

Homework 3

Linear classifiers for text

The Yelp restaurant reviews dataset is comprised of user-contributed reviews posted to Yelp for restaurants in Pittsburgh collected several years ago. Each review is accompanied by a binary indicator of whether or not the reviewer-assigned rating is at least four (on a five-point scale). The text of the reviews has been processed to replace all non-alphanumeric symbols with whitespace, and all letters have been changed to lowercase. The prediction problem we consider here is predicting the binary indicator from the review text.

The training data is contained in the file `reviews_tr.csv`, and the test data is contained in the file `reviews_te.csv`. The following code can be used to load the training data.

```
from csv import DictReader

vocab = {}
vocab_size = 0
examples = []
with open('reviews_tr.csv', 'r') as f:
    reader = DictReader(f)
    for row in reader:
        label = row['rating'] == '1'
        words = row['text'].split(' ')
        for word in words:
            if word not in vocab:
                vocab[word] = vocab_size
                vocab_size += 1
        examples.append((label, [vocab[word] for word in words]))
```

This code assigns every word that appears in training data a unique non-negative integer, and the mapping is provided in the dictionary `vocab`. Each training example is represented as a binary label paired with a list of non-negative integers (which we'll call “word IDs”) corresponding to the words that appear in the review. Note that the list may contain repetitions of any given word ID (since a review may contain repetitions of a word). The examples are collected in the array `examples`. Throughout the rest of this section, we'll refer to the variables defined by this code.

Of course, the test data can be loaded in a similar way. Note that the test data may contain words that do not appear in the training data.

Different examples may have different number of words, so the lists of word IDs may have different lengths. It may not be obvious how they can be used with common machine learning algorithms. However, it is relatively easy to transform a list of word IDs into a “bag-of-words” vector (also called “term frequency” vector), which is a vector $x = (x_0, \dots, x_{d-1})$ where x_j is

the number of times word ID j appears in the list. **This is the representation you should use for the feature vectors throughout this section.** The dimension d of the vector is equal to the vocabulary size `vocab_size` (and the vector indices, like the word IDs, start from 0). Such a vector representation is readily usable by many machine learning algorithms.

The following code applies this transformation to the first training example.

```
from numpy import zeros

def bag_of_words_rep(word_ids, dim):
    bow_vector = zeros(dim) # creates a numpy.ndarray of shape (dim,)
    for word_id in word_ids:
        bow_vector[word_id] += 1
    return bow_vector

first_bow_vector = bag_of_words_rep(examples[0][1], vocab_size)
```

Note that applying this transformation to all training examples would result in a large collection of high-dimensional vectors (each of dimension `vocab_size`), ultimately consuming a lot of memory:

```
from sys import getsizeof
print('memory required for examples:', getsizeof(examples))
print('(estimate of) memory required for bag-of-words representation:',
      ↪ len(examples) * getsizeof(first_bow_vector))
```

The difference is quite a few orders-of-magnitude. Therefore, function `bag_of_words_rep` should not be used explicitly; its purpose here is just to explain the semantics of the “bag-of-words” representation. A memory-efficient learning algorithm should use `examples` directly as input and avoid explicitly forming the “bag-of-words” vectors for these examples.

Online Perceptron

The Online Perceptron algorithm is a variant of the standard Perceptron algorithm. Recall that the standard Perceptron algorithm maintains a weight vector w , and repeatedly updates the weight vector with any training example that is misclassified by the (homogeneous) linear classifier with parameter w . This continues until there are no misclassified training examples. Of course, if the training data is not linearly separable, this will go on forever—some examples will be used to update the weight vector infinitely often. The Online Perceptron, shown below, rectifies this by only considering each training example exactly once, in some fixed ordering over the training examples, and then halting.

