice that if we apply this technique to the word “because” in our example, and specifically for the negative class, that $\Pr(\text{because}|\text{Negative}) = \frac{0+1}{13+8} > 0$ which solved our original problem of getting a divide by zero in the formula for Naive Bayes $\prod_{i=1}^m \frac{\Pr(w_i|\text{Pos})}{\Pr(w_i|\text{Neg})}$.

Log likelihoods Words can have many shades of emotional meaning, but for the purpose of sentiment classification they can be simplified into three categories: neutral, positive, and negative. A word can be taxonomized according to its conditional probabilities. We simply calculate for each word

$$\text{ratio}(w_i) = \frac{\Pr(w_i|\text{Pos})}{\Pr(w_i|\text{Neg})} \approx \frac{\text{freq}(w_i, 1) + 1}{\text{freq}(w_i, 0) + 1}.$$

If this ratio is identically unit valued, the word is neutral. Words that are more positive tend to have higher ratios (larger than one), and words that are more negative tend to have lower ratios (less than one). Observe that the ratio can lie in $[0, \infty)$.

It turns out that in our previous formulation of Naive Bayes, we assumed balanced class sizes. The correct formula for the likelihood includes the prior ratio, which becomes important for unbalanced datasets (where e.g. the number of positive and negative tweets is not equal):

$$\frac{\Pr(\text{Pos})}{\Pr(\text{Neg})} \prod_{i=1}^m \frac{\Pr(w_i|\text{pos})}{\Pr(w_i|\text{neg})} \quad (6)$$

Recognize that this computation involves the product of many probabilities that lie in $(0, 1]$, and we run the risk of numerical underflow if the number returned “is so small it can’t be stored on your device”. There is a nice mathematical trick that avoids this pitfall, and that’s to use properties of logarithms: $\log(a \times b) = \log(a) + \log(b)$.

$$\log \frac{\Pr(\text{Pos})}{\Pr(\text{Neg})} \prod_{i=1}^m \frac{\Pr(w_i|\text{pos})}{\Pr(w_i|\text{neg})} \rightsquigarrow \underbrace{\log \frac{\Pr(\text{Pos})}{\Pr(\text{Neg})}}_{\text{log prior}} + \underbrace{\sum_{i=1}^m \log \frac{\Pr(w_i|\text{Pos})}{\Pr(w_i|\text{Neg})}}_{\text{log likelihood}}.$$

Let $\lambda(w) = \log \frac{\Pr(w|\text{Pos})}{\Pr(w|\text{Neg})}$; we calculate this for each word in our vocabulary. Realize that neutral words (i.e. ones where $\Pr(w|\text{Pos}) = \Pr(w|\text{Neg})$) have $\lambda(w) = \log(1) = 0$. A positive sentiment is indicated by $\lambda(w) > 0$, and correspondingly $\lambda < 0$ indicates a negative sentiment. By using logarithms, we can reduce the risk of numerical underflow. Realize that our log-likelihood term can be expressed as $\sum_{i=1}^m \log \frac{\Pr(w_i|\text{Pos})}{\Pr(w_i|\text{Neg})} = \sum_{i=1}^m \lambda(w_i) \in (-\infty, \infty)$; we emphasize that our decision boundary is zero with our log-likelihood formula.

Training Naive Bayes In the context of Naive Bayes, “train” means something different than in logistic regression or deep learning: there’s no gradient descent; we’re just counting word frequencies in a corpus. There are five steps for training a Naive Bayes model:

1. Collect and annotate corpus (e.g. with positive and negative tweets)
2. Preprocessing (e.g. `process_tweet(tweet) ~> [w1, w2, ...]`)
 - Lowercase
 - Remove punctuation, urls, names, etc.
 - Remove stop words
 - Stemming
 - Tokenize sentences
3. Compute word counts, i.e. `freq(w, class)` and N_{class} .
4. Apply Laplacian smoothing to compute $\Pr(w|\text{class}) = \frac{\text{freq}(w, \text{class}) + 1}{N_{\text{class}} + |V_{\text{class}}|}$.
5. Calculate $\lambda(w) = \log \frac{\Pr(w|\text{Pos})}{\Pr(w|\text{Neg})}$.
6. Get the log-prior, which involves first counting D_{Pos} = number of positive tweets and D_{Neg} = number of negative tweets, whereby $\log \text{ prior} = \log \frac{D_{\text{Pos}}}{D_{\text{Neg}}}$.²

