

cs134-pa1

September 13, 2019

1 Train a logistic regression model for text classification from scratch

Due September 27, 2019

1.1 Overview

Logistic regression (aka Maximum Entropy) has been one of the workhorses in natural language processing. It has also been used very widely in other fields. The model has many strengths. It is effective when you have a large number of features. It also handles the correlation between features really well. In addition, the time and space complexity for the model is not too taxing. All of these reasons make the Logistic Regression model a very versatile classifier for many NLP tasks. In this assignment, we are going to use the Logistic Regression model for sentiment analysis.

1.2 Getting ready

Download the starter code along with the datasets using the link on latte. We highly recommend reviewing the slides before starting implementing.

1.3 Datasets

- **Sentiment analysis on Yelp dataset.** We processed the data from Yelp dataset challenge and put it in the json format. Review class provided in corpus.py reads in the input for you. Each restaurant review is tagged with 'negative' (1-2 stars), 'neutral' (3 stars), or 'positive' (4-5 stars). We are building a classifier that labels a review as one of the three types of sentiments based on the review text.
- **Name classification.** This should be used to test and debug your code. You can just use the first character and the last character as your features.

1.4 What needs to be done

A Logistic Regression model has a lot of moving parts. The list below guides you through what needs to be done. The three things you need to focus on are: representation, learning, and inference.

1.4.1 Representation

- Choose the feature set. Start with the feature set that is small or the learning algorithm will take a long time to run and this makes the debugging process difficult. You can ramp up the number and variety of features when your code is thoroughly tested.
- Choose the data structure that holds the features. We recommend sparse feature vectors. Regardless of your choice, cache the features internally within each Document object as the algorithm is iterative. Featurization should be done only once.
- Choose the data structures that hold the parameters. We recommend using a $k \times p$ matrix where k is the number of labels, and p is the number of linguistic features.

1.4.2 Learning

- Compute the negative log-likelihood function given a minibatch of data. You will need this function to track progress of parameter fitting.
- Compute the gradient with respect to the parameters. You will need the gradient for updating the parameters given a minibatch of data.
- Implement the mini-batch gradient descent algorithm to train the model to obtain the best parameters.

1.4.3 Classification/Inference.

- Apply the model to unseen data. The provided test cases will evaluate the model on both datasets, and you must pass those.

1.4.4 Experiments

In addition to implementing the mini-batch gradient descent algorithm for the Logistic Regression model, you are asked to do the following experiments to better understand the behavior of the model. For all three experiments, use the Yelp review dataset as it is a more realistic one.

Experiment 1 -- Training set size:

Does the size of the training set really matter? The mantra of machine learning tells us that the bigger the training set the better the performance. We will investigate how true this is.

In this experiment, fix the feature set to something reasonable and fix the dev set and the test set. Vary the size of the training set $\{1000, 10000, 50000, 100000, \text{and all}\}$ and compare the (peak) accuracy from each training set size. Make a plot of size vs accuracy. Analyze and write up a short paragraph on what you learn or observe from this experiment.

Experiment 2--- Minibatch size:

Why are we allowed to use mini-batch instead of the whole training set when updating the parameters? This is indeed the dark art of this optimization algorithm, which works well for many complicated models, including neural networks. Computing gradient is always expensive, so we want to know how much we gain from each gradient computation.

In this experiment, try minibatch sizes $\{1, 10, 50, 100, 1000\}$, using the best training size from Experiment 1. For each mini-batch size, plot the number of datapoints that you compute the gradient for (x-axis) against the accuracy of the development set (y-axis). Analyze and write up a short paragraph on what you learn or observe from this experiment.

Experiment 3 -- Hyperparameter tuning:

Try different values of $\lambda = \{0.1, 0.5, 1, 10\}$ for $L2$ regularization and observe its effect on the accuracy of the model against the development set. Make a plot of lambda value vs accuracy on

the development set. Write a short paragraph summarizing what you have observed from this experiment.

