

1. Introduction

Alice, Bob, and Daniel are friends learning machine learning together. After watching a few lectures, they are very proud of having learned many useful tools, including linear and logistic regression, non-linear features, regularization, and kernel tricks. To see how these methods can be used to solve a real life problem, they decide to get their hands dirty with the famous digit recognition problem using the MNIST (Mixed National Institute of Standards and Technology) database.

Hearing that you are an excellent student in the MITx machine learning class with solid understanding of the material and great coding ability in Python, they decide to invite you to their team and help them with implementing these different algorithms.

The MNIST database contains binary images of handwritten digits commonly used to train image processing systems. The digits were collected from among Census Bureau employees and high school students. The database contains 60,000 training digits and 10,000 testing digits, all of which have been size-normalized and centered in a fixed-size image of 28×28 pixels. Many methods have been tested with this dataset and in this project, you will get a chance to experiment with the task of classifying these images into the correct digit using some of the methods you have learned so far.

Setup:

As with the last project, please use Python's **NumPy** numerical library for handling arrays and array operations; use **matplotlib** for producing figures and plots.

This project will be split in two parts. Project 2 (this project) consists in the first part and project 3 will cover the second part.

1. *Note on software:* For all the projects, we will use python 3.6 augmented with the **NumPy** numerical toolbox, the **matplotlib** plotting toolbox. In this project, we will also use the **scikit-learn** package, which you could install in the same way you installed other packages, as described in [project 0](#), e.g. by `conda install scikit-learn` or `pip install sklearn`
2. Download [mnist.tar.gz](#) and untar it in to a working directory. The archive contains the various data files in the Dataset directory, along with the following python files:
 - `part1/linear_regression.py` where you will implement linear regression
 - `part1/svm.py` where you will implement support vector machine
 - `part1/softmax.py` where you will implement multinomial regression
 - `part1/features.py` where you will implement principal component analysis (PCA) dimensionality reduction
 - `part1/kernel.py` where you will implement polynomial and Gaussian RBF kernels
 - `part1/main.py` where you will use the code you write for this part of the project

Important: The archive also contains files for the second part of the MNIST project. For this project, you will only work with the part1 folder.

To get warmed up to the MNIST data set run `python main.py`. This file provides code that reads the data from `mnist.pkl.gz` by calling the function `get_MNIST_data` that is provided for you in `utils.py`. The call to `get_MNIST_data` returns Numpy arrays:

1. `train_x` : A matrix of the training data. Each row of `train_x` contains the features of one image, which are simply the raw pixel values flattened out into a vector of length $784=28^2$.
2. `train_y` : The labels for each training datapoint, aka the digit shown in the corresponding image (a number between 0-9).
3. `test_x` : A matrix of the test data, formatted like `train_x`.
4. `test_y` : The labels for the test data, which should only be used to evaluate the accuracy of different classifiers in your report.

Next, we call the function `plot_images` to display the first 20 images of the training set. Look at these images and get a feel for the data (don't include these in your write-up).

Tip: Throughout the whole online grading system, you can assume the NumPy python library is already imported as `np`. In some problems you will also have access to python's random library, and other functions you've already implemented. Look out for the "Available Functions" Tip before the codebox, as you did in the last project.

This project will unfold both on MITx and on your local machine. However, we encourage you to first implement the functions locally and run the test scripts to validate basic functionality. Think of the online graders as a submission box to submit your code when it is ready. You should not have to use the online graders to debug your code.

Tip: You may find the tutorial on Scikit-learn in <https://github.com/Varal7/ml-tutorial/blob/master/Part1.ipynb>

2. Linear Regression with Closed Form Solution

After seeing the problem, your classmate Alice immediately argues that we can apply a linear regression model, as the labels are numbers from 0-9, very similar to the example we learned from the lecture. Though being a little doubtful, you decide to have a try and start simple by using the raw pixel values of each image as features.

Alice wrote a skeleton code `run_linear_regression_on_MNIST` in `main.py`, but she needs your help to complete the code and make the model work.

Closed Form Solution of Linear Regression

To solve the linear regression problem, you recall the linear regression has a closed form solution:

$$\theta = (X^T X + \lambda I)^{-1} X^T Y$$

Where I is the identity matrix.

Write a function `closed_form` that computes this closed form solution given the features X , labels Y and the regularization parameter λ .