Testing Naive Bayes Once you’ve trained your model, you test it by taking the conditional probabilities derived and using them to predict the sentiments of new unseen tweets. We can evaluate model performance using test set accuracy. In particular, suppose we are given a tweet “I passed the NLP interview!”, and then after preprocessing we end up with [I, pass, the, NLP, interview]. We then look up our $\lambda(w)$ ’s that were calculated when we “trained” our model and compute the score, i.e. the log-prior plus log-likelihood for the unseen test case and compare it to our threshold (of zero). The values of the words that aren’t in our vocabulary are treated as neutral (i.e. zeros) and do not contribute to the final score (likelihood) of the unseen word: `pred = 1score>0`. If we are given a bunch of unseen words, i.e. data set aside during training $(X_{\text{val}}, Y_{\text{val}})$, we (i) compute `score = predict(X_{val} , λ , log-prior)` and then (ii) predict `pred = 1score>0`, and then (iii) compute test accuracy given by $\frac{1}{m} \sum_{i=1}^m (\text{pred}_i == Y_{\text{val}_i})$.

Applications of Naive Bayes There’s more we can do than just sentiment analysis. For example, we could do author identification: if you had two large corpora each written by different authors, you could train the model to recognize whether a document was written by one author or the other. Or, if you had some works by Shakespeare and some works by Hemmingway, you could calculate the λ for each word to predict how likely a new word is to be used by Shakespeare or alternatively Hemmingway. Another common use is spam filtering: $\frac{\Pr(\text{spam}|\text{email})}{\Pr(\text{non-spam}|\text{email})}$. One of the earliest applications of Naive Bayes was to filter between relevant and irrelevant documents in a database. I.e. given a set of keywords in a query, in this case, you can calculate the likelihood of the documents given the query:

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

Then, we have a decision rule that suggests retrieval if $\Pr(\text{document}_k|\text{query}) > \text{threshold}$, and then *sort* the documents based on their likelihoods; perhaps we choose to keep the first m results or ones with a

²If the dataset is balanced, the log-prior is zero since $D_{\text{Pos}} = D_{\text{Neg}}$ and $\log(1) = 0$.

likelihood above a certain threshold. Lastly, we can also use Naive Bayes for word disambiguation, i.e. breaking words down for contextual clarity. Consider that you have two possible interpretations of a given word within a text: let's say you don't know if the word "bank" in a text is referring to the bank of a river or a financial institution. To disambiguate your word, calculate the score of the document: $\frac{\text{Pr}(\text{river}|\text{text})}{\text{Pr}(\text{money}|\text{text})}$.

2.2.1 Assumptions of Naive Bayes

Naive Bayes is a very simple model: it doesn't involve setting any parameters. The method is called "naive" because of the assumptions it makes about the data. The first assumption is independence between predictors within each class, and the second has to do class (im)-balance with your validation sets. Let's explore each in detail and how they can affect our results.

Independence To illustrate what independence between features looks like, let's consider the following example:

It is sunny and hot in the Sahara Desert.

Naive Bayes assumes the words in a piece of text are independent of one another, but as you can see this is not always the case: the words "sunny" and "hot" often appear together as they do in this example. When taken together, they might also be related to the thing they're describing like a beach or a desert; so the words in a sentence aren't really independent of one another. But, Naive Bayes assumes that they are. The implication is that we could end up under *or* over estimating the conditional probabilities of individual words. E.g. if your task was to complete the sentence: "It's always cold and snowy in {blank}", then Naive Bayes might assign equal probability to the words spring, summer, fall, and winter even though from the context winter is the most likely candidate.³

Distribution of training data A good data set will contain the same proportion of (e.g. positive and negative) classes as a random sample would. However, most available annotated corpora are artificially balanced. E.g. in a real tweet stream a positive tweet is more likely to be sent than a negative tweet. Part of this has to do with platform decisions to perhaps ban content that is inappropriate or contains offensive vocabulary. Assuming that reality behaves as your training corpus could result in a very optimistic *or* pessimistic model.

2.2.2 Error Analysis

No matter what NLP method you use, you'll one day find yourself faced with an error, e.g. a misclassified sentence. How can we analyze such errors? Let's consider some possible errors in the model prediction that can be caused by:

- Removing punctuation and stop words – semantic meaning can be lost in the preprocessing step.
- Word order – can affect the meaning of a sentence.
- Adversarial attacks – language quirks can confuse Naive Bayes classifiers.

