RECITATION 4 LOGISTIC REGRESSION

10-301/601: Introduction to Machine Learning 06/16/2022

1 Parameter Derivation Practice

We would like to fit a regression model to the dataset

$$\mathcal{D} = \left\{ \left(\mathbf{x}^{(1)}, y^{(1)}\right), \left(\mathbf{x}^{(2)}, y^{(2)}\right), \cdots, \left(\mathbf{x}^{(N)}, y^{(N)}\right) \right\}$$

with $\mathbf{x}^{(i)} \in \mathbb{R}^M$ by minimizing the loss function:

$$J(\mathbf{w}) = \sum_{i=1}^{N} \left(y^{(i)} - \sum_{l=1}^{L} w_l \sum_{j=1}^{M} w_j x_j^{(i)} \right).$$

To find a best model match, we want to minimize the loss function with respect to the weights \mathbf{w} . This is just a practice in taking the first order condition of $J(\mathbf{w})$ with respect to each $w_k \in \mathbf{w}$, setting them equal to zero, and solving for w_k . In other words, we set $\frac{\partial J(\mathbf{w})}{\partial w_k} = 0$ and solve for w_k .

(a) Solve for $\frac{\partial J(\mathbf{w})}{\partial w_k}$ and set it equal to zero.

(b) Fill in the blanks to solve for w_k .

$$0 = -\sum_{i=1}^{N} \left(\sum_{j=1}^{M} w_j x_j^{(i)} + x_k \sum_{l=1}^{L} w_l \right)$$

$$0 = -\sum_{i=1}^{N} \left(-\dots + \sum_{j=1,\dots}^{M} w_j x_j^{(i)} + \dots + x_k^{(i)} \sum_{l=1,\dots}^{L} w_l \right)$$

$$0 = \dots - \sum_{i=1}^{N} \left(\sum_{j=1,\dots}^{M} w_j x_j^{(i)} + x_k^{(i)} \sum_{l=1,\dots}^{L} w_l \right)$$

2 Binary Logistic Regression

Consider the following dataset,

$$\mathcal{D} = \{ (\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)}) \} \text{ where } \mathbf{x}^{(i)} \in \mathbb{R}^M, y^{(i)} \in \{0, 1\}.$$

Notice that instead of regressing on continuous variables, as we do in linear regression, we wish to predict on discrete variables, and, more specifically, binary outcomes (this process is called classification). Since we want to take noise into account when making predictions, we will predict the probability of each outcome given some inputs, each denoted $\mathbf{x}^{(i)}$. This means that the result of the linear combination of each input, $\mathbf{x}^{(i)}$ with the parameter vector, $\boldsymbol{\theta}$, or, $\boldsymbol{\theta}^T \mathbf{x}^{(i)}$, must be manipulated to fit between zero and one. We use the sigmoid function, $\sigma(\cdot)$, for this purpose.

Recall that

$$\sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x}^{(i)})} = \frac{\exp(\boldsymbol{\theta}^T \mathbf{x}^{(i)})}{1 + \exp(\boldsymbol{\theta}^T \mathbf{x}^{(i)})}.$$

The conditional probability of $y^{(i)}$ given $\mathbf{x}^{(i)}$ is,

$$p\left(y^{(i)} \mid \mathbf{x}^{(i)}, \boldsymbol{\theta}\right) = \begin{cases} \sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) & y^{(i)} = 1\\ 1 - \sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) & y^{(i)} = 0 \end{cases}.$$

We can rewrite this as a single statement,

$$p(y^{(i)} \mid \mathbf{x}^{(i)}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)})^{y^{(i)}} (1 - \sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)}))^{(1-y^{(i)})}.$$

Can you show why this is true?

We assume each observation of $y^{(i)}$ in the data is independent and identically distributed. This means we can write down the likelihood of the data as a product of negative conditional probabilities.

Let's plug in the expression for $p\left(y^{(i)} \mid \mathbf{x}^{(i)}, \boldsymbol{\theta}\right)$ and simplify.

In stochastic gradient descent, we use only a single $\mathbf{x}^{(i)}$. Given $\boldsymbol{\phi}^{(i)} = \sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)})$ and

$$J^{(i)}(\boldsymbol{\theta}) = -y^{(i)}\log(\boldsymbol{\phi}^{(i)}) - (1 - y^{(i)})\log(1 - \boldsymbol{\phi}^{(i)})$$

Show that the partial derivative of $J^{(i)}(\boldsymbol{\theta})$ with respect to the jth parameter θ_j is as follows:

$$\frac{\partial J^{(i)}(\boldsymbol{\theta})}{\partial \theta_i} = (\sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) - y^i) x_j^{(i)}$$

Remember,

$$\frac{\partial \boldsymbol{\phi}^{(i)}}{\partial \boldsymbol{\theta}_{i}} = \boldsymbol{\phi}^{(i)} * (1 - \boldsymbol{\phi}^{(i)}) * \frac{\partial \boldsymbol{\theta}^{T} \mathbf{x}^{(i)}}{\partial \theta_{i}}$$

.

