

## Homework 2

605 796 227.

$$Q1) \quad D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$$

$$x^{(i)} \in \mathbb{R}^d$$

↳ feature vector.

$$y^{(i)} \in \mathbb{R} \rightarrow \text{price}$$

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

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

Noise vectors.

$$\mathbb{E}_{\delta \sim N} \tilde{L}(\theta) = \mathbb{E} \left( \frac{1}{N} \sum_{i=1}^N (y^{(i)} - (x^{(i)} + \delta^{(i)})^T \theta)^2 \right)$$

$$= \mathbb{E} \left( \frac{1}{N} \sum_{i=1}^N (y^{(i)} - x^{(i)T} \theta - \delta^{(i)T} \theta)^2 \right)$$

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

$$= L(\theta) + \mathbb{E} \left( \frac{1}{N} \sum_{i=1}^N 2(y^{(i)} - x^{(i)T} \theta)(\delta^{(i)T} \theta) + (\delta^{(i)T} \theta)^2 \right)$$

$$= L(\theta) + 2(y^{(i)} - x^{(i)T} \theta) \mathbb{E} \left( \delta^{(i)T} \theta \right) + \sum_{i=1}^N \mathbb{E} \left( (\delta^{(i)T} \theta)^2 \right)$$

$\mathbb{E}(\delta^{(i)}) = 0$  as it is the mean.

$$= L(\theta) + \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left( \theta^T \delta^{(i)} \delta^{(i)T} \theta \right) = L(\theta) + \left( \theta^T \mathbb{E} \left( \frac{1}{N} \sum_{i=1}^N \delta^{(i)} \delta^{(i)T} \right) \theta \right)$$

$$= L(\theta) + \frac{1}{N} \sum_{i=1}^N \theta^T \sigma^2 I \theta = \theta^T \theta$$

$$= L(\theta) + \frac{1}{N} \cdot N (\theta^T \theta) \cdot \sigma^2$$



$$\mathbb{E}[\tilde{L}(\theta)] = L(\theta) + \sigma^2 \theta^T \theta$$

b) Under  $L_1$  Norm,  $\mathbb{E}[\tilde{L}(\theta)] = L(\theta) + \lambda \sum_{i=1}^N |w_i|$

$$\|W\|_1 = \sum_{i=1}^p |w_i|$$

for  $L_2$  Norm  $L(\theta) + \lambda \sum (w_i)^2$   $\|W\|_2 = \sqrt{\sum_{i=1}^p (w_i)^2}$

plugging in  $\theta^T \theta$  in the  $L_2$  Norm, it applies regularization  $(\theta^T \theta) = \|\theta\|_2^2$

c) when  $\sigma \rightarrow 0$ , this would make the  $\sigma(\theta^T \theta)$  term zero, thus no regularization is applied. Leads to overfitting.  $L(\theta)$  is the dominant term.

d) when  $\sigma \rightarrow \infty$ , the term  $\sigma(\theta^T \theta) \rightarrow \infty$ , unless  $\theta = 0$ , infinite regularization applied, and it becomes a straight line. (parallel to x axis, no fitting) underfitting. Thus, no learning occurs, as parameter  $\theta$  would minimize loss when  $\theta = 0$ .

(Q3)  $\theta = \{w_i, b_i\}_{i=1 \dots c}$  (= Number of classes)

$$x = \begin{bmatrix} x \\ 1 \end{bmatrix} \quad \tilde{w}_i = \begin{bmatrix} w_i \\ b_i \end{bmatrix}$$

Taking log of softmax

$$= \log \left( \frac{e^{w_i^T x + b_i}}{\sum_{k=1}^c e^{w_k^T x + b_k}} \right)$$

$$= (w_i^T x + b_i) - \log \left( \sum_{k=1}^c e^{w_k^T x + b_k} \right)$$



$$= (\tilde{w}_i^T \tilde{x}) - \log \left( \sum_{k=1}^c e^{\tilde{w}_k^T \tilde{x}} \right)$$

Taking derivative

$$\nabla_{\tilde{w}_i} \mathcal{L} = \frac{\partial}{\partial \tilde{w}_i} \left( \tilde{w}_i^T \tilde{x} - \log \left( \sum_{k=1}^c e^{\tilde{w}_k^T \tilde{x}} \right) \right)$$

$$= \frac{\partial}{\partial \tilde{w}_i} (\tilde{w}_i^T \tilde{x}) - \frac{\partial}{\partial \tilde{w}_i} \log \left( \sum_{k=0}^{p-1} e^{\tilde{w}_k^T \tilde{x}} + e^{\tilde{w}_i^T \tilde{x}} + \sum_{k=i+1}^c e^{\tilde{w}_k^T \tilde{x}} \right)$$

when  $y_j = p$

$$\nabla_{\tilde{w}_i} \mathcal{L} = \tilde{x} - \frac{e^{\tilde{w}_i^T \tilde{x}} \cdot \tilde{x}}{\sum_{k=1}^c e^{\tilde{w}_k^T \tilde{x}}} = \tilde{x} - \tilde{x} \cdot \text{softmax}(\tilde{x})$$