Available Functions: You have access to the NumPy python library as `np`; No need to import anything.

```
1 def closed_form(X, Y, lambda_factor):
2     """
3     Computes the closed form solution of linear regression with L2 regularization
4
5     Args:
6         X - (n, d + 1) NumPy array (n datapoints each with d features plus the bias feature in the
7         Y - (n, ) NumPy array containing the labels (a number from 0-9) for each
8         data point
9         lambda_factor - the regularization constant (scalar)
10    Returns:
11        theta - (d + 1, ) NumPy array containing the weights of linear regression. Note that theta[
12        represents the y-axis intercept of the model and therefore X[0] = 1
13    """
14    # YOUR CODE HERE
15    raise NotImplementedError
```

Test Error on Linear Regression

Apply the linear regression model on the test set. For classification purpose, you decide to round the predicted label into numbers 0-9.

Note: For this project we will be looking at the error rate defined as the fraction of labels that don't match the target labels, also known as the "gold labels" or ground truth. (In other context, you might want to consider other performance measures such as [precision and recall](#), which we have not discussed in this course).

Please enter the **test error** of your linear regression algorithm for different

λ

(copy the output from the main.py run).

Grading note: The tolerance of this problem has been lowered to accept answers from both the correct and incorrect versions of run_linear_regression_on_MNIST.

Error | $\lambda=1$ =

Error | $\lambda=0.1$ =

Error | $\lambda=0.01$ =

What went Wrong?

Alice and you find that no matter what λ factor you try, the test error is large. With some thinking, you realize that something is wrong with this approach.

- A. Gradient descent should be used instead of the closed form solution.
- B. The loss function related to the closed-form solution is inadequate for this problem.
- C. Regularization should not be used here.

3. Support Vector Machine

Bob thinks it is clearly not a regression problem, but a classification problem. He thinks that we can change it into a binary classification and use the support vector machine we learned in Lecture 4 to solve the problem. In order to do so, he suggests that we can build an one vs. rest model for every digit. For example, classifying the digits into two classes: 0 and not 0.

Bob wrote a function run_svm_one_vs_rest_on_MNIST where he changed the labels of digits 1-9 to 1 and keeps the label 0 for digit 0. He also found that sklearn package contains an SVM model that you can use directly. He gave you the link to this model and hopes you can tell him how to use that.

You will be working in the file part1/svm.py in this problem

Important: For this problem, you will need to use the [scikit-learn](https://scikit-learn.org/) library. If you don't have it, install it using pip install sklearn

One vs. Rest SVM

Use the sklearn package and build the SVM model on your local machine. Use random_state = 0, C=0.1 and default values for other parameters.

Available Functions: You have access to the sklearn's implementation of the linear SVM as LinearSVC; No need to import anything.

```
1 def one_vs_rest_svm(train_x, train_y, test_x):
2     """
3     Trains a linear SVM for binary classification
4
5     Args:
6         train_x - (n, d) NumPy array (n datapoints each with d features)
7         train_y - (n, ) NumPy array containing the labels (0 or 1) for each training data point
8         test_x - (m, d) NumPy array (m datapoints each with d features)
9     Returns:
10        pred_test_y - (m,) NumPy array containing the labels (0 or 1) for each test data point
11    """
12    raise NotImplementedError
13
```

Binary classification error

Report the test error by running run_svm_one_vs_rest_on_MNIST.

Error =

Implement C-SVM

Play with the C parameter of SVM, what statement is true about the C parameter?
(Choose all that apply.)

- A. Larger C gives larger tolerance of violation.
- B. Larger C gives smaller tolerance of violation.
- C. Larger C gives a larger-margin separating hyperplane.
- D. Larger C gives a smaller-margin separating hyperplane.

Multiclass SVM

In fact, sklearn already implements a multiclass SVM with a one-vs-rest strategy. Use LinearSVC to build a multiclass SVM model

Available Functions: You have access to the sklearn's implementation of the linear SVM as LinearSVC; No need to import anything.

```

1 def multi_class_svm(train_x, train_y, test_x):
2     """
3     Trains a linear SVM for multiclass classification using a one-vs-rest strategy
4
5     Args:
6         train_x - (n, d) NumPy array (n datapoints each with d features)
7         train_y - (n, ) NumPy array containing the labels (int) for each training data point
8         test_x - (m, d) NumPy array (m datapoints each with d features)
9     Returns:
10        pred_test_y - (m,) NumPy array containing the labels (int) for each test data point
11    """
12    raise NotImplementedError
13

```

Multiclass SVM error

Report the overall test error by running `run_multiclass_svm_on_MNIST`.

Error=

4. Multinomial (Softmax) Regression and Gradient Descent

Daniel suggests that instead of building ten models, we can expand a single logistic regression model into a multinomial regression and solve it with similar gradient descent algorithm.

The main function which you will call to run the code you will implement in this section is `run_softmax_on_MNIST` in `main.py` (already implemented). In the appendix at the bottom of this page, we describe a number of the methods that are already implemented for you in `softmax.py` that will be useful.

In order for the regression to work, you will need to implement three methods. Below we describe what the functions should do. We have included some test cases in `test.py` to help you verify that the methods you have implemented are behaving sensibly.

You will be working in the file `part1/softmax.py` in this problem

Computing Probabilities for Softmax

Write a function `compute_probabilities` that computes, for each data point $x^{(i)}$, the probability that $x^{(i)}$ is labeled as j for $j=0, 1, \dots, k-1$.

The softmax function h for a particular vector x requires computing

$$h(x) = \frac{1}{\sum_{j=0}^{k-1} e^{\theta_j \cdot x / \tau}} \begin{bmatrix} e^{\theta_0 \cdot x / \tau} \\ e^{\theta_1 \cdot x / \tau} \\ \vdots \\ e^{\theta_{k-1} \cdot x / \tau} \end{bmatrix},$$

where $\tau > 0$ is the temperature parameter. When computing the output probabilities (they should always be in the range $[0, 1]$), the terms $e^{\theta_j \cdot x / \tau}$ may be very large or very small, due to the use of the exponential function. This can cause numerical or overflow errors. To deal with this, we can simply subtract some fixed amount c from each exponent to keep the resulting number from getting too large. Since

$$\begin{aligned} h(x) &= \frac{e^{-c}}{e^{-c} \sum_{j=0}^{k-1} e^{\theta_j \cdot x / \tau}} \begin{bmatrix} e^{\theta_0 \cdot x / \tau} \\ e^{\theta_1 \cdot x / \tau} \\ \vdots \\ e^{\theta_{k-1} \cdot x / \tau} \end{bmatrix} \\ &= \frac{1}{\sum_{j=0}^{k-1} e^{[\theta_j \cdot x / \tau] - c}} \begin{bmatrix} e^{[\theta_0 \cdot x / \tau] - c} \\ e^{[\theta_1 \cdot x / \tau] - c} \\ \vdots \\ e^{[\theta_{k-1} \cdot x / \tau] - c} \end{bmatrix}, \end{aligned}$$

subtracting some fixed amount c from each exponent will not change the final probabilities. A suitable choice for this fixed amount is $c = \max_j \theta_j \cdot x / \tau$.

Available Functions: You have access to the NumPy python library as np; No need to import anything.


```

1 def compute_probabilities(X, theta, temp_parameter):
2     """
3     Computes, for each datapoint X[i], the probability that X[i] is labeled as j
4     for j = 0, 1, ..., k-1
5
6     Args:
7         X - (n, d) NumPy array (n datapoints each with d features)
8         theta - (k, d) NumPy array, where row j represents the parameters of our model for label j
9         temp_parameter - the temperature parameter of softmax function (scalar)
10    Returns:
11        H - (k, n) NumPy array, where each entry H[j][i] is the probability that X[i] is labeled as
12    """
13    #YOUR CODE HERE
14    raise NotImplementedError
15

```

Cost Function

Write a function `compute_cost_function` that computes the total cost over every data point.

The cost function $J(\theta)$ is given by: (Use natural log)

$$J(\theta) = -\frac{1}{n} \left[\sum_{i=1}^n \sum_{j=0}^{k-1} [[y^{(i)} == j]] \log \frac{e^{\theta_j \cdot X^{(i)} / \tau}}{\sum_{l=0}^{k-1} e^{\theta_l \cdot X^{(i)} / \tau}} \right] + \frac{\lambda}{2} \sum_{j=0}^{k-1} \sum_{i=0}^{d-1} \theta_{ji}^2$$

Available Functions: You have access to the NumPy python library as `np` and the previous function as `compute_probabilities`

```

1 def compute_cost_function(X, Y, theta, lambda_factor, temp_parameter):
2     """
3     Computes the total cost over every datapoint.
4
5     Args:
6         X - (n, d) NumPy array (n datapoints each with d features)
7         Y - (n, ) NumPy array containing the labels (a number from 0-9) for each
8           data point
9         theta - (k, d) NumPy array, where row j represents the parameters of our
10        model for label j
11        lambda_factor - the regularization constant (scalar)
12        temp_parameter - the temperature parameter of softmax function (scalar)
13
14    Returns
15        c - the cost value (scalar)

```


Gradient Descent

Solution to this problem available before due date: The function `run_gradient_descent_iteration` is necessary for the rest of the project. Hence, once you have either submitted the correct function or finished your attempts for this problem, the solution to this function will be available. We are sorry we acted this late, but hope that it will still help.

Now, in order to run the gradient descent algorithm to minimize the cost function, we need to take the derivative of $J(\theta)$ wrt a particular θ_m . Notice that within $J(\theta)$, we have:

$$\frac{e^{\theta_j \cdot x^{(i)} / \tau}}{\sum_{l=0}^{k-1} e^{\theta_l \cdot x^{(i)} / \tau}} = p(y^{(i)} = j | x^{(i)}, \theta)$$

so we first compute: $\frac{\partial p(y^{(i)} = j | x^{(i)}, \theta)}{\partial \theta_m}$,
when $m = j$,

$$\frac{\partial p(y^{(i)} = j | x^{(i)}, \theta)}{\partial \theta_m} = \frac{x^{(i)}}{\tau} p(y^{(i)} = m | x^{(i)}, \theta) [1 - p(y^{(i)} = m | x^{(i)}, \theta)]$$

when $m \neq j$,

$$\frac{\partial p(y^{(i)} = j | x^{(i)}, \theta)}{\partial \theta_m} = -\frac{x^{(i)}}{\tau} p(y^{(i)} = m | x^{(i)}, \theta) p(y^{(i)} = j | x^{(i)}, \theta)$$

Now we compute

$$\begin{aligned} \frac{\partial}{\partial \theta_m} \left[\sum_{j=0}^{k-1} [[y^{(i)} == j]] \log \frac{e^{\theta_j \cdot x^{(i)} / \tau}}{\sum_{l=0}^{k-1} e^{\theta_l \cdot x^{(i)} / \tau}} \right] &= \sum_{j=0, j \neq m}^{k-1} \left[[[y^{(i)} == j]] \left[-\frac{x^{(i)}}{\tau} p(y^{(i)} = m | x^{(i)}, \theta) \right] \right. \\ &\quad \left. + [[y^{(i)} == m]] \frac{x^{(i)}}{\tau} [1 - p(y^{(i)} = m | x^{(i)}, \theta)] \right] \\ &= \frac{x^{(i)}}{\tau} \left[[[y^{(i)} == m]] - p(y^{(i)} = m | x^{(i)}, \theta) \sum_{j=0}^{k-1} [[y^{(i)} == j]] \right] \\ &= \frac{x^{(i)}}{\tau} \left[[[y^{(i)} == m]] - p(y^{(i)} = m | x^{(i)}, \theta) \right] \end{aligned}$$

Plug this into the derivative of $J(\theta)$, we have

$$\begin{aligned}\frac{\partial J(\theta)}{\partial \theta_m} &= \frac{\partial}{\partial \theta_m} \left[-\frac{1}{n} \left[\sum_{i=1}^n \sum_{j=0}^{k-1} [[y^{(i)} == j]] \log p(y^{(i)} = j | x^{(i)}, \theta) \right] + \frac{\lambda}{2} \sum_{j=0}^{k-1} \sum_{i=0}^{d-1} \theta_{ji}^2 \right] \\ &= -\frac{1}{n} \sum_{i=1}^n [x^{(i)} ([[y^{(i)} == m]] - p(y^{(i)} = m | x^{(i)}, \theta))] + \lambda \theta_m\end{aligned}$$

To run gradient descent, we will update θ at each step with $\theta \leftarrow \theta - \alpha \nabla J(\theta)$, where α is the learning rate.

Write a function `run_gradient_descent_iteration` that runs one step of the gradient descent algorithm.

Available Functions: You have access to the NumPy python library as `np`, `compute_probabilities` which you previously implemented and `scipy.sparse` as `sparse`

You should use `sparse.coo_matrix` so that your function can handle larger matrices efficiently (and not time out for the online graders). The sparse matrix representation can handle sparse matrices efficiently.

Hint

This is how to use `scipy`'s `sparse.coo_matrix` function to create a sparse matrix of 0's and 1's:

```
M = sparse.coo_matrix(([1]*n, (Y, range(n))), shape=(k,n)).toarray()
```

This will create a normal numpy array with 1s and 0s.

On larger inputs (i.e. MNIST), this is 10x faster than using a naive for loop. (See example code if interested).

Note: As a personal challenge, try to see if you can use special numpy functions to add 1 in-place. This would be even faster.

```
import time
import numpy as np
import scipy.sparse as sparse

ITER = 100
K = 10
N = 10000

def naive(indices, k):
    mat = [[1 if i == j else 0 for j in range(k)] for i in indices]
    return np.array(mat).T

def with_sparse(indices, k):
    n = len(indices)
    M = sparse.coo_matrix(([1]*n, (Y, range(n))), shape=(k,n)).toarray()
    return M

Y = np.random.randint(0, K, size=N)

t0 = time.time()
for i in range(ITER):
    naive(Y, K)
print(time.time() - t0)

t0 = time.time()
for i in range(ITER):
    with_sparse(Y, K)
print(time.time() - t0)
```

```
1 def run_gradient_descent_iteration(X, Y, theta, alpha, lambda_factor, temp_parameter):
2     """
3     Runs one step of batch gradient descent
4
5     Args:
6         X - (n, d) NumPy array (n datapoints each with d features)
7         Y - (n, ) NumPy array containing the labels (a number from 0-9) for each
8             data point
9         theta - (k, d) NumPy array, where row j represents the parameters of our
10             model for label j
11         alpha - the learning rate (scalar)
12         lambda_factor - the regularization constant (scalar)
13         temp_parameter - the temperature parameter of softmax function (scalar)
14
15     Returns:
```

Test Error on Softmax Regression

Finally, report the final test error by running the main.py file, using the temperature parameter $\tau=1$. If you have implemented everything correctly, the error on the test set should be around 0.1, which implies the linear softmax regression model is able to recognize MNIST digits with around 90 percent accuracy.

Note: For this project we will be looking at the error rate defined as the fraction of labels that don't match the target labels, also known as the "gold labels" or ground truth. (In other contexts, you might want to consider other performance measures such as precision and recall, which we have not discussed in this class.

Please enter the test error of your Softmax algorithm (copy the output from the main.py run).

5. Temperature

We will now explore the effects of the temperature parameter in our algorithm.

You will be working in the files `part1/main.py` and `part1/softmax.py` in this problem

Effects of Adjusting Temperature

Explain how the temperature parameter affects the probability of a sample $x(i)$ being assigned a label that has a large θ . What about a small θ ?

- A. Larger temperature leads to less variance
- B. Smaller temperature leads to less variance
- C. Smaller temperature makes the distribution more uniform

Reporting Error Rates

Set the temperature parameter to be 0.5, 1, and 2; re-run `run_softmax_on_MNIST` for each one of these (add your code to the specified part in `main.py`).

Error| $T=0.5$ =

Error| $T=1$ =

Error| $T=2$ =

6. Changing Labels

We now wish to classify the digits by their (mod 3) value, such that the new label $y(i)$ of sample i is the old $y(i) \pmod{3}$. (Reminder: Return the `temp_parameter` to be 1 if you changed it for the last section)

You will be working in the file `part1/main.py` and `part1/softmax.py` in this problem

Using the Current Model - update target

Given that we already classified every $x(i)$ as a digit, we could use the model we already trained and just calculate our estimations (mod 3).

Implement `update_y` function, which changes the old digit labels for the training and test set for the new (mod 3) labels.

Available Functions: You have access to the NumPy python library as `np`

```
1 def update_y(train_y, test_y):
2     """
3     Changes the old digit labels for the training and test set for the new (mod 3)
4     labels.
5
6     Args:
7         train_y - (n, ) NumPy array containing the labels (a number between 0-9)
8                 for each datapoint in the training set
9         test_y - (n, ) NumPy array containing the labels (a number between 0-9)
10                for each datapoint in the test set
11
12     Returns:
13         train_y_mod3 - (n, ) NumPy array containing the new labels (a number between 0-2)
14                       for each datapoint in the training set
15         test_y_mod3 - (n, ) NumPy array containing the new labels (a number between 0-2)
```

Using the Current Model - compute test error

Implement `compute_test_error_mod3` function, which takes the test points X , their correct labels Y (digits (mod 3) from 0-2), θ , and the `temp_parameter`, and returns the error.

	Estimated Y	Estimated Y (mod 3)	Correct Y	Correct Y (mod 3)
x_1	9	0	8	2
x_2	6	0	6	0
x_3	5	2	8	2

The error of the regression with the original labels would be 0.66667

However, the error of the regression when comparing the (mod 3) of the labels would be 0.33333

Available Functions: You have access to the NumPy python library as `np` and to the `get_classification` function from the project release

```

1 def compute_test_error_mod3(X, Y, theta, temp_parameter):
2     """
3     Returns the error of these new labels when the classifier predicts the digit. (mod 3)
4
5     Args:
6         X - (n, d - 1) NumPy array (n datapoints each with d - 1 features)
7         Y - (n, ) NumPy array containing the labels (a number from 0-2) for each
8           data point
9         theta - (k, d) NumPy array, where row j represents the parameters of our
10              model for label j
11         temp_parameter - the temperature parameter of softmax function (scalar)
12
13     Returns:
14         test_error - the error rate of the classifier (scalar)
15     """

```

Using the Current Model - test error

Find the error rate of the new labels (call these two functions at the end of `run_softmax_on_MNIST`). See the functions' documentation for detailed explanations of the inputs and outputs.

Error rate for labels mod 3:

Retrain with New Labels

Now suppose that instead we want to retrain our classifier with the new labels. In other words, rather than training the model to predict the original digits and then taking those predictions modulo 3, we explicitly train the model to predict the digits modulo 3 from the original image. How do you expect the performance to change using the new labels?

[Increase](#)
[Decrease](#)
[Stay the same](#)

Implement `run_softmax_on_MNIST_mod3` in **main.py** to perform this new training; report the new error rate.

[Error rate when trained on labels mod 3](#)

7. Classification Using Manually Crafted Features

The performance of most learning algorithms depends heavily on the representation of the training data. In this section, we will try representing each image using different features in

place of the raw pixel values. Subsequently, we will investigate how well our regression model from the previous section performs when fed different representations of the data.

Dimensionality Reduction via PCA

Principal Components Analysis (PCA) is the most popular method for linear dimension reduction of data and is widely used in data analysis. For an in-depth exposition see: http://snobear.colorado.edu/Markw/BioMath/Otis/PCA/principal_components.ps.

Briefly, this method finds (orthogonal) directions of maximal variation in the data. By projecting an $n \times d$ dataset X onto $k \leq d$ of these directions, we get a new dataset of lower dimension that reflects more variation in the original data than any other k -dimensional linear projection of X . By going through some linear algebra, it can be proven that these directions are equal to the k eigenvectors corresponding to the k largest eigenvalues of the covariance matrix $\tilde{X}^T \tilde{X}$, where \tilde{X} is a centered version of our original data.

Remark: The best implementations of PCA actually use the Singular Value Decomposition of \tilde{X} rather than the more straightforward approach outlined here, but these concepts are beyond the scope of this course.

Cubic Features

In this section, we will also work with a cubic feature mapping which maps an input vector $x = [x_1, \dots, x_d]$ into a new feature vector $\phi(x)$, defined so that for any $x, x' \in \mathbb{R}^d$:

$$\phi(x)^T \phi(x') = (x^T x' + 1)^3$$

8. Dimensionality Reduction Using PCA

PCA finds (orthogonal) directions of maximal variation in the data. In this problem we're going to project our data onto the principal components and explore the effects on performance. **You will be working in the files `part1/main.py` and `part1/features.py` in this problem**

Project onto Principal Components

Fill in function `project_onto_PC` in `features.py` that implements PCA dimensionality reduction of dataset X .

Note that to project a given $n \times d$ dataset X into its k -dimensional PCA representation, one can use matrix multiplication, after first centering X :

$$\tilde{X}V$$

where \tilde{X} is the centered version of the original data X and V is the $d \times k$ matrix whose columns are the top k eigenvectors of $\tilde{X}^T \tilde{X}$. This is because the eigenvectors are of unit-norm, so there is no need to divide by their length.

You are given the full principal component matrix V' as `pcs` in this function.

Available Functions: You have access to the NumPy python library as `np` and the function `center_data` which returns a centered version of the data, where each feature now has mean = 0

Correction on test.py:: We have updated the test for this function in test.py to remove some ambiguity. Please download a corrected version of the project archive `mnist.tar.gz` or correct the function `check_project_onto_PC` in test.py to the following:

```
def check_project_onto_PC():
    ex_name = "Project onto PC"
    X = np.array([
        [1, 2, 3],
        [2, 4, 6],
        [3, 6, 9],
        [4, 8, 12],
    ])
    pcs = features.principal_components(X)
    exp_res = np.array([
        [5.61248608, 0],
        [1.87082869, 0],
        [-1.87082869, 0],
        [-5.61248608, 0],
    ])
    n_components = 2
    if check_array(
        ex_name, features.project_onto_PC,
        exp_res, X, pcs, n_components):
        return
    log(green("PASS"), ex_name, "")
```

```
1 def project_onto_PC(X, pcs, n_components, feature_means):
2     """
3     Given principal component vectors pcs = principal_components(X)
4     this function returns a new data array in which each sample in X
5     has been projected onto the first n_components principal components.
6     """
7     # TODO: first center data using the feature_means
8     # TODO: Return the projection of the centered dataset
9     #         on the first n_components principal components.
10    #         This should be an array with dimensions: n x n_components.
11    # Hint: these principal components = first n_components columns
12    #         of the eigenvectors returned by principal_components().
13    #         Note that each eigenvector is already be a unit-vector,
14    #         so the projection may be done using matrix multiplication.
15    raise NotImplementedError
```

9. Cubic Features

In this section, we will work with a cubic feature mapping which maps an input vector $x=[x_1, \dots, x_d]$ into a new feature vector $\phi(x)$, defined so that for any $x, x' \in \mathbb{R}^d$:

$$\phi(x)^T \phi(x') = (x^T x' + 1)^3$$

You will be working in the files `part1/main.py` and `part1/features.py` in this problem

Computing Cubic Features

In 2-D, let $x=[x_1, x_2]$. Write down the explicit cubic feature mapping $\phi(x)$ as a vector; i.e., $\phi(x)=[f_1(x_1, x_2), \dots, f_N(x_1, x_2)]$

Hint

Hint: $\phi(x)$ should be a 10-dimensional vector.

The `cubic_features` function in `features.py` is already implemented for you. That function can handle input with an arbitrary dimension and compute the corresponding features for the cubic Kernel. Note that here we don't leverage the kernel properties that allow us to do a more efficient computation with the kernel function (without computing the features themselves). Instead, here we do compute the cubic features explicitly and apply the PCA on the output features.

Applying to MNIST

If we explicitly apply the cubic feature mapping to the original 784-dimensional raw pixel features, the resulting representation would be of massive dimensionality. Instead, we will apply the quadratic feature mapping to the 10-dimensional PCA representation of our training data which we will have to calculate just as we calculated the 18-dimensional representation in the previous problem. After applying the cubic feature mapping to the PCA representations for both the train and test datasets, retrain the softmax regression model using these new features and report the resulting test set error below.

Important: You will probably get a runtime warning for getting the log of 0, ignore. Your code should still run and perform correctly.

Note: Use the same training parameters as the first softmax model given in `main.py` file and temperature 1.

If you have done everything correctly, softmax regression should perform better (on the test set) using these features than either the 18-dimensional principal components or raw pixels. The error on the test set using cubic features should only be around 0.08, demonstrating the power of nonlinear classification models.

Error rate for 10-dimensional cubic PCA features =

10. Kernel Methods

As you can see, implementing a direct mapping to the high-dimensional features is a lot of work (imagine using an even higher dimensional feature mapping.) This is where the kernel trick becomes useful.

Recall the kernel perceptron algorithm we learned in the lecture. The weights θ can be represented by a linear combination of features:

$$\theta = \sum_{i=1}^n \alpha^{(i)} y^{(i)} \phi(x^{(i)})$$

In the softmax regression formulation, we can also apply this representation of the weights:

$$\theta_j = \sum_{i=1}^n \alpha_j^{(i)} y^{(i)} \phi(x^{(i)}) .$$

$$h(x) = \frac{1}{\sum_{j=1}^k e^{[\theta_j \cdot \phi(x)] - c}} \begin{bmatrix} e^{[\theta_1 \cdot \phi(x)] - c} \\ e^{[\theta_2 \cdot \phi(x)] - c} \\ \vdots \\ e^{[\theta_k \cdot \phi(x)] - c} \end{bmatrix}$$

$$h(x) = \frac{1}{\sum_{j=1}^k e^{[\sum_{i=1}^n \alpha_j^{(i)} y^{(i)} \phi(x^{(i)}) \cdot \phi(x)] - c}} \begin{bmatrix} e^{[\sum_{i=1}^n \alpha_1^{(i)} y^{(i)} \phi(x^{(i)}) \cdot \phi(x)] - c} \\ e^{[\sum_{i=1}^n \alpha_2^{(i)} y^{(i)} \phi(x^{(i)}) \cdot \phi(x)] - c} \\ \vdots \\ e^{[\sum_{i=1}^n \alpha_k^{(i)} y^{(i)} \phi(x^{(i)}) \cdot \phi(x)] - c} \end{bmatrix}$$

We actually do not need the real mapping $\phi(x)$, but the inner product between two features after mapping: $\phi(x_i) \cdot \phi(x)$, where x_i is a point in the training set and x is the new data point for which we want to compute the probability. If we can create a kernel function $K(x,y)=\phi(x) \cdot \phi(y)$, for any two points x and y , we can then kernelize our softmax regression algorithm.

You will be working in the files `part1/main.py` and `part1/kernel.py` in this problem

Implementing Polynomial Kernel

In the last section, we explicitly created a cubic feature mapping. Now, suppose we want to map the features into d dimensional polynomial space,

$$\phi(x) = \langle x_d^2, \dots, x_1^2, \sqrt{2}x_dx_{d-1}, \dots, \sqrt{2}x_dx_1, \sqrt{2}x_{d-1}x_{d-2}, \dots, \sqrt{2}x_{d-1}x_1, \dots, \sqrt{2}x_2x_1, \sqrt{2c}x_d, \dots, \sqrt{2c}x_1, c \rangle$$

Write a function `polynomial_kernel` that takes in two matrix X and Y and computes the polynomial kernel $K(x,y)$ for every pair of rows x in X and y in Y .

```

1
2 def polynomial_kernel(X, Y, c, p):
3     """
4         Compute the polynomial kernel between two matrices X and Y::
5             K(x, y) = (<x, y> + c)^p
6         for each pair of rows x in X and y in Y.
7
8     Args:
9         X - (n, d) NumPy array (n datapoints each with d features)
10        Y - (m, d) NumPy array (m datapoints each with d features)
11        c - a coefficient to trade off high-order and low-order terms (scalar)
12        p - the degree of the polynomial kernel
13
14    Returns:
15        kernel_matrix - (n, m) Numpy array containing the kernel matrix

```

Gaussian RBF Kernel

Another commonly used kernel is the Gaussian RBF kernel. Similarly, write a function `rbf_kernel` that takes in two matrices X and Y and computes the RBF kernel $K(x,y)$ for every pair of rows x in X and y in Y .

Available Functions: You have access to the NumPy python library as `np`

```

1
2 def rbf_kernel(X, Y, gamma):
3     """
4     Compute the Gaussian RBF kernel between two matrices X and Y::
5         K(x, y) = exp(-gamma ||x-y||^2)
6     for each pair of rows x in X and y in Y.
7
8     Args:
9         X - (n, d) NumPy array (n datapoints each with d features)
10        Y - (m, d) NumPy array (m datapoints each with d features)
11        gamma - the gamma parameter of gaussian function (scalar)
12
13    Returns:
14        kernel_matrix - (n, m) Numpy array containing the kernel matrix
15    """

```

Now, try implementing the softmax regression using kernelized features. You will have to rewrite the `softmax_regression` function in `softmax.py`, as well as the auxiliary functions `compute_cost_function`, `compute_probabilities`, `run_gradient_descent_iteration`. How does the test error change?

In this project, you have been familiarized with the MNIST dataset for digit recognition, a popular task in computer vision.

You have implemented a linear regression which turned out to be inadequate for this task. You have also learned how to use scikit-learn's SVM for binary classification and multiclass classification.

Then, you have implemented your own softmax regression using gradient descent.

Finally, you experimented with different hyperparameters, different labels and different features, including kernelized features.

In the next project, you will apply neural networks to this task.