- Input: training examples $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$
- Initialize: $w^{(0)} = (0, \dots, 0)$
- For $i = 1, 2, \dots, n$:

- If $(x^{(i)}, y^{(i)})$ is misclassified by the linear classifier with parameter $w^{(i-1)}$, then

$$w^{(i)} \leftarrow \begin{cases} w^{(i-1)} + x^{(i)} & \text{if } y^{(i)} = \text{true} \\ w^{(i-1)} - x^{(i)} & \text{if } y^{(i)} = \text{false} \end{cases}$$

- Else $w^{(i)} \leftarrow w^{(i-1)}$

- Return: $w^{(n)}$

Here, we have assumed that each label is either “true” or “false”, and that a linear classifier with parameter w predicts “true” on an input feature vector x if and only if $w^\top x > 0$.

Problem 1. Implement the Online Perceptron algorithm and apply it to the training data (in the order that they appear in `reviews_tr.csv`). Briefly explain how you avoid explicitly forming the bag-of-words feature vectors in your implementation. What is the training error rate of the linear classifier returned by Online Perceptron? And what is the test error rate (based on the test data from `reviews_te.csv`)? Please report three significant digits for all computed error rates.

Recall, that the test data may contain words that do not appear in the training data. It may seem that this could be cumbersome to deal with. However, it turns out to be easy to handle because the initial weight vector in Online Perceptron is the zero vector.

Problem 2. To get a sense of the behavior of the linear classifier that you obtained above, determine the 10 words that have the highest (i.e., most positive) weights in the weight vector, and also determine the 10 words that have the lowest (i.e., most negative) weights in the weight vector. Report the actual words, not the word IDs. You may find it helpful to invert the mapping given in `vocab`.

Upgrades

There are many ways to upgrade Online Perceptron. Below, we discuss two possibilities.

The first is the “Averaged Perceptron” variant, which is exactly the same as Online Perceptron, except that the weight vector returned is $w^{\text{avg}} := \frac{1}{n} \sum_{i=1}^n w^{(i)}$ (instead of just $w^{(n)}$). This is more “stable” than Online Perceptron, since the average of the weight vectors is much less sensitive to a single training example than the current weight vector maintained by Online Perceptron. (Of course, $w^{\text{sum}} := \sum_{i=1}^n w^{(i)}$ works just as well.)

The second is to improve the feature representation. So far, we have only used the bag-of-words representation of the reviews. It is called bag-of-words because it only depends on the number of times words appear in the review; the order of the words doesn’t matter at all. One way to take the order of words into account (in a very limited way) is to also consider the number of times “bigrams” appear in the review. A bigram is a pair of words that appear consecutively. For example, in “a rose is a rose”, the bigrams that appear are: (a, rose), which appears twice; (rose, is); and (is, a). If there are d words in a vocabulary, then there are d^2 possible bigrams, which can be enormous when d is even just moderately large.

This makes it even more important to avoid explicitly forming the bag-of-words-and-bigrams feature vectors.

Problem 3. Implement one of the above suggested improvements, or any other improvement you can think of. (If you come up with your own improvement, give a high-level description of it and explain its motivation.) Apply the new algorithm to the training data. What are the training and test error rates of the classifier you obtain? Which 10 words (or bigrams) have the highest weights in the weight vector (à la Problem 2), and which have the lowest weights?

Logistic regression

Problem 4. Let $LL: \mathbb{R}^d \rightarrow \mathbb{R}$ be the log-likelihood function based on training data $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$ from $\mathbb{R}^d \times \{0, 1\}$ under the logistic regression model:

$$LL(w) = \sum_{i=1}^n \left(y^{(i)} w^\top x^{(i)} - \ln(1 + e^{w^\top x^{(i)}}) \right).$$

Suppose you have a weight vector $\hat{w} \in \mathbb{R}^d$ such that

$$\hat{w}^\top x^{(i)} > 0 \quad \text{if and only if} \quad y^{(i)} = 1, \quad \text{for all } i = 1, \dots, n.$$

A statistician friend asks you if \hat{w} (approximately) maximizes the log-likelihood function. However, in a computing accident, you lose access to the training data and cannot compute $LL(\hat{w})$. The system administrator is only able to recover two pieces of information about the training data for you: the number of training examples n , and $\gamma = \min\{|\hat{w}^\top x^{(1)}|, \dots, |\hat{w}^\top x^{(n)}|\}$. Assume that $\gamma > 0$. Write a procedure that, given \hat{w} , n , γ , and an arbitrary positive number $\epsilon > 0$, returns another weight vector $v \in \mathbb{R}^d$ such that $LL(v) \geq -\epsilon$. Explain why your procedure works.