Removing punctuation Let's consider an example tweet: "My beloved grandmother :(". The sad face punctuation in this case is *very* important to the sentiment of the tweet because it tells you what's happening; but, if we remove punctuation then the processed tweet will leave behind a different (positive) sentiment. After processing, we may end up with [belov, grandmoth] which appear positive in nature.

³More sophisticated methods can deal with this issue.

Removing (stop) words It’s not just about punctuation either, consider as an example “this is not good, because your attitude is not even close to being nice”. If we remove stop words, we’re left with [good, attitude, close, nice]. From this set of words, any classifier would infer that the sentiment is positive. There are techniques we will learn about later to handle “nots” and word-order. For now, the takeaway is to look at the processed data to make sure your model can get an accurate read.

Word order The input pipeline isn’t the only source of trouble. E.g. consider the following two tweets:

I am happy because I did not go.
I am not happy because I did go.

The first is purely positive, the latter is negative. In this case, the “not” is important to the sentiment but gets missed by the Naive Bayes classifier: word order can be as important as spelling.

Adversarial attacks Lastly, let’s discuss adversarial attacks which essentially describe a language phenomenon like sarcasm, irony, and euphemism. Humans pick these up quickly but machines are terrible at it. The tweet, “This is a ridiculously powerful movie. The plot was gripping and I cried right through until the ending” contains a somewhat positive movie review, but pre-processing might suggest otherwise. I.e. if we pre-process, you’ll get a list of mostly negative words, but these words were in fact used to describe a movie that the author enjoyed. Applying Naive Bayes to this list of words would yield a negative score, unfortunately.

3 Vector Space Models

3.1 Motivating Vector Space Models

Let’s talk about vector spaces, and what types of information vectors can encode. We’ll talk about vector space models, their advantages, and some applications.

Suppose you have two questions:

Where are you heading?
Where are you from?

The sentences are identical but for the last word. However, the meaning is entirely different. On the other hand, the two questions:

What is your age?
How old are you?

contain no overlapping words but have identical meaning. Vector space models can identify similarity even when sentences don’t share common words; they can also be used for answering paraphrasing, and summarization. One advantage of vector space models is they allow us to capture dependencies between words. E.g.

You eat cereal from a bowl.
You buy something and someone else sells it.

In the first sentence, the words *cereal* and *bowl* are related (we’re eating one thing out of the other). In the second sentence, the latter half of it depends on the first half. Vector space models can capture this and many other types of relationships among different sets of words. When using a vector space models, the way that representations are made are by identifying the *context* around each word in the text, and this captures the relative meaning. “You shall know a word by the company it keeps” (John Firth).

3.2 Word by Word and Word by Doc

Let's discuss how to construct vectors based off of a co-occurrence matrix. Depending on the task you're trying to solve, you can have several possible designs.

Word by word To get a vector space model using a word by word design, you'll make a co-occurrence matrix and extract vector representations for the words in your corpus. We'll also discuss how to find relationships between words and vectors, also known as their similarity. We define the **co-occurrence** of two different words as the number of times they occur together within a certain distance k . E.g. suppose our corpus has the following two sentences:

I like simple data.
I prefer simple raw data.

Then, the row of the co-occurrence matrix corresponding to the word **data** with a value $k = 2$ would be populated with the following values:

| | simple | raw | like | I |
|------|--------|-----|------|---|
| data | 2 | 1 | 1 | 0 |

Table 1: The word “data” appears within distance one of the word “simple” in the first sentence, and within distance two of the word simple in the second example. Both of these are within our tolerance of $k = 2$, so the word by word co-occurrence entry for “data” and “simple” in this example is the value 2. We can proceed to compare the word “data” with other words appearing in our vocabulary, populating a representation that is based on word by word co-occurrence counts within a certain window size.

With a word by word design, you get a representation with $|V|$ entries, where V is our vocabulary.

Word by document With a word by document design, you'll count the number of times that words from your vocabulary appear in documents that belong to specific categories.