When  $y_j \neq p$ :

$$\nabla_{\tilde{w}_i} \mathcal{L} = - \frac{e^{\tilde{w}_i^T \tilde{x}} \cdot \tilde{x}}{\sum_{k=1}^c e^{\tilde{w}_k^T \tilde{x}}} = - \tilde{x} \cdot \text{softmax}(\tilde{x})$$

$$(Q4) \text{ hinge}_{y^{(i)}}(x^{(i)}) = \max(0, 1 - y^{(i)}(w^T x^{(i)} + b))$$

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

★ when  $y^{(i)}(w^T x^{(i)} + b) > 1$

$$L(w, b) = \frac{1}{K} \sum_{i=1}^K 0$$

$$\text{let } \tilde{w} = \begin{bmatrix} w \\ b \end{bmatrix} \quad x = \begin{bmatrix} x \\ 1 \end{bmatrix}$$

$$\nabla_{\tilde{w}} L = \frac{\partial}{\partial \tilde{w}} \frac{1}{K} \sum_{i=1}^K 0 = \frac{\partial}{\partial \tilde{w}} 0 = \underline{\underline{0}}$$

\* When  $y^{(i)} (\tilde{w}^T \tilde{x}_i) < 1$

$$L = \frac{1}{K} \sum_{i=0}^K 1 - y^{(i)} (\tilde{w}^T \tilde{x}_i)$$

$$\nabla_{\tilde{w}} L = \frac{1}{K} \times K - \frac{1}{K} \sum_{i=0}^K y^{(i)} (\tilde{w}^T \tilde{x}_i)$$

$$= \frac{\partial}{\partial \tilde{w}} \left( 1 - \frac{1}{K} \sum_{i=0}^K y^{(i)} \tilde{w}^T \tilde{x}_i \right) = \frac{\partial a^T x^T b}{\partial x} = b a^T$$

$$\nabla_{\tilde{w}} L = -x^{(i)} y^{(i)}$$

$$\nabla_{\begin{bmatrix} w \\ b \end{bmatrix}} L = - \begin{bmatrix} x^{(i)} \\ 1 \end{bmatrix} \begin{bmatrix} y^{(i)} \end{bmatrix} = - \begin{bmatrix} x^{(i)} y^{(i)} \\ -y^{(i)} \end{bmatrix}$$

$$\begin{cases} \nabla_w L = -x^{(i)} y^{(i)} & \text{if } (1 - y^{(i)} (\tilde{w}^T x^{(i)} + b)) > 0 \\ \nabla_w L = 0 & \text{else} \end{cases}$$

$$\begin{cases} \nabla_b L = -y^{(i)} & \text{if } (1 - y^{(i)} (\tilde{w}^T x^{(i)} + b)) > 0 \\ \nabla_b L = 0 & \text{else} \end{cases}$$



# 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.

## 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
%reload_ext autoreload
%autoreload 2

# Set the path to the CIFAR-10 data
cifar10_dir = './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)

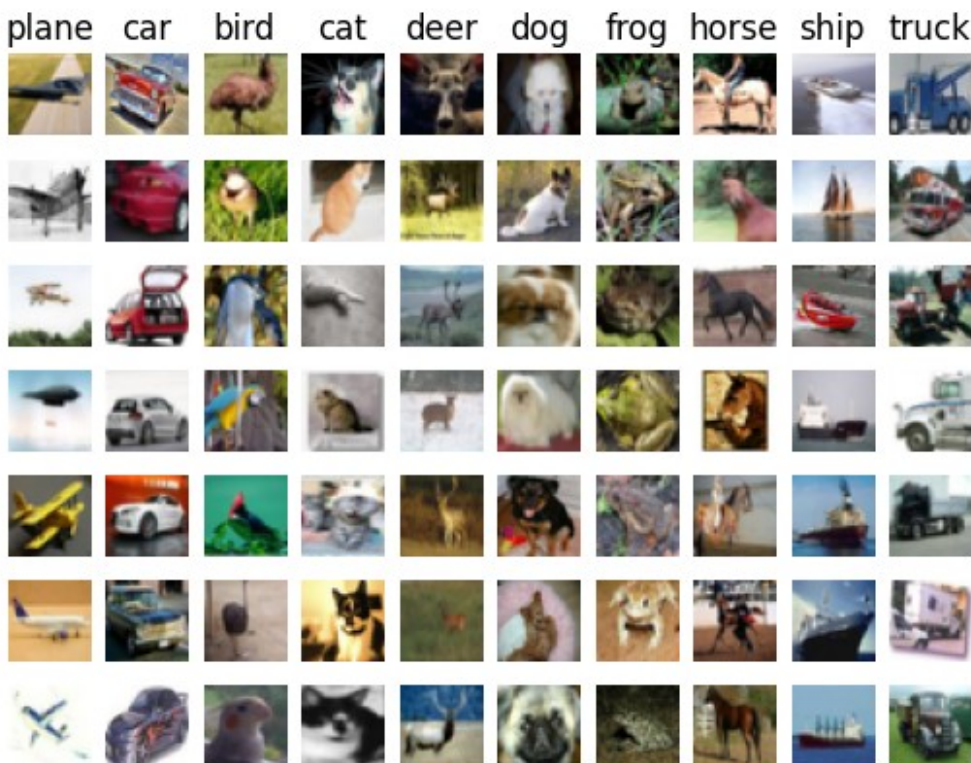
Training data shape: (50000, 32, 32, 3)
Training labels shape: (50000,)
Test data shape: (10000, 32, 32, 3)
Test labels shape: (10000,)