Feature maps

Problem 5. Specify a feature map $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^p$, with p as small as possible, so that homogeneous linear classifiers in the feature space are able to represent all binary classifiers $f_{c,r}: \mathbb{R}^d \rightarrow \{0, 1\}$ of the following form:

$$f_{c,r}(x) = \begin{cases} 1 & \text{if } \|x - c\| < r \\ 0 & \text{otherwise.} \end{cases}$$

Here, $c \in \mathbb{R}^d$ and $r > 0$ are parameters of the classifier. In other words, for every $c \in \mathbb{R}^d$ and $r > 0$, there should exist $w \in \mathbb{R}^p$ such that

$$\varphi(x)^\top w > 0 \quad \text{if and only if} \quad f_{c,r}(x) = 1, \quad \text{for all } x \in \mathbb{R}^d.$$

Explain why your feature map works.

Problem 6. Suppose the distribution P over $\mathbb{R}^2 \times \{0, 1\}$ comes from the bivariate normal generative model.

- It turns out the log-odds function $x \mapsto \ln \frac{\Pr(Y=1|X=x)}{\Pr(Y=0|X=x)}$ for $(X, Y) \sim P$ can be expressed as a polynomial in x . What is the degree of this polynomial?
- Let $f: \mathbb{R}^2 \rightarrow \{0, 1\}$ be the classifier that has smallest possible error rate under P . If φ is a polynomial feature map, how large must its degree be so that f can be represented

as a homogeneous linear classifier in the feature space given by φ ?

- (c) Repeat the previous parts in the case where the bivariate normal generative model comes with additional restriction that the covariance matrix for each of the two class conditional distributions is the identity matrix.

Problem 7. Consider training data $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})$ from $\mathbb{R}^d \times \{0, 1\}$ for a binary classification problem about individuals from a population. Within this population, there are d well-defined subgroups (e.g., subgroup 1 is “people born in New York”, subgroup 2 is “people with blue eyes”). Let (X, Y) be a random example whose distribution is the *empirical distribution* on these n training examples, i.e., for each $(x, y) \in \mathbb{R}^d \times \{0, 1\}$,

$$\Pr[(X, Y) = (x, y)] = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{(x^{(i)}, y^{(i)}) = (x, y)\}.$$

Each feature vector $x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})$ is comprised of $\{0, 1\}$ -valued features that indicate membership in these subgroups: for each subgroup $j \in \{1, 2, \dots, d\}$,

$$x_j^{(i)} = \begin{cases} 1 & \text{if individual } i \text{ belongs to subgroup } j; \\ 0 & \text{otherwise.} \end{cases}$$

Specify a feature map $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^p$ (with p as small as possible) so that the maximum likelihood estimate $\hat{w} \in \mathbb{R}^p$ of the weight vector parameter in logistic regression (with feature map φ) ensures all the following simultaneously:

1. $\Pr[Y = 1] = \mathbb{E}[\text{logistic}(\varphi(X)^\top \hat{w})]$;
2. $\Pr[Y = 1 \mid X_j = 1 \wedge X_k = 1] = \mathbb{E}[\text{logistic}(\varphi(X)^\top \hat{w}) \mid X_j = 1 \wedge X_k = 1]$ for each $j \in \{1, \dots, d\}$ and $k \in \{1, \dots, d\}$.

(You should assume that \hat{w} is a maximizer of the log-likelihood objective in the logistic regression model.)

Solutions

- https://drive.google.com/file/d/1l5zEQmb5eiJ2rXaXrosSMh6PIZy7pdza/view?usp=drive_link
- https://drive.google.com/file/d/1Vu1car7p5j8ALat6fL2FUPhIY5B-eK4v/view?usp=drive_link