3 Feature Representation for Sentiment Classification

In many machine learning problems, we will want to find appropriate representations for the inputs of the algorithm we are developing. In Homework 4, we will work on using logistic regression for a sentiment classification task, where our algorithm takes a paragraph of movie review as the input and outputs a binary value denoting whether the review is positive or not. To build an appropriate representation for the input (aka. the review text), we consider a representation built using Word2vec¹ word embeddings.

In this section, consider a scenario where we are interested in representing the following text:

```
a hot dog is not a sandwich because it is not square (1)
```

We consider the following dictionary (denoted below as **Vocab**) as the set of vocabulary that we will consider. Note that the vocabulary dictionary might not contain all words in the text shown above.

```
dictionary = {
    "the": 0,
    "square": 1,
    "hot": 2,
    "is": 3,
    "not": 4,
    "a": 5,
    "happy": 6,
    "sandwich": 7
}
```

¹You can read more about Word2vec on this blog (or the original research paper for the enthusiasts!)

1. Word Embedding Based Representation

(a) Word embeddings are reduced dimension vector representations (features) of words. Given a single word in the dictionary, word embeddings can convert it to a vector of fixed dimension. In Homework 4, we will provide a dictionary file specifying pre-computed mappings between every word in Vocab and their corresponding word embeddings. To facilitate better understanding towards word embeddings, we produce a plot showing the spatial relationship between several sample words from the vocabulary used in Homework 4, with their corresponding word embeddings (reduced to 2D vectors from 300D vectors using a technique called PCA we will learn about later in this course!):

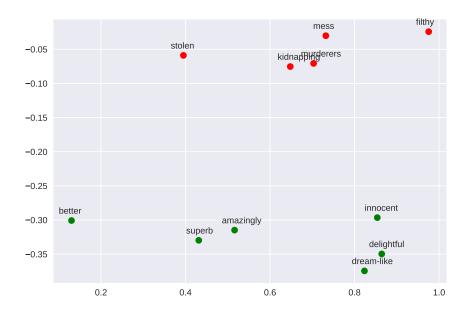


Figure 1: Visualization of word embeddings. We select a few positive words (shown in green) and a few negative words (shown in red). To make the plot, we map the highdimensional word representations of these words to 2D space using PCA and then visualize them in the scatter plot above.

Please comment on your observations and findings based on this plot.

(b) Now, we much translate these word embeddings to sentence embeddings (a vector representing the sentence as a whole). One approach to building a sentence embedding is to average out the vector representation of every word in the sentence that is in the dictionary. For example, given text "a hot dog flies like a sandwich", we can find the sentence embedding for this text by taking the average of the vector representation of the words "a", "hot", "a", and "sandwich".

Now suppose we have the following word embedding dictionary for building sentence embeddings (this is a toy example used for illustrative purposes; actual word embeddings will have higher dimensions than this example):

```
dictionary = {
    "the": [0.2, 0.3],
    "square": [0.8, 0.9],
    "hot": [0.1, -0.2],
    "is": [0.1, 0.1],
    "not": [-0.2, -0.3],
    "a": [0.0, 0.0],
    "happy": [0.4, 0.4],
    "sandwich": [0.2, -0.3]
}
```

Write the word embedding based representation of the sample text define above, repeated here for convenience:

```
a hot dog is not a sandwich because it is not square (2)
```

4 Gradient Descent and Stochastic Gradient Descent

Now we will compare two different optimization methods using pseudocode. Consider a model with parameter $\boldsymbol{\theta} \in \mathbb{R}^M$ being trained with a design matrix $\mathbf{X} \in \mathbb{R}^{N \times M}$ and labels $\mathbf{y} \in \mathbb{R}^N$. Say we update $\boldsymbol{\theta}$ using the objective function $J(\boldsymbol{\theta}|\mathbf{X},\mathbf{y}) = \frac{1}{N} \sum_{i=1}^{N} J^{(i)}(\boldsymbol{\theta}|\mathbf{x}^{(i)},y^{(i)}) \in \mathbb{R}$. Recall that an epoch refers to one complete cycle through the dataset.