# 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)
```

## 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 nnml 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)
```

## Questions

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

## Answers

- (1) We feed in the inputs `X` and their correct class `y`, used the features in `x` to plot the data points in the feature space. These points are being stored in a vector space.
- (2) This is a very simple training step, as not many computations such as gradient computation, optimizations, etc need to be performed, and thus is not computationally expensive. However the cons of this training is that it requires a lot of space (memory) to store these data points.

## 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: 26.681670904159546
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

## 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: 0.43439197540283203
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.



## 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,1)
errors = (y_test-y_pred)
error = np.count_nonzero(errors)/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%.

## 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.

### 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]
# ===== #
batch = X_train.shape[0]//num_folds
# for i in range(num_folds-1):
#     start_idx = i*batch
#     end_idx = (i+1)*batch
#     X_fold= X_train[start_idx : end_idx]
#     y_fold = y_train[start_idx : end_idx]

#     X_train_folds.append(X_fold)
#     y_train_folds.append(y_fold)
X_train_folds = np.split(X_train, num_folds)
y_train_folds = np.split(y_train, num_folds)
# ===== #
# END YOUR CODE HERE
# ===== #
print(X_train_folds)

[array([[ 59.,  62.,  63., ..., 123.,  92.,  72.],
        [154., 177., 187., ..., 143., 133., 144.],
        [255., 255., 255., ...,  80.,  86.,  84.],
        ...,
        [145., 148., 157., ..., 126., 160.,  91.],
        [146., 146., 146., ..., 238., 238., 238.],
        [203., 206., 208., ..., 132., 131., 126.])], array([[242.,
243., 250., ..., 105., 123., 135.],
        [ 56.,  50.,  28., ..., 131., 112.,  86.],
        [100.,  86.,  89., ...,  44.,  49.,  48.],
        ...,
        [ 41.,  47.,  35., ..., 161., 149.,  89.],
        [ 66., 101., 131., ..., 171., 176., 186.],
        [124., 190., 225., ..., 138., 145., 110.]])], array([[255.,
255., 247., ...,  53.,  52.,  45.],
        [119., 103.,  92., ...,  95., 113., 126.],
        [255., 255., 255., ..., 159., 160., 164.],
        ...,
        [ 29.,  32.,  32., ..., 212., 215., 207.],
        [171., 151., 119., ..., 166., 147., 117.],
        [213., 219., 244., ...,  52.,  54.,  44.]])], array([[254.,
254., 254., ..., 217., 215., 213.],
        [175., 247., 159., ..., 110., 110., 136.],
        [ 91.,  67.,  69., ...,   5.,   2.,   3.],
        ...,
        [164., 150., 127., ..., 161., 138., 103.],
        [228., 236., 240., ...,  92., 105., 113.],

```