Here, we'd count the number of times that our words appear in the document that belong to each of the three categories. In this example, suppose that the word data appears 500 times in documents from your corpus related to entertainment, 6,620 times in economy documents, and 9,320 in documents related to machine learning. The word film might appear in each document's category 7,000, 4,000, and 1,000 times respectively.

| | Entertainment | Economy | Machine Learning |
|------|---------------|---------|------------------|
| data | 500 | 6620 | 9320 |
| film | 7000 | 4000 | 1000 |

Once you've constructed representations for multiple sets of documents or words, you'll get your vector space. For example, we could take a representation for the *words* “data” and “film” by looking at the rows

of the table (or matrix). However, we can also get a representation for each *category* of document type by looking at the columns! In our toy example, the vector space will have two dimensions: the number of times that the words “data” and “film” appear on the type of document.



Figure 2: In this vector representation, the document categories Economy and Machine Learning are more similar to each other than they are to entertainment. Measures of similarity include angle (cosine similarity) and distance (Euclidean, for example).

3.3 Distance and Similarity Metrics

Euclidean distance Euclidean distance is a similarity metric: it identifies how far apart two points (or vectors) are from each other.



Figure 3: Here we visualize the distance between two points, i.e. the length of the line segment connecting them. In two dimensions, this is given by $d(B, A) = \sqrt{(B_1 - A_1)^2 + (B_2 - A_2)^2}$; the first term is the horizontal distance squared, and the second term is the vertical distance squared.

| | data | boba | ice-cream |
|--------|------|------|-----------|
| AI | 6 | 0 | 1 |
| drinks | 0 | 4 | 6 |
| food | 0 | 6 | 8 |

Table 2: Suppose boba is word \vec{w} and ice-cream word \vec{v} , and we wish to calculate the Euclidean distance between these two word-vectors. This is given by $d(\vec{v}, \vec{w}) = \sqrt{(1 - 0)^2 + (6 - 4)^2 + (8 - 6)^2} = \sqrt{1 + 4 + 4} = \sqrt{9} = 3$.

In higher dimensions, i.e. an n dimensional vector, the process is similar.

$$d(\vec{v}, \vec{w}) = \sqrt{\sum_{i=1}^n (v_i - w_i)^2} = \|\vec{v} - \vec{w}\|_2.$$

Cosine similarity Another type of similarity function is cosine similarity. Euclidean distance can be problematic: in particular it can be biased by the size difference between the representations; to see why, consider the following example.



Figure 4: Realize that $d_2 < d_1$ due to the fact that, as measured by Euclidean distance, they are more similar in *magnitude*. In this case, the Agriculture and History corpora are much larger than the Food corpus; this causes their frequency counts along the dimensions of “disease” and “eggs” to be quite high, and their Euclidean distance is resultingly smaller than the distance between Food and Agriculture, which isn’t intuitive. What about considering the cosine of the angle between the vectors as a measure of similarity? In this case $\alpha < \beta$, and it happens to be that cosine similarity is a better metric than Euclidean distance under this representation.

An alternative measure to Euclidean distance is Cosine similarity: i.e. the cosine of the inner angle between two vectors. If the angle is small, the cosine would be close to unit value, whereas if the angle is close to 90 degrees, the cosine approaches zero. Recall that a vector norm is defined as $\|v\| = \sqrt{\sum_{i=1}^n v_i^2}$ and an inner-product is defined by $\langle u, v \rangle = u \cdot v = \sum_{i=1}^n u_i \cdot v_i$. Then, we can define cosine similarity as

$$\text{Cosine-similarity}(u, v) = \frac{\langle u, v \rangle}{\|u\| \|v\|}$$

If the input arguments (vectors) are perpendicular, the cosine of their inner angle is zero (since $\cos(90^\circ) = 0$), whereas if they point in the same direction their cosine angle is unit valued (since $\cos(0^\circ) = 1$, even if they have differing magnitudes).

Manipulating words in vector spaces How can we manipulate vectors using vector arithmetic to make predictions? I.e. we will use manipulate vector representations to infer unknown relations among words. Suppose we have a vector space with countries and their capital cities. Perhaps we know that the capital of the U.S. is Washington D.C., but that we don’t know the capital of Russia, and we further we’d like to use our knowledge of Washington D.C. with the U.S. to figure it out: it’s as easy as applying linear algebra.

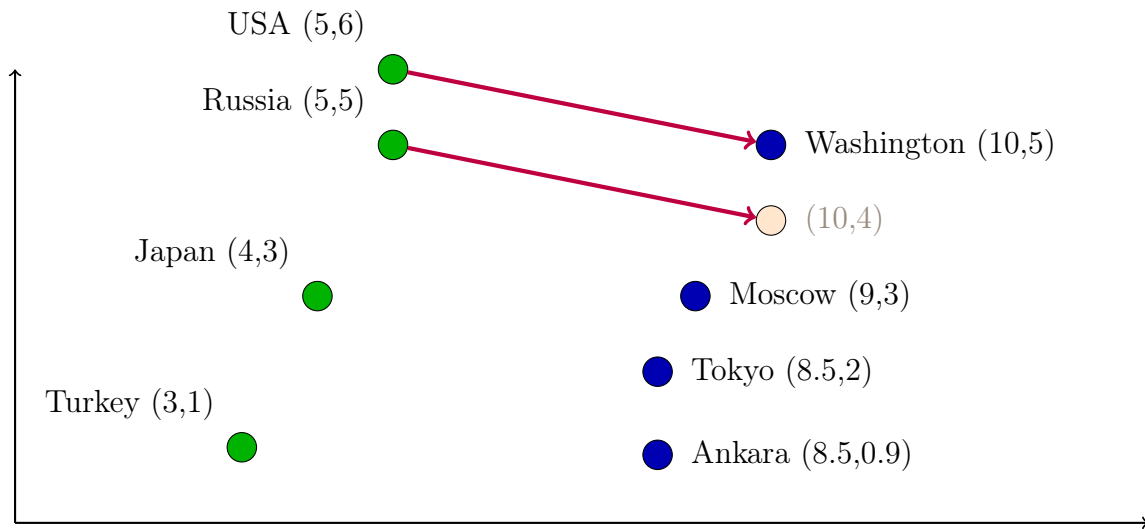
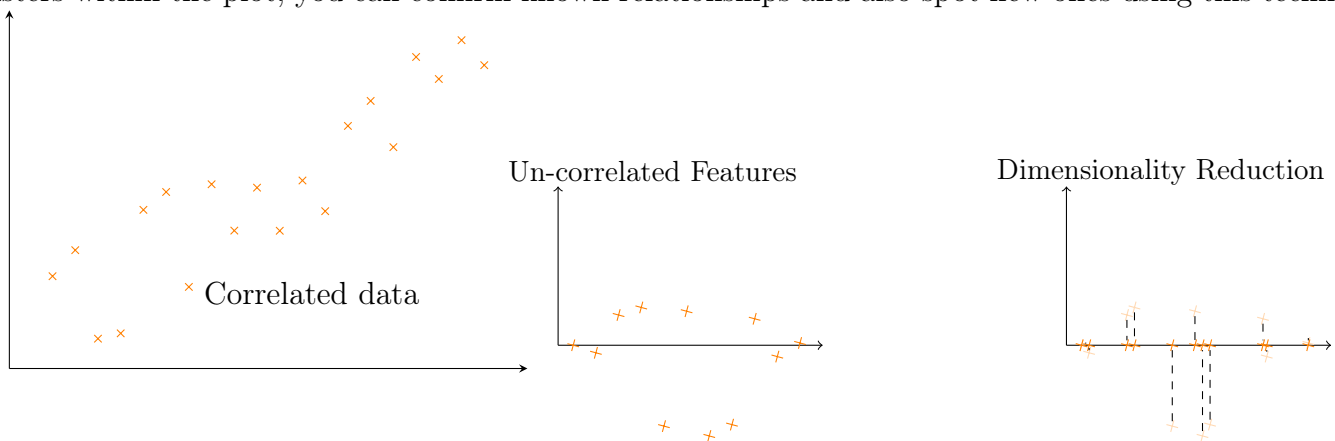


Figure 5: From Mikolov et. al. “Distributed Representations of Words and Phrases and their Compositionality”: in this representation, the vector between capitals Washington and USA is given by $[\text{Washington} - \text{USA}] = [5, -1]$. The easiest way to predict the capital of Russia is to simply take $\text{Russia} + [5, -1] = [10, 4]$. Because there are no cities with that representation, we’ll search for the capital city that is most similar by comparing each vector with the Euclidean distances (or Cosine similarities). The catch is that we need a vector space where the representations capture the relative meaning of words.

3.4 Principal Components Analysis

Visualization and PCA Often, we’ll have vector representations in very high dimension, but we’ll want to find a way to reduce the dimensions of the vectors to two dimensions so that we may plot it on a 2-D cartesian plane. What’s the intuition? Principal Component Analysis is the key: it allows us to take a representation with $d > 2$ dimensions and project it down to a lower dimension, e.g. $d = 2$, which we can then plot. If you perform PCA on your data, you might find that your word representation captures relationship between pairs of words that you expect to be related, in so far as these related words appear in clusters within the plot; you can confirm known relationships and also spot new ones using this technique.



Using PCA, we first find a set of uncorrelated features. Then, we can project our data into a lower dimensional space, trying to retain as much information as possible.

PCA algorithm How do we get uncorrelated features? Realize that the eigenvectors of the covariance matrix of our data correspond to the uncorrelated features for our data, and the eigenvalue corresponds to

the amount of information retained by each feature. Note that it's essential that our eigenvectors are sorted according to their eigenvalues from largest to smallest (in magnitude): this means that if we use the first k eigenvectors to get a representation we are greedily selecting vectors that retain the most information from our data.

1. Mean-normalize data: $x_i = \frac{x_i - \mu_{x_i}}{\sigma_{x_i}}$.
2. Obtain covariance matrix Σ .
3. Perform Singular Value Decomposition on Σ to get three matrices, where the first stores the eigenvectors U (columnwise), the second S stores the eigenvalues along the diagonal.
4. Dot product to Project your data: $X' = XU[:, 0 : k]$, where k is the number of dimensions we wish to reduce to.
5. Calculate Percentage of Retained Variance: $\frac{\sum_{i=0}^{k-1} S_{ii}}{\sum_{j=0}^{d-1} S_{jj}}$.

Eigenvectors give the direction of uncorrelated features, whereas the eigenvalues are the variance of the new features. The dot product gives the projection on uncorrelated features.

4 Machine Translation

4.1 Transforming Word Vectors

Word vectors can capture important properties of words. Let's make use of them to learn to align words into different languages, which will give us a basic translation program; locality sensitive hashing will speed it up.

Overview of translation In order to translate an English word to a French word, one naive way could be to generate an extensive list of English words and their associated French word. If you ask a human to do this, you would find someone who knows both languages to start making a list. If, however, you want a machine to do this, you need to first calculate word embeddings associated with English and word embeddings associated with French. Next, retrieve the English word embedding of a particular English word such as "cats", then find some way to transform the English word embedding into word embedding that has a meaning in the French word vector space; we'll see how to convert between vector spaces in a moment. The next step is to search for word vectors in the French word vector space that are most similar to the transformed vector representation: the most similar words are candidate words for our translation.

Transforming vectors using matrices Define a randomly selected matrix R and see how well it works to transform our English vector space representation X into our French vector space representation Y , i.e. compute $XR \approx Y$. In order for this to work, you'll need to first get a subset of English words and their French equivalence, get the respective word vectors, and stack the word vectors in their respective matrices X and Y . The key is to keep the rows lined up (or to align the row vectors): i.e. if the first row of the English matrix X contains the word representation for "cat", then so should the French matrix Y contain the word representation for the translated word. It's natural at this point to ask: if we already have the mappings from $X \rightsquigarrow Y$ (i.e. the translation pairs), then why not just store them in an associative data-structure like a dictionary...why train a model at all? The nice property of our transformation matrix R is that it will *generalize* to translate unseen example inputs. So, we only need to train on a subset of the English-French vocabulary, and *not* the entire vocabulary.

Finding a good transformation matrix We define our loss function as

$$\text{Loss} = \|XR - Y\|_F.$$

where $\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |A_{ij}|^2}$. If we start with a random matrix R , we can gradually improve upon it in an iterative fashion. We first compute the gradient by taking the gradient of the loss function with respect to matrix R , i.e. $g = \frac{\partial}{\partial R} \text{Loss}$, then update using gradient descent, i.e. $R \leftarrow R - \alpha g$. We can either pick a fixed number of times to go through the loop, or exit the loop when the loss falls below a certain threshold. In practice, it's easier to minimize the the square of the norm instead. What about the gradient of the loss function, which is defined as the derivative of the loss with respect to the matrix R . With our loss:

$$g = \frac{\partial}{\partial R} \text{Loss} = \frac{2}{m} (X^T (XR - Y)).$$