(a) Complete the pseudocode for gradient descent.

```
def dJ(theta, X, y, i):
    (omitted) # Returns \partial J^{(i)}(\boldsymbol{\theta}|\mathbf{x}^{(i)},y^{(i)})/\partial \boldsymbol{\theta}
    # You may call this function in your pseudocode.

def GD(theta, X, y, learning_rate):
    for epoch in range(num_epoch):
        Complete this section with the update rule
    return theta # return the updated theta
```

(b) Complete the pseudocode for stochastic gradient descent that samples without replacement.

5 The Need For Speed: Vectorization and Numpy

Performing mathematical operations on vectors and matrices is ubiquitous in most machine learning algorithms. Whether it's a simple similarity measure that works by calculating the dot product between two vectors, or deep neural networks, they all involve repeated matrix operations. This makes it imperative that our underlying code design to perform matrix operations is efficient.

5.1 The Perils of Python

While Python is widely the language of choice for machine learning researchers across the globe (thanks to the speed of development and code readability it offers and the support it enjoys from the open-source community), Python as a high-level language on average is much slower than a lower level language like C++. To combat this, libraries like numpy and scipy implement most of the back-end operations they perform in C/C++, while providing wrappers in Python to be able to call underlying C code seamlessly from a Python script.

5.2 Speed Comparison: Numpy and Python

We highly recommend you to use *numpy* extensively in this course, it will be difficult to pass the programming portion of Homework 4 without writing most of your matrix operations in numpy. In this section, we'll see why. Please refer to this Google Colab notebook for demos and example operations.

Consider you have two vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$. To see how similar they are, as measured by the cosine angle between them, you want to compute their dot product. This translates to the following operation:

$$\mathbf{a} \cdot \mathbf{b} = a_1 b_1 + a_2 b_2 + ... + a_n b_n$$

When translated to code, notice how the dot product in NumPy is a whopping 100x faster than the native Python!

```
from timeit import timeit
import numpy as np
import array

VECTOR_SIZE = int(1e8)

# NumPy arrays
a = np.random.rand(VECTOR_SIZE)
b = np.random.rand(VECTOR_SIZE)

# Python arrays
aArr = array.array('d', a)
bArr = array.array('d', b)
```

```
def test_np():
    return np.dot(a, b)

# faster than multiprocessing, python lists, or numpy arrays with python loops
# faster than using a range and indexing
def test_py_arr():
    return sum(x * y for x, y in zip(aArr, bArr))

def time_dot_product(f):
    return timeit(f, setup=f, number=5) / 5

if __name__ == "__main__":
    print(f"NumPy = {time_dot_product(test_np):.2f}") # 0.05s
    print(f"Python on an array = {time_dot_product(test_py_arr):.2f}") # 5.45s
```

5.3 Useful Numpy Operations

Some operations in numpy that you will find really useful in your assignments are:

- np.matmul: Matrix multiplication of two matrices
- np.unique: Returns unique elements along an axis.
- np.hstack: Stack two arrays horizontally (column-wise)
- np.expand_dims: Convert a row vector of size n into a matrix of size n * 1 or 1 * n
- np.log, np.sum, np.exp, and so on...

Other than the Colab notebook, you can also read these two tutorials (beginner, intermediate) from the official numpy website. For instance, understanding broadcasting is recommended. It will help you debug the shape errors you might face in all future homeworks.

6 Logistic Regression: Toy Example

Let's go through a toy problem.

Y	X_1	X_2	X_3
1	1	2	1
1	1	1	-1
0	1	-2	1

(a) What is
$$J(\boldsymbol{\theta})$$
 of above data given initial $\boldsymbol{\theta} = \begin{bmatrix} -2 \\ 2 \\ 1 \end{bmatrix}$?

- (b) Calculate $\frac{\partial J^{(1)}(\boldsymbol{\theta})}{\partial \theta_1}$, $\frac{\partial J^{(1)}(\boldsymbol{\theta})}{\partial \theta_2}$ and $\frac{\partial J^{(1)}(\boldsymbol{\theta})}{\partial \theta_3}$ for first training example. Note that $\sigma(3) \approx 0.95$.
- (c) Calculate $\frac{\partial J^{(2)}(\boldsymbol{\theta})}{\partial \theta_1}$, $\frac{\partial J^{(2)}(\boldsymbol{\theta})}{\partial \theta_2}$ and $\frac{\partial J^{(2)}(\boldsymbol{\theta})}{\partial \theta_3}$ for second training example. Note that $\sigma(-1) \approx 0.25$.
- (d) Assuming we are doing stochastic gradient descent with a learning rate of 1.0, what are the updated parameters θ if we update θ using the second training example?
- (e) What is the new $J(\theta)$ after doing the above update? Should it decrease or increase?
- (f) Given a test example where $(X_1 = 1, X_2 = 3, X_3 = 4)$, what will the classifier output following this update?