As you are doing more experiments, the number of experimental settings starts to multiply. Use your best settings from your Experiment 1 and your Experiment 2 for the tuning of the $L2$ regularization parameters. It's not advisable to vary more than one experimental variable at a time, and that'll make it hard to interpret your results. You can set up a grid search procedure to do this experiment automatically without manual intervention.

Experiment 4 -- Feature Engineering (Extra credit for creative and effective features): In addition to bag-of-words features, experiment with additional features (bigrams, trigrams, etc.) to push up the performance of the model as much as you can. The addition of new features should be driven by error analysis. This process is similar to the machine learning process itself, only that it involves actual humans looking at the errors made of the machine learning model and trying to come up with new features to fix or reduce those errors. Briefly describe what new features you have tried if they are useful.

1.5 Submission

Submit the following on Latte:

1.5.1 All your code.

But don't include the datasets as we already have those.

1.5.2 Report.

Please include the following sections in your report: 1. A brief explanation of your code structure

2. How to run your code, and what output to expect
3. Experimental settings (Explain clearly what feature set is being used and how you set up mini-batch gradient descent because there can be quite a bit of variations.)
4. Experimental results

Please keep this report no more than two pages single-spaced including graphs.

1.6 More on Mini-batch Gradient Descent

In this assignment, we will train Logistic Regression models using mini-batch gradient descent. Gradient descent learns the parameter by iterative updates given a chunk of data and its gradient.

If a chunk of data is the entire training set, we call it batch gradient descent.

```
In [ ]: while not converged:
        gradient = compute_gradient(parameters, training_set)
        parameters -= gradient * learning_rate
```

Batch gradient descent is much slower. The gradient from the entire dataset needs to be computed for each update. This is usually not necessary. Computing gradient from a smaller subset of the data at a time usually gives the same results if done repeatedly.

If a subset of the training set is used to compute the gradient, we call it mini-batch gradient descent. This approximates the gradient of batch gradient descent.

```
In [ ]: while not converged:
        minibatches = chop_up(training_set)
        for minibatch in minibatches:
            gradient = compute_gradient(parameters, minibatch)
            parameters -= gradient * learning_rate
```

If a chunk of data is just one instance from the training set, we call it stochastic gradient descent (SGD). Each update only requires the computation of the gradient of one data instance.

```
In [ ]: while not converged:
        for datapoint in training_set:
            gradient = compute_gradient(parameters, datapoint)
            parameters -= gradient * learning_rate
```

1.6.1 Practical issues with mini-batch gradient descent

- How should I initialize the parameters at the first iteration?

Set them all to zero. This is generally not advisable for more complicated models. But for the Logistic Regression model, zero initialization works perfectly.

- How do I introduce the bias term?

Include a feature that fires in ALL data instances. And treat it as a normal feature and proceed as usual.

- Why do the posterior $P(Y|X)$ become NaN?

It is very likely that you exponentiate some big number and divide by the same amount i.e. if `unnormalized_score` is a vector of unnormalized scores (the sum of lambdas), then:

```
posterior = exp(unnormalized_score) / sum(exp(unnormalized_score))
```

This is no good. We have to get around by using some math tricks:

```
posterior = exp(unnormalized_score - scipy.misc.logsumexp(unnormalized_score))
```

If this confuses you or you are not sure why this is correct, think about it more or ask the TAs. But we are quite confident that you will need to use the `logsumexp` function.

- How do you know that it converges?

It is extremely difficult to know. If you stop too early, the model has not reached its peak yet i.e. *underfitting*. If you stop too late, the model will fit too well to the training set and not generalize to the unseen data i.e. *overfitting*. But there are multiple ways to guess the convergence. We suggest this method called *early stopping*.

Every once in a while evaluate the model on the development set during gradient descent.

- If the performance is better than last evaluation, then save this set of parameters and keep going for a little more.

- If the performance stops going up after a few updates, stop and use the last saved parameters. (How many is a few? Up to you)
- How often should I run evaluation on the dev set during training?
Up to you. It is actually OK to run the evaluation on the dev at every update you make to the parameters.
- How do I know that my implementation is correct?
Look at the average negative log-likelihood. It should keep going down monotonically i.e. at every single update. You should also see that the gradient should get closer and closer to zero.

In []: