

1. (10 points) **Noisy linear regression**

Brayr Tilban Elberier

A real estate company have assigned us the task of building a model to predict the house prices in Westwood. For this task, the company has provided us with a dataset \mathcal{D} :

$$\mathcal{D} = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$$

where $x^{(i)} \in \mathbb{R}^d$ is a feature vector of the i^{th} house and $y^{(i)} \in \mathbb{R}$ is the price of the i^{th} house. Since we just learned about linear regression, so we have decided to use a linear regression model for this task. Additionally, the IT manager of the real estate company has requested us to design a model with small weights. In order to accommodate his request, we will design a linear regression model with parameter regularization. In this problem, we will navigate through the process of achieving regularization by adding noise to the feature vectors. Recall, that we define the cost function in a linear regression problem as:

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - (x^{(i)})^T \theta)^2$$

where $\theta \in \mathbb{R}^d$ is the parameter vector. As mentioned earlier, we will be adding noise to the feature vectors in the dataset. Specifically, we will be adding zero-mean gaussian noise of known variance σ^2 from the distribution

$$\mathcal{N}(0, \sigma^2 \mathbf{I})$$

where $\mathbf{I} \in \mathbb{R}^{d \times d}$ and $\sigma \in \mathbb{R}$. With the addition of gaussian noise the modified cost function is given by,

$$\tilde{\mathcal{L}}(\theta) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - (x^{(i)} + \delta^{(i)})^T \theta)^2$$

where $\delta^{(i)}$ are i.i.d noise vectors with $\delta^{(i)} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$.

- (a) (6 points) Express the expectation of the modified loss over the gaussian noise, $\mathbb{E}_{\delta \sim \mathcal{N}}[\tilde{\mathcal{L}}(\theta)]$, in terms of the original loss plus a term independent of the dataset \mathcal{D} . To be precise, your answer should be of the form:

$$\mathbb{E}_{\delta \sim \mathcal{N}}[\tilde{\mathcal{L}}(\theta)] = \mathcal{L}(\theta) + R$$

where R is not a function of \mathcal{D} . For answering this part, you might find the following result useful:

$$\mathbb{E}_{\delta \sim \mathcal{N}}[\delta \delta^T] = \sigma^2 \mathbf{I}$$

- (b) (2 points) Based on your answer to (a), under expectation what regularization effect would the addition of the noise have on the model?
- (c) (1 point) Suppose $\sigma \rightarrow 0$, what effect would this have on the model?
- (d) (1 point) Suppose $\sigma \rightarrow \infty$, what effect would this have on the model?

$$\textcircled{1} a. E[\mathcal{L}(\theta)]$$

$$= E\left[\frac{1}{N} \sum_{i=1}^N (y^{(i)} - (x^{(i)} + f^{(i)})^T \theta)^2\right]$$

$$= \frac{1}{N} \sum_{i=1}^N E[(y^{(i)} - (x^{(i)} + f^{(i)})^T \theta)^2]$$

$$= \underbrace{(y^{(i)} - x^{(i)T} \theta)^2}_{\text{blue}} - \underbrace{2(y^{(i)} - x^{(i)T} \theta)(f^{(i)T} \theta)}_{\text{red}} + \underbrace{(f^{(i)T} \theta)^2}_{\text{green}}$$

$$E[(y^{(i)} - x^{(i)T} \theta)^2] = (y^{(i)} - x^{(i)T} \theta)^2$$

$$E[-2(y^{(i)} - x^{(i)T} \theta)(f^{(i)T} \theta)]$$

$$= -2(y^{(i)} - x^{(i)T} \theta) E[f^{(i)T} \theta]$$

$$= 0$$

$$E[(f^{(i)T} \theta)^2] = E[(f^{(i)T} \theta)^T (f^{(i)T} \theta)]$$

$$= E[\theta^T f^{(i)} f^{(i)T} \theta]$$

$$= \theta^T E[f^{(i)} f^{(i)T}] \theta = \theta^T \sigma^2 I \theta$$

$$= \sigma^2 \theta^T \theta = \sigma^2 \|\theta\|_2^2$$

$$E_{\theta \sim N}[\mathcal{L}(\theta)] = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - x^{(i)T} \theta)^2 + \sigma^2 \|\theta\|_2^2$$

↳ mean square error ↳ L2 reg.

$$E_{f \sim N[\tilde{L}(\theta)]} = \mathcal{L}(\theta) + \sigma^2 \|\theta\|_2^2$$

- b. It would have L2 regularization effect
- c. Close to no effect at all, overfit?
- d. Model could be heavily regularized and potentially underfit.

3. (30 points) **Softmax classifier gradient.** For softmax classifier, derive the gradient of the log likelihood.

Concretely, assume a classification problem with c classes

- Samples are $(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})$, where $\mathbf{x}^{(j)} \in \mathbb{R}^n$, $y^{(j)} \in \{1, \dots, c\}$, $j = 1, \dots, m$
- Parameters are $\theta = \{\mathbf{w}_i, b_i\}_{i=1, \dots, c}$
- Probabilistic model is

$$\Pr(y^{(j)} = i \mid \mathbf{x}^{(j)}, \theta) = \text{softmax}_i(\mathbf{x}^{(j)})$$

where

$$\text{softmax}_i(\mathbf{x}) = \frac{e^{\mathbf{w}_i^T \mathbf{x} + b_i}}{\sum_{k=1}^c e^{\mathbf{w}_k^T \mathbf{x} + b_k}}$$

Derive the log-likelihood \mathcal{L} , and its gradient w.r.t. the parameters, $\nabla_{\mathbf{w}_i} \mathcal{L}$ and $\nabla_{b_i} \mathcal{L}$, for $i = 1, \dots, c$.

Note: We can group \mathbf{w}_i and b_i into a single vector by augmenting the data vectors with an additional dimension of constant 1. Let $\tilde{\mathbf{x}} = \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$, $\tilde{\mathbf{w}}_i = \begin{bmatrix} \mathbf{w}_i \\ b_i \end{bmatrix}$, then $a_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + b_i = \tilde{\mathbf{w}}_i^T \tilde{\mathbf{x}}$.

This unifies $\nabla_{\mathbf{w}_i} \mathcal{L}$ and $\nabla_{b_i} \mathcal{L}$ into $\nabla_{\tilde{\mathbf{w}}_i} \mathcal{L}$.

⑤. Log likelihood ...

from lecture ... $\text{softmax}_i(\mathbf{x}) = \frac{e^{a_i(\mathbf{x})}}{\sum_{j=1}^c e^{a_j(\mathbf{x})}}$

for $a_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + b_i$ *

$c = \#$ of classes

if we let $\Theta = \{\mathbf{w}_i, b_i\}_{i=1, \dots, c}$, then $\text{softmax}_i(\mathbf{x})$ can be interpreted as the probability that belongs to class i

$\Pr(y^{(j)} = i \mid \mathbf{x}^{(j)}, \theta) = \text{softmax}_i(\mathbf{x}^{(j)})$

$\rightarrow \mathcal{L} = \log \prod_{i=1}^m \Pr(y^{(i)} \mid \mathbf{x}^{(i)}, \theta)$

$$= \log \prod_{i=1}^m \text{softmax}_{y^{(i)}}(x^{(i)})$$

$$= \log \prod_{i=1}^m \left[\frac{e^{w_{y^{(i)}}^T x^{(i)} + b_{y^{(i)}}}}{\sum_{j=1}^C e^{w_j^T x^{(i)} + b_j}} \right]$$

$$= \log \prod_{i=1}^m \left[\frac{e^{(a_{y^{(i)}}(x^{(i)}))}}{\sum_{j=1}^C e^{(a_j(x^{(i)}))}} \right]$$

$$= \sum_{i=1}^m \left[a_{y^{(i)}}(x^{(i)}) - \log \sum_{j=1}^C e^{(a_j(x^{(i)}))} \right]$$

Gradient ↴

$$\nabla_{w_i} L = \frac{dL}{dw_i} = \sum_{j=1}^m x^{(j)} (1 \{y^{(j)} = i\} - \text{softmax}_i(x^{(j)}))$$

$$\nabla_{b_i} L = \frac{dL}{db_i} = \sum_{j=1}^m (1 \{y^{(j)} = i\} - \text{softmax}_i(x^{(j)}))$$

$$\text{softmax}_i(x^{(i)}) = \frac{e^{(a_i(x^{(i)}))}}{\sum_{k=1}^C e^{(a_k(x^{(i)}))}}$$

4. (10 points) **Hinge loss gradient.**

Owing to the drastic changes in climate throughout the world, a weather forecasting organization wants our help to build a model that can classify the observed weather patterns as severe or not severe. They have accumulated data on various attributes of the weather pattern such as temperature, precipitation, humidity, wind speed, air pressure, and geographical location along with severity of weather. However, the contribution of the attributes to the classification of weather as severe or not is unknown.

We choose to use a binary support vector machine (SVM) classification model. The SVM model parameters are learned by optimizing a hinge loss. The company has provided us with a data-set

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(K)}, y^{(K)})\}$$

where $\mathbf{x}^{(i)} \in \mathbb{R}^d$ is a feature vector of the i^{th} data sample and $y^{(i)} \in \{-1, 1\}$. We define the hinge loss per training sample as

$$\text{hinge}_{y^{(i)}}(\mathbf{x}^{(i)}) = \max(0, 1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b)) \quad (1)$$

, where $\mathbf{w} \in \mathbb{R}^d$ and bias $b \in \mathbb{R}$ are the model parameters. With the hinge loss per sample defined, we can then formulate the average loss for our model as:

$$\mathcal{L}(\mathbf{w}, b) = \frac{1}{K} \sum_{i=1}^K \text{hinge}_{y^{(i)}}(\mathbf{x}^{(i)}) \quad (2)$$

Find the gradient of the loss function $\mathcal{L}(\mathbf{w}, b)$ with respect to the parameters i.e $\nabla_{\mathbf{w}} \mathcal{L}$ and $\nabla_b \mathcal{L}$.

Hint: An Indicator function, also known as a characteristic function, takes on the value of 1 at certain designated points and 0 at all other points. Mathematically, we can represent it as follows:

$$\mathbb{1}_{\{p < 1\}} = \begin{cases} 1, & \text{if } p < 1 \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

$$\begin{aligned} \text{hinge}_{y^{(i)}}(\mathbf{x}^{(i)}) &= \max(0, 1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b)) \\ \mathcal{L}(\mathbf{w}, b) &= \frac{1}{K} \sum_{i=1}^K \text{hinge}_{y^{(i)}}(\mathbf{x}^{(i)}) \end{aligned}$$

$$\left[\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L} &= \frac{1}{K} \sum_{i=1}^K \nabla_{\mathbf{w}} \max(0, 1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b)) \\ \text{when } 1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) &> 0 \\ \text{when } 1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) &\leq 0, \mathcal{L} = 0 \end{aligned} \right]$$

$\nabla_{\omega} L$ follow same process as above \uparrow

$$\nabla_{\omega} L = \frac{1}{K} \sum_{i=1}^K -y^{(i)} x^{(i)} \mid \sum 1 - y^{(i)} (\omega^T x^{(i)} + b) > 0 \mid$$

$$\nabla_b L = \frac{1}{K} \sum_{i=1}^K -y^{(i)} \mid \sum 1 - y^{(i)} (\omega^T x^{(i)} + b) > 0 \mid$$

knn_nosol

January 30, 2024

0.1 This is the k-nearest neighbors workbook for ECE C147/C247 Assignment #2

Please follow the notebook linearly to implement k-nearest neighbors.

Please print out the workbook entirely when completed.

The goal of this workbook is to give you experience with the data, training and evaluating a simple classifier, k-fold cross validation, and as a Python refresher.

0.2 Import the appropriate libraries

```
[ ]: import numpy as np # for doing most of our calculations
import matplotlib.pyplot as plt # for plotting
from utils.data_utils import load_CIFAR10 # function to load the CIFAR-10
↳ dataset.

# Load matplotlib images inline
%matplotlib inline

# These are important for reloading any code you write in external .py files.
# see http://stackoverflow.com/questions/1907993/
↳ autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2
```

The autoreload extension is already loaded. To reload it, use:

```
%reload_ext autoreload
```

```
[ ]: # Set the path to the CIFAR-10 data
cifar10_dir = '/Users/tilboon/Documents/UCLA/Courses/C247/HW2/HW2_code/
↳ cifar-10-batches-py' # You need to update this line
X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

# As a sanity check, we print out the size of the training and test data.
print('Training data shape: ', X_train.shape)
print('Training labels shape: ', y_train.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)
```



```

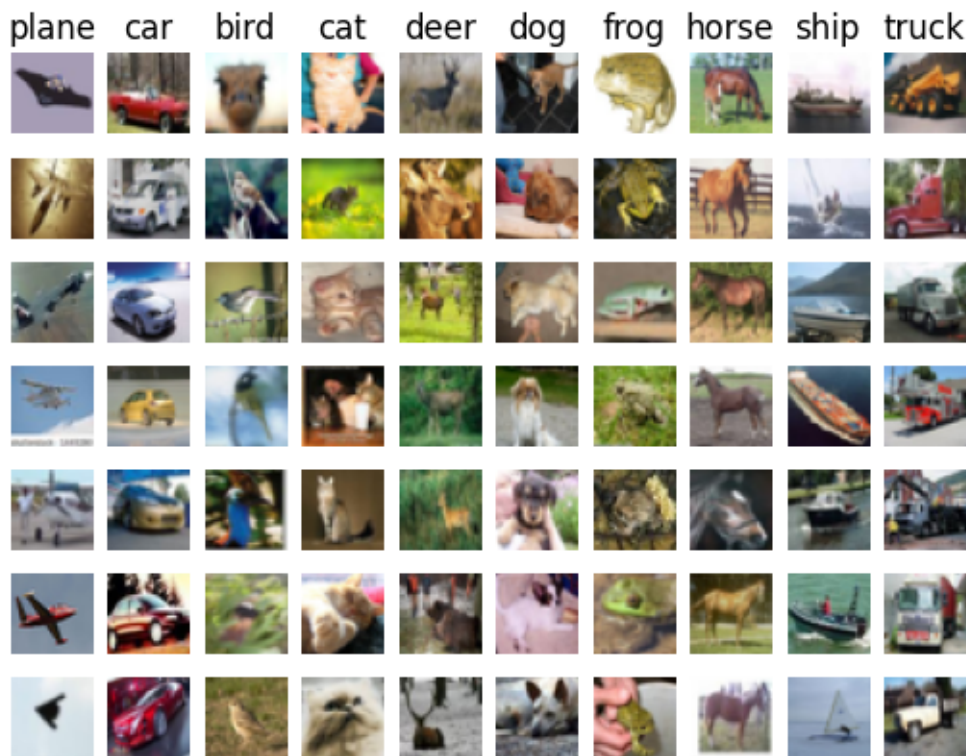
Training data shape: (50000, 32, 32, 3)
Training labels shape: (50000,)
Test data shape: (10000, 32, 32, 3)
Test labels shape: (10000,)

```

```

[ ]: # Visualize some examples from the dataset.
# We show a few examples of training images from each class.
classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship', 'truck']
num_classes = len(classes)
samples_per_class = 7
for y, cls in enumerate(classes):
    idxs = np.flatnonzero(y_train == y)
    idxs = np.random.choice(idxs, samples_per_class, replace=False)
    for i, idx in enumerate(idxs):
        plt_idx = i * num_classes + y + 1
        plt.subplot(samples_per_class, num_classes, plt_idx)
        plt.imshow(X_train[idx].astype('uint8'))
        plt.axis('off')
    if i == 0:
        plt.title(cls)
plt.show()

```



```
[ ]: # Subsample the data for more efficient code execution in this exercise
num_training = 5000
mask = list(range(num_training))
X_train = X_train[mask]
y_train = y_train[mask]

num_test = 500
mask = list(range(num_test))
X_test = X_test[mask]
y_test = y_test[mask]

# Reshape the image data into rows
X_train = np.reshape(X_train, (X_train.shape[0], -1))
X_test = np.reshape(X_test, (X_test.shape[0], -1))
print(X_train.shape, X_test.shape)
```

(5000, 3072) (500, 3072)

1 K-nearest neighbors

In the following cells, you will build a KNN classifier and choose hyperparameters via k-fold cross-validation.

```
[ ]: # Import the KNN class

from nndl import KNN
```

```
[ ]: # Declare an instance of the knn class.
knn = KNN()

# Train the classifier.
# We have implemented the training of the KNN classifier.
# Look at the train function in the KNN class to see what this does.
knn.train(X=X_train, y=y_train)
```

1.1 Questions

- (1) Describe what is going on in the function `knn.train()`.
- (2) What are the pros and cons of this training step?

1.2 Answers

- (1) This is a pretty simple function. It is simply storing the training data in the class instance.
- (2) This is a very simple and easy to implement training function. But, this could lead to a large amount of memory usage and eventually slow down your computations.

1.3 KNN prediction

In the following sections, you will implement the functions to calculate the distances of test points to training points, and from this information, predict the class of the KNN.

```
[ ]: # Implement the function compute_distances() in the KNN class.
# Do not worry about the input 'norm' for now; use the default definition of
    ↳ the norm
# in the code, which is the 2-norm.
# You should only have to fill out the clearly marked sections.

import time
time_start =time.time()

dists_L2 = knn.compute_distances(X=X_test)

print('Time to run code: {}'.format(time.time()-time_start))
print('Frobenius norm of L2 distances: {}'.format(np.linalg.norm(dists_L2,
    ↳ 'fro'))))
```

Time to run code: 45.17984080314636

Frobenius norm of L2 distances: 7906696.077040902

Really slow code Note: This probably took a while. This is because we use two for loops. We could increase the speed via vectorization, removing the for loops.

If you implemented this correctly, evaluating `np.linalg.norm(dists_L2, 'fro')` should return: ~7906696

1.3.1 KNN vectorization

The above code took far too long to run. If we wanted to optimize hyperparameters, it would be time-expensive. Thus, we will speed up the code by vectorizing it, removing the for loops.

```
[ ]: # Implement the function compute_L2_distances_vectorized() in the KNN class.
# In this function, you ought to achieve the same L2 distance but WITHOUT any
    ↳ for loops.
# Note, this is SPECIFIC for the L2 norm.

time_start =time.time()
dists_L2_vectorized = knn.compute_L2_distances_vectorized(X=X_test)
print('Time to run code: {}'.format(time.time()-time_start))
print('Difference in L2 distances between your KNN implementations (should be
    ↳ 0): {}'.format(np.linalg.norm(dists_L2 - dists_L2_vectorized, 'fro')))
```

Time to run code: 0.21989989280700684

Difference in L2 distances between your KNN implementations (should be 0): 0.0

Speedup Depending on your computer speed, you should see a 10-100x speed up from vectorization. On our computer, the vectorized form took 0.36 seconds while the naive implementation

took 38.3 seconds.

1.3.2 Implementing the prediction

Now that we have functions to calculate the distances from a test point to given training points, we now implement the function that will predict the test point labels.

```
[ ]: # Implement the function predict_labels in the KNN class.
# Calculate the training error (num_incorrect / total_samples)
# from running knn.predict_labels with k=1

error = 1

# ===== #
# YOUR CODE HERE:
# Calculate the error rate by calling predict_labels on the test
# data with k = 1. Store the error rate in the variable error.
# ===== #

y_pred = knn.predict_labels(dists_L2_vectorized)
error = np.count_nonzero(y_test-y_pred)/float(len(y_test))

# ===== #
# END YOUR CODE HERE
# ===== #

print(error)
```

0.726

If you implemented this correctly, the error should be: 0.726.

This means that the k-nearest neighbors classifier is right 27.4% of the time, which is not great, considering that chance levels are 10%.

2 Optimizing KNN hyperparameters

In this section, we'll take the KNN classifier that you have constructed and perform cross-validation to choose a best value of k , as well as a best choice of norm.

2.0.1 Create training and validation folds

First, we will create the training and validation folds for use in k-fold cross validation.

```
[ ]: # Create the dataset folds for cross-validation.
num_folds = 5

X_train_folds = []
y_train_folds = []
```

```

# ===== #
# YOUR CODE HERE:
#   Split the training data into num_folds (i.e., 5) folds.
#   X_train_folds is a list, where X_train_folds[i] contains the
#       data points in fold i.
#   y_train_folds is also a list, where y_train_folds[i] contains
#       the corresponding labels for the data in X_train_folds[i]
# ===== #

# Randomization
np.random.seed(0)
indices = np.arange(num_training)
np.random.shuffle(indices)

X_train_shuffled = X_train[indices]
y_train_shuffled = y_train[indices]

X_train_folds = np.array_split(X_train_shuffled, num_folds)
y_train_folds = np.array_split(y_train_shuffled, num_folds)

# ===== #
# END YOUR CODE HERE
# ===== #

```

2.0.2 Optimizing the number of nearest neighbors hyperparameter.

In this section, we select different numbers of nearest neighbors and assess which one has the lowest k-fold cross validation error.

```

[ ]: time_start =time.time()

ks = [1, 2, 3, 5, 7, 10, 15, 20, 25, 30]

# ===== #
# YOUR CODE HERE:
#   Calculate the cross-validation error for each k in ks, testing
#   the trained model on each of the 5 folds. Average these errors
#   together and make a plot of k vs. cross-validation error. Since
#   we are assuming L2 distance here, please use the vectorized code!
#   Otherwise, you might be waiting a long time.
# ===== #

final_errors = []

for k in ks:

```

```

errors = []

for i in range(num_folds):
    # Find validation set
    X_current_validation = X_train_folds[i]
    y_current_validation = y_train_folds[i]

    # Find training set and concatenate folds
    X_current_train = np.concatenate(X_train_folds[:i] + X_train_folds[(i + 1):
↪])
    y_current_train = np.concatenate(y_train_folds[:i] + y_train_folds[(i + 1):
↪])

    # Train model with current 'k' value
    knn = KNN()
    knn.train(X=X_current_train, y=y_current_train)
    current_dists_arr = knn.
↪compute_L2_distances_vectorized(X=X_current_validation)
    y_prediction = knn.predict_labels(current_dists_arr, k=k)

    # Error computation
    errors.append(np.count_nonzero(y_current_validation - y_prediction)/
↪float(len(y_current_validation)))

    # Add average errors to final_errors array
    final_errors.append(np.mean(errors))

    print('K: %d, Avg Error: %f' % (k, final_errors[-1]))

# Plot
plt.plot(ks, final_errors)
plt.xlabel('K')
plt.ylabel('Cross-Validation Error')
plt.show()

# ===== #
# END YOUR CODE HERE
# ===== #

print('Computation time: %.2f'%(time.time()-time_start))

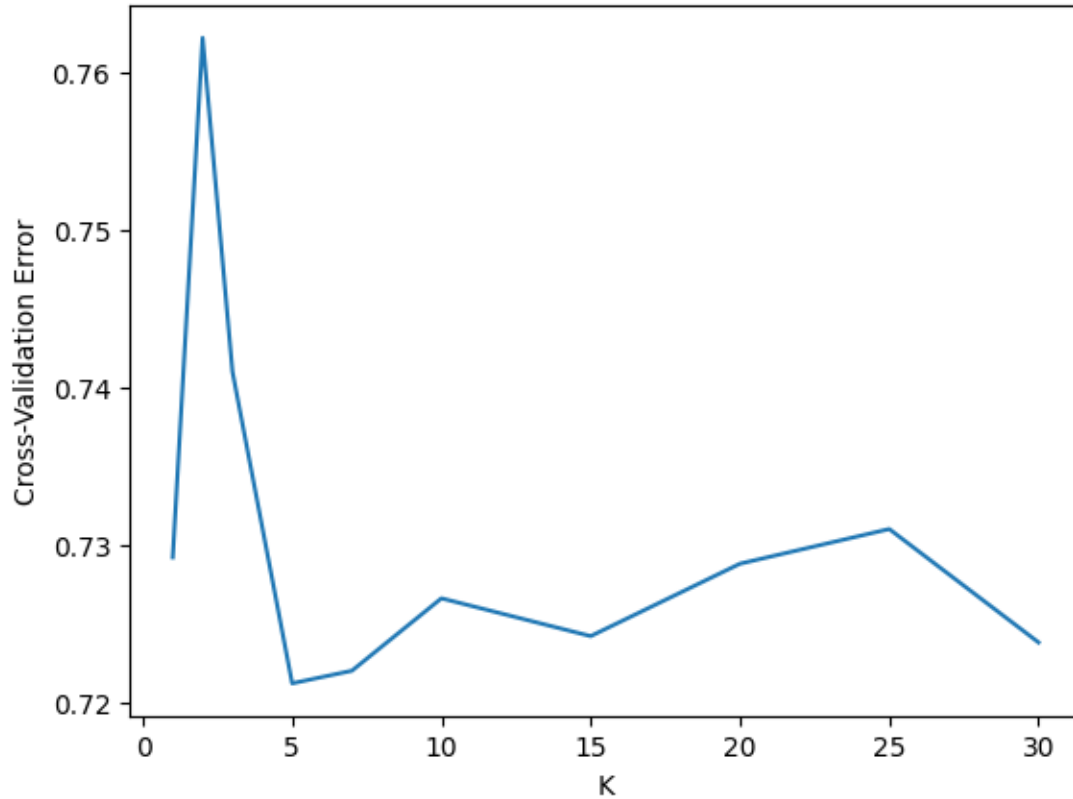
```

```

K: 1, Avg Error: 0.729200
K: 2, Avg Error: 0.762200
K: 3, Avg Error: 0.741000
K: 5, Avg Error: 0.721200
K: 7, Avg Error: 0.722000
K: 10, Avg Error: 0.726600

```

K: 15, Avg Error: 0.724200
K: 20, Avg Error: 0.728800
K: 25, Avg Error: 0.731000
K: 30, Avg Error: 0.723800



Computation time: 20.98

2.1 Questions:

- (1) What value of k is best amongst the tested k 's?
- (2) What is the cross-validation error for this value of k ?

2.2 Answers:

- (1) This value varies dependent on the random seed selected, but for this example (0) - $k = 5$
- (2) This value varies dependent on the random seed selected, but for this example (0) - cross-validation error = 0.721200

2.2.1 Optimizing the norm

Next, we test three different norms (the 1, 2, and infinity norms) and see which distance metric results in the best cross-validation performance.

```

[ ]: time_start =time.time()

L1_norm = lambda x: np.linalg.norm(x, ord=1)
L2_norm = lambda x: np.linalg.norm(x, ord=2)
Linf_norm = lambda x: np.linalg.norm(x, ord= np.inf)
norms = [L1_norm, L2_norm, Linf_norm]

# ===== #
# YOUR CODE HERE:
# Calculate the cross-validation error for each norm in norms, testing
# the trained model on each of the 5 folds. Average these errors
# together and make a plot of the norm used vs the cross-validation error
# Use the best cross-validation k from the previous part.
#
# Feel free to use the compute_distances function. We're testing just
# three norms, but be advised that this could still take some time.
# You're welcome to write a vectorized form of the L1- and Linf- norms
# to speed this up, but it is not necessary.
# ===== #

final_norm_errors = []

for norm in norms:

    errors = []

    for i in range(num_folds):
        # Find validation set
        X_current_validation = X_train_folds[i]
        y_current_validation = y_train_folds[i]

        # Find training set and concatenate folds
        X_current_train = np.concatenate(X_train_folds[:i] + X_train_folds[(i + 1):
↪])
        y_current_train = np.concatenate(y_train_folds[:i] + y_train_folds[(i + 1):
↪])

        # Train model with current 'k' value
        knn = KNN()
        knn.train(X=X_current_train, y=y_current_train)
        current_dists_arr = knn.compute_distances(X_current_validation, norm)
        y_prediction = knn.predict_labels(current_dists_arr, k=5)

        # Error computation
        errors.append(np.count_nonzero(y_current_validation - y_prediction)/
↪float(len(y_current_validation)))

```



```

# Add average errors to final_errors array
final_norm_errors.append(np.mean(errors))

print('L1_norm: Avg Error = %f' % (final_norm_errors[0]))
print('L2_norm: Avg Error = %f' % (final_norm_errors[1]))
print('Linf_norm: Avg Error = %f' % (final_norm_errors[2]))

# Plot
plt.plot(final_norm_errors)
plt.title("Norm vs. Cross-Validation Error")
plt.xlabel('Norm')
plt.ylabel('Cross-Validation Error')
plt.xticks(np.arange(3), ['L1_norm', 'L2_norm', 'Linf_norm'])
plt.show()

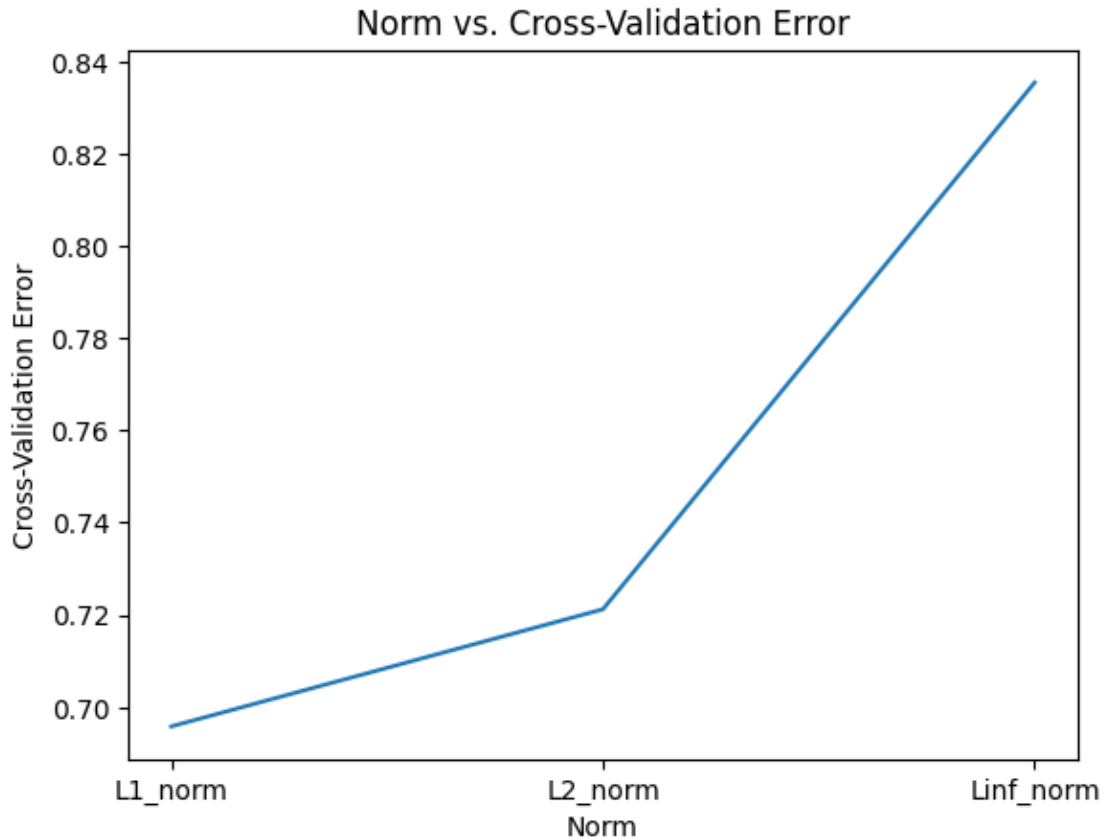
# ===== #
# END YOUR CODE HERE
# ===== #
print('Computation time: %.2f'%(time.time()-time_start))

```

```

L1_norm: Avg Error = 0.695800
L2_norm: Avg Error = 0.721200
Linf_norm: Avg Error = 0.835400

```



Computation time: 531.58

2.3 Questions:

- (1) What norm has the best cross-validation error?
- (2) What is the cross-validation error for your given norm and k ?

2.4 Answers:

- (1) L1 norm
- (2) $k = 5$ and $CVE = 0.695800$

3 Evaluating the model on the testing dataset.

Now, given the optimal k and norm you found in earlier parts, evaluate the testing error of the k -nearest neighbors model.

```
[ ]: error = 1

# ===== #
```

```

# YOUR CODE HERE:
# Evaluate the testing error of the k-nearest neighbors classifier
# for your optimal hyperparameters found by 5-fold cross-validation.
# ===== #

k = 5

knn = KNN()
knn.train(X=X_train, y=y_train)

new_dists = knn.compute_distances(X_test, L1_norm)
y_pred = knn.predict_labels(new_dists, k)
error = np.count_nonzero(y_test - y_pred)/float(len(y_test))

# ===== #
# END YOUR CODE HERE
# ===== #

print('Error rate achieved: {}'.format(error))

```

Error rate achieved: 0.698

3.1 Question:

How much did your error improve by cross-validation over naively choosing $k = 1$ and using the L2-norm?

3.2 Answer:

$0.726 - 0.698 = 0.028 \sim 2.8\%$

```
import numpy as np
import pdb
```

```
class KNN(object):
```

```
    def __init__(self):
        pass
```

```
    def train(self, X, y):
        """
        Inputs:
        - X is a numpy array of size (num_examples, D)
        - y is a numpy array of size (num_examples, )
        """
        self.X_train = X
        self.y_train = y
```

```
    def compute_distances(self, X, norm=None):
        """
        Compute the distance between each test point in X and each training point
        in self.X_train.

        Inputs:
        - X: A numpy array of shape (num_test, D) containing test data.
        - norm: the function with which the norm is taken.

        Returns:
        - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
          is the Euclidean distance between the ith test point and the jth training
          point.
        """
        if norm is None:
            norm = lambda x: np.sqrt(np.sum(x**2))
            #norm = 2

        num_test = X.shape[0]
        num_train = self.X_train.shape[0]
        dists = np.zeros((num_test, num_train))
        for i in np.arange(num_test):

            for j in np.arange(num_train):
                # ===== #
                # YOUR CODE HERE:
                #   Compute the distance between the ith test point and the jth
                #   training point using norm(), and store the result in dists[i, j].
                # ===== #

                dists[i, j] = norm(X[i] - self.X_train[j])

                # ===== #
                # END YOUR CODE HERE
                # ===== #

        return dists
```

```
    def compute_L2_distances_vectorized(self, X):
        """
        Compute the distance between each test point in X and each training point
        in self.X_train WITHOUT using any for loops.

        Inputs:
        - X: A numpy array of shape (num_test, D) containing test data.

        Returns:
        - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
          is the Euclidean distance between the ith test point and the jth training
```

```

        point.
    """
    num_test = X.shape[0]
    num_train = self.X_train.shape[0]
    dists = np.zeros((num_test, num_train))

    # ===== #
    # YOUR CODE HERE:
    # Compute the L2 distance between the ith test point and the jth
    # training point and store the result in dists[i, j]. You may
    # NOT use a for loop (or list comprehension). You may only use
    # numpy operations.
    #
    # HINT: use broadcasting. If you have a shape (N,1) array and
    # a shape (M,) array, adding them together produces a shape (N, M)
    # array.
    # ===== #

    test_squared_sum = X.dot(X).sum(axis=1).reshape((-1, 1))
    train_squared_sum = self.X_train.dot(self.X_train).sum(axis=1).reshape((1, -1))

    test_squared_sum = (X**2).sum(axis=1).reshape((-1, 1))
    train_squared_sum = (self.X_train**2).sum(axis=1).reshape((1, -1))

    test_dot_trainT = X.dot(self.X_train.T)

    dists = np.sqrt(test_squared_sum + train_squared_sum - (2 * test_dot_trainT))

    # ===== #
    # END YOUR CODE HERE
    # ===== #

    return dists

def predict_labels(self, dists, k=1):
    """
    Given a matrix of distances between test points and training points,
    predict a label for each test point.

    Inputs:
    - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
      gives the distance between the ith test point and the jth training point.

    Returns:
    - y: A numpy array of shape (num_test,) containing predicted labels for the
      test data, where y[i] is the predicted label for the test point X[i].
    """
    num_test = dists.shape[0]
    y_pred = np.zeros(num_test)
    for i in np.arange(num_test):
        # A list of length k storing the labels of the k nearest neighbors to
        # the ith test point.
        closest_y = []
        # ===== #
        # YOUR CODE HERE:
        # Use the distances to calculate and then store the labels of
        # the k-nearest neighbors to the ith test point. The function
        # numpy.argsort may be useful.
        #
        # After doing this, find the most common label of the k-nearest
        # neighbors. Store the predicted label of the ith training example
        # as y_pred[i]. Break ties by choosing the smaller label.
        # ===== #

        sortedIdxs = np.argsort(dists[i])
        closest_y = self.y_train[sortedIdxs[:k]]

```

```
y_pred[i] = np.argmax(np.bincount(closest_y))

# ===== #
# END YOUR CODE HERE
# ===== #

return y_pred
```

softmax_nosol

January 30, 2024

0.1 This is the softmax workbook for ECE C147/C247 Assignment #2

Please follow the notebook linearly to implement a softmax classifier.

Please print out the workbook entirely when completed.

The goal of this workbook is to give you experience with training a softmax classifier.

```
[ ]: import random
import numpy as np
from utils.data_utils import load_CIFAR10
import matplotlib.pyplot as plt

%matplotlib inline
%load_ext autoreload
%autoreload 2
```

The autoreload extension is already loaded. To reload it, use:

```
%reload_ext autoreload
```

```
[ ]: def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000,
↪ num_dev=500):
    """
    Load the CIFAR-10 dataset from disk and perform preprocessing to prepare
    it for the linear classifier. These are the same steps as we used for the
    SVM, but condensed to a single function.
    """
    # Load the raw CIFAR-10 data
    cifar10_dir = '/Users/tilboon/Documents/UCLA/Courses/C247/HW2/HW2_code/
↪ cifar-10-batches-py' # You need to update this line
    X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

    # subsample the data
    mask = list(range(num_training, num_training + num_validation))
    X_val = X_train[mask]
    y_val = y_train[mask]
    mask = list(range(num_training))
    X_train = X_train[mask]
    y_train = y_train[mask]
    mask = list(range(num_test))
```

```

X_test = X_test[mask]
y_test = y_test[mask]
mask = np.random.choice(num_training, num_dev, replace=False)
X_dev = X_train[mask]
y_dev = y_train[mask]

# Preprocessing: reshape the image data into rows
X_train = np.reshape(X_train, (X_train.shape[0], -1))
X_val = np.reshape(X_val, (X_val.shape[0], -1))
X_test = np.reshape(X_test, (X_test.shape[0], -1))
X_dev = np.reshape(X_dev, (X_dev.shape[0], -1))

# Normalize the data: subtract the mean image
mean_image = np.mean(X_train, axis = 0)
X_train -= mean_image
X_val -= mean_image
X_test -= mean_image
X_dev -= mean_image

# add bias dimension and transform into columns
X_train = np.hstack([X_train, np.ones((X_train.shape[0], 1))])
X_val = np.hstack([X_val, np.ones((X_val.shape[0], 1))])
X_test = np.hstack([X_test, np.ones((X_test.shape[0], 1))])
X_dev = np.hstack([X_dev, np.ones((X_dev.shape[0], 1))])

return X_train, y_train, X_val, y_val, X_test, y_test, X_dev, y_dev

# Invoke the above function to get our data.
X_train, y_train, X_val, y_val, X_test, y_test, X_dev, y_dev = ↳
get_CIFAR10_data()
print('Train data shape: ', X_train.shape)
print('Train labels shape: ', y_train.shape)
print('Validation data shape: ', X_val.shape)
print('Validation labels shape: ', y_val.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)
print('dev data shape: ', X_dev.shape)
print('dev labels shape: ', y_dev.shape)

```

```

Train data shape: (49000, 3073)
Train labels shape: (49000,)
Validation data shape: (1000, 3073)
Validation labels shape: (1000,)
Test data shape: (1000, 3073)
Test labels shape: (1000,)
dev data shape: (500, 3073)
dev labels shape: (500,)

```


0.2 Training a softmax classifier.

The following cells will take you through building a softmax classifier. You will implement its loss function, then subsequently train it with gradient descent. Finally, you will choose the learning rate of gradient descent to optimize its classification performance.

```
[ ]: from nndl import Softmax

[ ]: # Declare an instance of the Softmax class.
# Weights are initialized to a random value.
# Note, to keep people's first solutions consistent, we are going to use a
    ↪ random seed.

np.random.seed(1)

num_classes = len(np.unique(y_train))
num_features = X_train.shape[1]

softmax = Softmax(dims=[num_classes, num_features])
```

Softmax loss

```
[ ]: ## Implement the loss function of the softmax using a for loop over
# the number of examples

loss = softmax.loss(X_train, y_train)

[ ]: print(loss)
```

2.3277607028048966

0.3 Question:

You'll notice the loss returned by the softmax is about 2.3 (if implemented correctly). Why does this make sense?

0.4 Answer:

$-\ln(0.1) = 2.3$. This shows that the probability of choosing the correct class is 10%, which makes sense because the weights are random.

Softmax gradient

```
[ ]: ## Calculate the gradient of the softmax loss in the Softmax class.
# For convenience, we'll write one function that computes the loss
# and gradient together, softmax.loss_and_grad(X, y)
# You may copy and paste your loss code from softmax.loss() here, and then
# use the appropriate intermediate values to calculate the gradient.

loss, grad = softmax.loss_and_grad(X_dev, y_dev)
```

```
# Compare your gradient to a gradient check we wrote.
# You should see relative gradient errors on the order of 1e-07 or less if you
↳ implemented the gradient correctly.
softmax.grad_check_sparse(X_dev, y_dev, grad)
```

```
numerical: -3.441406 analytic: -3.441406, relative error: 7.371486e-09
numerical: 0.566164 analytic: 0.566164, relative error: 6.158833e-08
numerical: -1.509496 analytic: -1.509496, relative error: 2.205557e-08
numerical: 0.480927 analytic: 0.480927, relative error: 8.133125e-08
numerical: -1.119525 analytic: -1.119526, relative error: 5.768046e-08
numerical: -0.419041 analytic: -0.419041, relative error: 1.806362e-09
numerical: -3.639235 analytic: -3.639235, relative error: 5.895431e-09
numerical: 0.387843 analytic: 0.387843, relative error: 1.610869e-07
numerical: 2.392608 analytic: 2.392608, relative error: 2.146916e-08
numerical: 1.320036 analytic: 1.320036, relative error: 5.961130e-08
```

0.5 A vectorized version of Softmax

To speed things up, we will vectorize the loss and gradient calculations. This will be helpful for stochastic gradient descent.

```
[ ]: import time

[ ]: ## Implement softmax.fast_loss_and_grad which calculates the loss and gradient
# WITHOUT using any for loops.

# Standard loss and gradient
tic = time.time()
loss, grad = softmax.loss_and_grad(X_dev, y_dev)
toc = time.time()
print('Normal loss / grad_norm: {} / {} computed in {}s'.format(loss, np.linalg.
↳ norm(grad, 'fro'), toc - tic))

tic = time.time()
loss_vectorized, grad_vectorized = softmax.fast_loss_and_grad(X_dev, y_dev)
toc = time.time()
print('Vectorized loss / grad: {} / {} computed in {}s'.format(loss_vectorized,
↳ np.linalg.norm(grad_vectorized, 'fro'), toc - tic))

# The losses should match but your vectorized implementation should be much
↳ faster.
print('difference in loss / grad: {} / {} '.format(loss - loss_vectorized, np.
↳ linalg.norm(grad - grad_vectorized)))

# You should notice a speedup with the same output.
```

```
Normal loss / grad_norm: 2.33872090924829 / 321.98589570406955 computed in
0.03820514678955078s
```

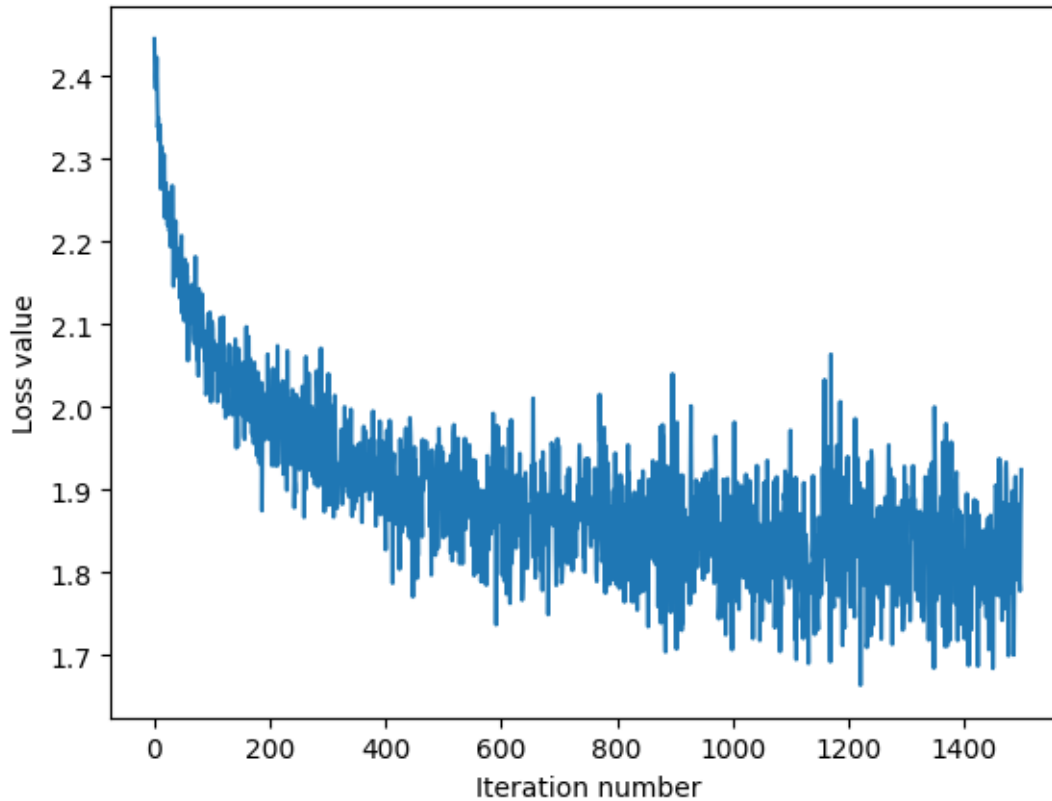
Vectorized loss / grad: 2.338720909248289 / 321.98589570406955 computed in
0.002907991409301758s
difference in loss / grad: 8.881784197001252e-16 / 2.0394590171326893e-13

0.6 Stochastic gradient descent

We now implement stochastic gradient descent. This uses the same principles of gradient descent we discussed in class, however, it calculates the gradient by only using examples from a subset of the training set (so each gradient calculation is faster).

```
[ ]: # Implement softmax.train() by filling in the code to extract a batch of data  
# and perform the gradient step.  
import time  
  
tic = time.time()  
loss_hist = softmax.train(X_train, y_train, learning_rate=1e-7,  
                           num_iters=1500, verbose=True)  
toc = time.time()  
print('That took {}s'.format(toc - tic))  
  
plt.plot(loss_hist)  
plt.xlabel('Iteration number')  
plt.ylabel('Loss value')  
plt.show()
```

```
iteration 0 / 1500: loss 2.44515538257882  
iteration 100 / 1500: loss 2.102722970807078  
iteration 200 / 1500: loss 1.982153270491871  
iteration 300 / 1500: loss 1.9716279917796629  
iteration 400 / 1500: loss 1.8271622721011465  
iteration 500 / 1500: loss 1.8417495237518748  
iteration 600 / 1500: loss 1.8891403759151801  
iteration 700 / 1500: loss 1.8608596252429663  
iteration 800 / 1500: loss 1.8291336459127479  
iteration 900 / 1500: loss 1.832292639851959  
iteration 1000 / 1500: loss 1.8595635065276859  
iteration 1100 / 1500: loss 1.9715495308251019  
iteration 1200 / 1500: loss 1.8029841704072975  
iteration 1300 / 1500: loss 1.8654046891563791  
iteration 1400 / 1500: loss 1.7808285404094022  
That took 3.496454954147339s
```



0.6.1 Evaluate the performance of the trained softmax classifier on the validation data.

```
[ ]: ## Implement softmax.predict() and use it to compute the training and testing
      error.

y_train_pred = softmax.predict(X_train)
print('training accuracy: {}'.format(np.mean(np.equal(y_train,y_train_pred), )))
y_val_pred = softmax.predict(X_val)
print('validation accuracy: {}'.format(np.mean(np.equal(y_val, y_val_pred)), ))
```

```
training accuracy: 0.3812857142857143
validation accuracy: 0.379
```

0.7 Optimize the softmax classifier

```
[ ]: np.finfo(float).eps
```

```
[ ]: 2.220446049250313e-16
```

```

[ ]: # ===== #
# YOUR CODE HERE:
#   Train the Softmax classifier with different learning rates and
#   evaluate on the validation data.
#   Report:
#       - The best learning rate of the ones you tested.
#       - The best validation accuracy corresponding to the best validation error.
#
#   Select the SVM that achieved the best validation error and report
#   its error rate on the test set.
# ===== #

# Array of learning rates
learning_rates = np.logspace(-10, -1, num=30)

best_lr = None
best_val_accuracy = 0
best_softmax = None

for lr in learning_rates:
    softmax = Softmax(dims=[num_classes, num_features])
    softmax.train(X_train, y_train, learning_rate=lr, num_iters=1500,
↳ verbose=False)

    # Prediction
    y_pred = softmax.predict(X_val)

    accuracy = np.mean(y_val == y_pred)
    val_error = 1 - accuracy
    print(f"Learning Rate: {lr}, Validation Accuracy: {accuracy}")

    # Update
    if accuracy > best_val_accuracy:
        best_val_accuracy = accuracy
        best_lr = lr
        best_val_error = val_error
        #best_softmax = softmax

    # Prediction
    y_pred_test = softmax.predict(X_test)
    test_accuracy = np.mean(y_pred_test == y_test)

    best_test_error_rate = 1 - test_accuracy

print(f"Best Learning Rate: {best_lr}")
print(f"Best Validation Accuracy: {best_val_accuracy}")
print(f"Best Validation Error: {best_val_error}")

```

```
print(f"Test Error Rate for the best classifier: {best_test_error_rate}")
```

```
# ===== #  
# END YOUR CODE HERE  
# ===== #
```

```
Learning Rate: 1e-10, Validation Accuracy: 0.1  
Learning Rate: 2.0433597178569396e-10, Validation Accuracy: 0.155  
Learning Rate: 4.175318936560409e-10, Validation Accuracy: 0.133  
Learning Rate: 8.531678524172814e-10, Validation Accuracy: 0.14  
Learning Rate: 1.7433288221999873e-09, Validation Accuracy: 0.19  
Learning Rate: 3.5622478902624368e-09, Validation Accuracy: 0.26  
Learning Rate: 7.278953843983161e-09, Validation Accuracy: 0.284  
Learning Rate: 1.4873521072935118e-08, Validation Accuracy: 0.302  
Learning Rate: 3.039195382313195e-08, Validation Accuracy: 0.341  
Learning Rate: 6.210169418915617e-08, Validation Accuracy: 0.37  
Learning Rate: 1.2689610031679235e-07, Validation Accuracy: 0.397  
Learning Rate: 2.592943797404667e-07, Validation Accuracy: 0.408  
Learning Rate: 5.298316906283712e-07, Validation Accuracy: 0.416  
Learning Rate: 1.0826367338740541e-06, Validation Accuracy: 0.403  
Learning Rate: 2.2122162910704504e-06, Validation Accuracy: 0.41  
Learning Rate: 4.520353656360241e-06, Validation Accuracy: 0.381  
Learning Rate: 9.236708571873865e-06, Validation Accuracy: 0.316  
Learning Rate: 1.8873918221350995e-05, Validation Accuracy: 0.323  
Learning Rate: 3.856620421163472e-05, Validation Accuracy: 0.293  
Learning Rate: 7.880462815669921e-05, Validation Accuracy: 0.274  
Learning Rate: 0.00016102620275609394, Validation Accuracy: 0.267  
Learning Rate: 0.0003290344562312671, Validation Accuracy: 0.241  
Learning Rate: 0.0006723357536499335, Validation Accuracy: 0.352  
Learning Rate: 0.0013738237958832637, Validation Accuracy: 0.273  
Learning Rate: 0.002807216203941181, Validation Accuracy: 0.276  
Learning Rate: 0.005736152510448681, Validation Accuracy: 0.329  
Learning Rate: 0.011721022975334793, Validation Accuracy: 0.319  
Learning Rate: 0.02395026619987491, Validation Accuracy: 0.273  
Learning Rate: 0.04893900918477499, Validation Accuracy: 0.275  
Learning Rate: 0.1, Validation Accuracy: 0.227  
Best Learning Rate: 5.298316906283712e-07  
Best Validation Accuracy: 0.416  
Best Validation Error: 0.5840000000000001  
Test Error Rate for the best classifier: 0.612
```

```
import numpy as np
```

```
class Softmax(object):
```

```
    def __init__(self, dims=[10, 3073]):
        self.init_weights(dims=dims)
```

```
    def init_weights(self, dims):
        """
        Initializes the weight matrix of the Softmax classifier.
        Note that it has shape (C, D) where C is the number of
        classes and D is the feature size.
        """
        self.W = np.random.normal(size=dims) * 0.0001
```

```
    def loss(self, X, y):
        """
        Calculates the softmax loss.

        Inputs have dimension D, there are C classes, and we operate on minibatches
        of N examples.

        Inputs:
        - X: A numpy array of shape (N, D) containing a minibatch of data.
        - y: A numpy array of shape (N,) containing training labels; y[i] = c means
            that X[i] has label c, where 0 <= c < C.

        Returns a tuple of:
        - loss as single float
        """
```

```
    # Initialize the loss to zero.
    loss = 0.0
```

```
    # ===== #
    # YOUR CODE HERE:
    #   Calculate the normalized softmax loss. Store it as the variable loss.
    #   (That is, calculate the sum of the losses of all the training
    #   set margins, and then normalize the loss by the number of
    #   training examples.)
    # ===== #
```

```
    shape = X.shape
    N = shape[0]
```

```
    for i in range(N):
        # Calculate score
        scores = X[i].dot(self.W.T)
        scores -= np.max(scores)

        # Softmax probabilities Computation
        softmax_probs = np.exp(scores) / np.sum(np.exp(scores))

        # Calculate the loss for the correct class
        correct_class = y[i]
        loss += -np.log(softmax_probs[correct_class])
```

```
    # Normalize the loss by the number of training examples
    loss /= N
```

```
    # ===== #
    # END YOUR CODE HERE
    # ===== #
```

```
    return loss
```

```

def loss_and_grad(self, X, y):
    """
    Same as self.loss(X, y), except that it also returns the gradient.

    Output: grad -- a matrix of the same dimensions as W containing
    the gradient of the loss with respect to W.
    """

    # Initialize the loss and gradient to zero.
    loss = 0.0
    grad = np.zeros_like(self.W)

    # ===== #
    # YOUR CODE HERE:
    # Calculate the softmax loss and the gradient. Store the gradient
    # as the variable grad.
    # ===== #

    shape = X.shape
    N = shape[0]

    w_shape = self.W.shape
    K = w_shape[0]

    for i in range(N):
        # Calculate score
        scores = X[i].dot(self.W.T)
        scores -= np.max(scores)

        # Softmax probabilities Computation
        softmax_probs = np.exp(scores) / np.sum(np.exp(scores))

        # Calculate the loss for the correct class
        correct_class = y[i]
        loss += -np.log(softmax_probs[correct_class])

        # Gradient computation
        softmax_probs[correct_class] = (softmax_probs[correct_class] - 1)
        for j in range(K):
            grad[j, :] += softmax_probs[j] * X[i]

    # Normalize the loss by the number of training examples
    loss /= N
    grad /= N

    # ===== #
    # END YOUR CODE HERE
    # ===== #

    return loss, grad

def grad_check_sparse(self, X, y, your_grad, num_checks=10, h=1e-5):
    """
    sample a few random elements and only return numerical
    in these dimensions.
    """

    for i in np.arange(num_checks):
        ix = tuple([np.random.randint(m) for m in self.W.shape])

        oldval = self.W[ix]
        self.W[ix] = oldval + h # increment by h
        fxph = self.loss(X, y)
        self.W[ix] = oldval - h # decrement by h
        fxmh = self.loss(X, y) # evaluate f(x - h)
        self.W[ix] = oldval # reset

```



```

grad_numerical = (fxph - fxmh) / (2 * h)
grad_analytic = your_grad[ix]
rel_error = abs(grad_numerical - grad_analytic) / (abs(grad_numerical) +
abs(grad_analytic))
print('numerical: %f analytic: %f, relative error: %e' % (grad_numerical, grad_analytic,
rel_error))

def fast_loss_and_grad(self, X, y):
    """
    A vectorized implementation of loss_and_grad. It shares the same
    inputs and outputs as loss_and_grad.
    """
    loss = 0.0
    grad = np.zeros(self.W.shape) # initialize the gradient as zero

    # ===== #
    # YOUR CODE HERE:
    # Calculate the softmax loss and gradient WITHOUT any for loops.
    # ===== #

    shape = X.shape
    N = shape[0]

    # Calculate scores
    scores = X.dot(self.W.T)
    scores -= np.max(scores, axis=1, keepdims=True)

    # Calculate softmax prob.
    softmax_probs = np.exp(scores) / np.sum(np.exp(scores), axis=1, keepdims=True)

    # Loss calculation
    loss = (-np.sum(np.log(softmax_probs[np.arange(N), y]))) / N

    # Gradient computation
    softmax_probs[np.arange(N), y] = (softmax_probs[np.arange(N), y] - 1)
    grad = (softmax_probs.T.dot(X) / N)

    # ===== #
    # END YOUR CODE HERE
    # ===== #

    return loss, grad

def train(self, X, y, learning_rate=1e-3, num_iters=100,
          batch_size=200, verbose=False):
    """
    Train this linear classifier using stochastic gradient descent.

    Inputs:
    - X: A numpy array of shape (N, D) containing training data; there are N
        training samples each of dimension D.
    - y: A numpy array of shape (N,) containing training labels; y[i] = c
        means that X[i] has label 0 <= c < C for C classes.
    - learning_rate: (float) learning rate for optimization.
    - num_iters: (integer) number of steps to take when optimizing
    - batch_size: (integer) number of training examples to use at each step.
    - verbose: (boolean) If true, print progress during optimization.

    Outputs:
    A list containing the value of the loss function at each training iteration.
    """
    num_train, dim = X.shape
    num_classes = np.max(y) + 1 # assume y takes values 0...K-1 where K is number of classes

    self.init_weights(dims=[np.max(y) + 1, X.shape[1]]) # initializes the weights of self.W

    # Run stochastic gradient descent to optimize W

```

```

loss_history = []

for it in np.arange(num_iters):
    X_batch = None
    y_batch = None

    # ===== #
    # YOUR CODE HERE:
    # Sample batch_size elements from the training data for use in
    # gradient descent. After sampling,
    # - X_batch should have shape: (batch_size, dim)
    # - y_batch should have shape: (batch_size,)
    # The indices should be randomly generated to reduce correlations
    # in the dataset. Use np.random.choice. It's okay to sample with
    # replacement.
    # ===== #

    shape = X.shape
    N = shape[0]

    index = np.random.choice(N, batch_size)
    X_batch = X[index]
    y_batch = y[index]

    # ===== #
    # END YOUR CODE HERE
    # ===== #

    # evaluate loss and gradient
    loss, grad = self.fast_loss_and_grad(X_batch, y_batch)
    loss_history.append(loss)

    # ===== #
    # YOUR CODE HERE:
    # Update the parameters, self.W, with a gradient step
    # ===== #

    self.W = self.W - (grad * learning_rate)

    # ===== #
    # END YOUR CODE HERE
    # ===== #

    if verbose and it % 100 == 0:
        print('iteration {} / {}: loss {}'.format(it, num_iters, loss))

return loss_history

def predict(self, X):
    """
    Inputs:
    - X: N x D array of training data. Each row is a D-dimensional point.

    Returns:
    - y_pred: Predicted labels for the data in X. y_pred is a 1-dimensional
      array of length N, and each element is an integer giving the predicted
      class.
    """
    y_pred = np.zeros(X.shape[1])
    # ===== #
    # YOUR CODE HERE:
    # Predict the labels given the training data.
    # ===== #

    scores = X.dot(self.W.T)
    y_pred = np.argmax(scores, axis=1)

```

```
# ===== #  
# END YOUR CODE HERE  
# ===== #  
  
return y_pred
```