```

[ 90., 109., 89., ..., 83., 153., 57.])), array([[ 86.,
138., 179., ..., 75., 123., 163.],
[158., 156., 178., ..., 61., 65., 66.],
[185., 190., 180., ..., 144., 113., 74.],
...,
[167., 163., 145., ..., 42., 78., 84.],
[154., 152., 125., ..., 194., 247., 114.],
[ 45., 32., 21., ..., 156., 142., 100.]])]]

```

## 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]
cross_val_errors = []

def CE_err(num_training, num_foldss, kk, normm = None):
    error = 0
    for i in range(num_foldss):
        X_fold_validation = X_train_folds[i]
        y_fold_validation = y_train_folds[i]
        X_fold_train = np.vstack(np.delete(X_train_folds, i, axis=0))
        y_fold_train = np.hstack(np.delete(y_train_folds, i, axis=0))
        knn.train(X = X_fold_train, y = y_fold_train)

        if(normm):
            dist = knn.compute_distances(X=X_fold_validation,
norm=normm)
        else:
            dist =
knn.compute_L2_distances_vectorized(X=X_fold_validation)
            test_fold_num = num_training/num_foldss
            y_pred = knn.predict_labels(dists = dist, k=kk)
            num_wrong = test_fold_num -
np.count_nonzero(y_pred==y_fold_validation)
            error += num_wrong/test_fold_num
    return error/num_foldss

# ===== #
# 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.
# ===== #

avg_errors = [CE_err(num_training,num_folds, ks[k]) for k in
range(len(ks))]

[print("k = {}, Average cross validation error: {}".format(
    ks[k], avg_errors[k])) for k in range(len(ks))]

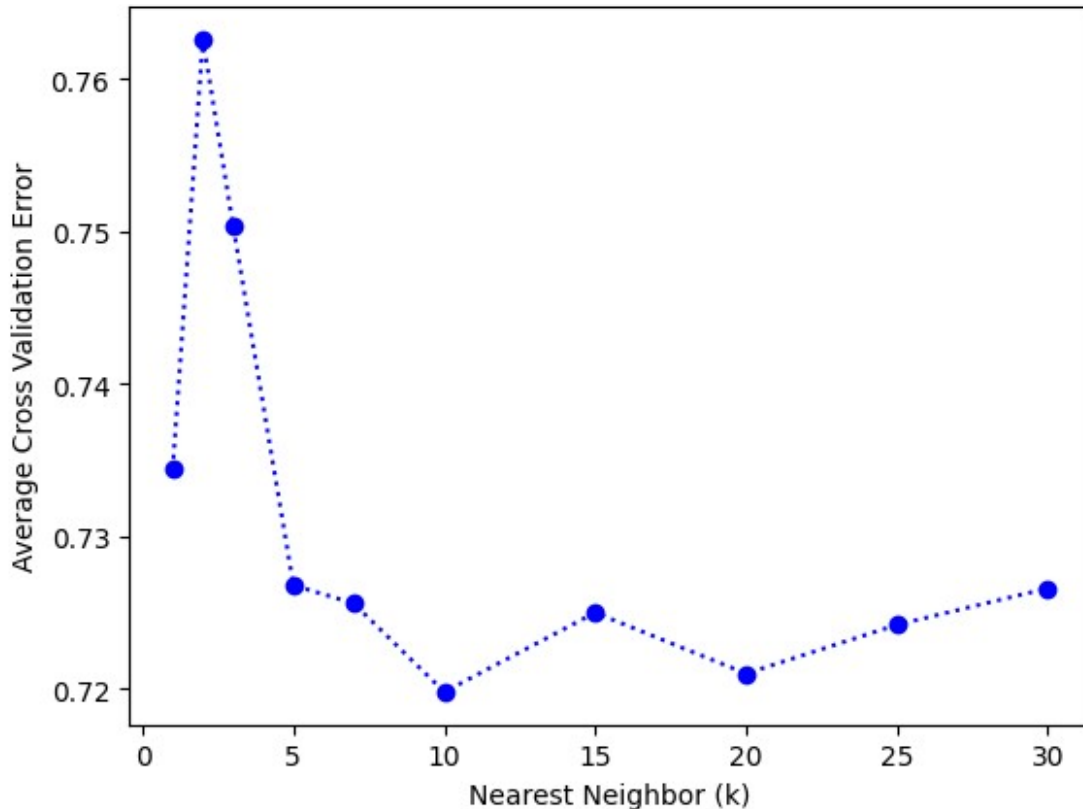
# Plotting the results
fig, ax = plt.subplots(1,1)
ax.plot(ks, avg_errors, "bo:")
ax.set_xlabel("Nearest Neighbor (k)")
ax.set_ylabel("Average Cross Validation Error")
# ===== #
# END YOUR CODE HERE
# ===== #

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

k = 1, Average cross validation error: 0.7344
k = 2, Average cross validation error: 0.76260000000000002
k = 3, Average cross validation error: 0.75040000000000001
k = 5, Average cross validation error: 0.72679999999999999
k = 7, Average cross validation error: 0.7256
k = 10, Average cross validation error: 0.7198
k = 15, Average cross validation error: 0.725
k = 20, Average cross validation error: 0.721
k = 25, Average cross validation error: 0.7242
k = 30, Average cross validation error: 0.7266
Computation time: 57.21

```





## Questions:

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

## Answers:

- (1) Based on the results above from the graph, the best value of  $k$  is 10, which gives the best accuracy.
- (2) The cross validation for when  $k = 10$  is 71.98%

## 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.
# ===== #

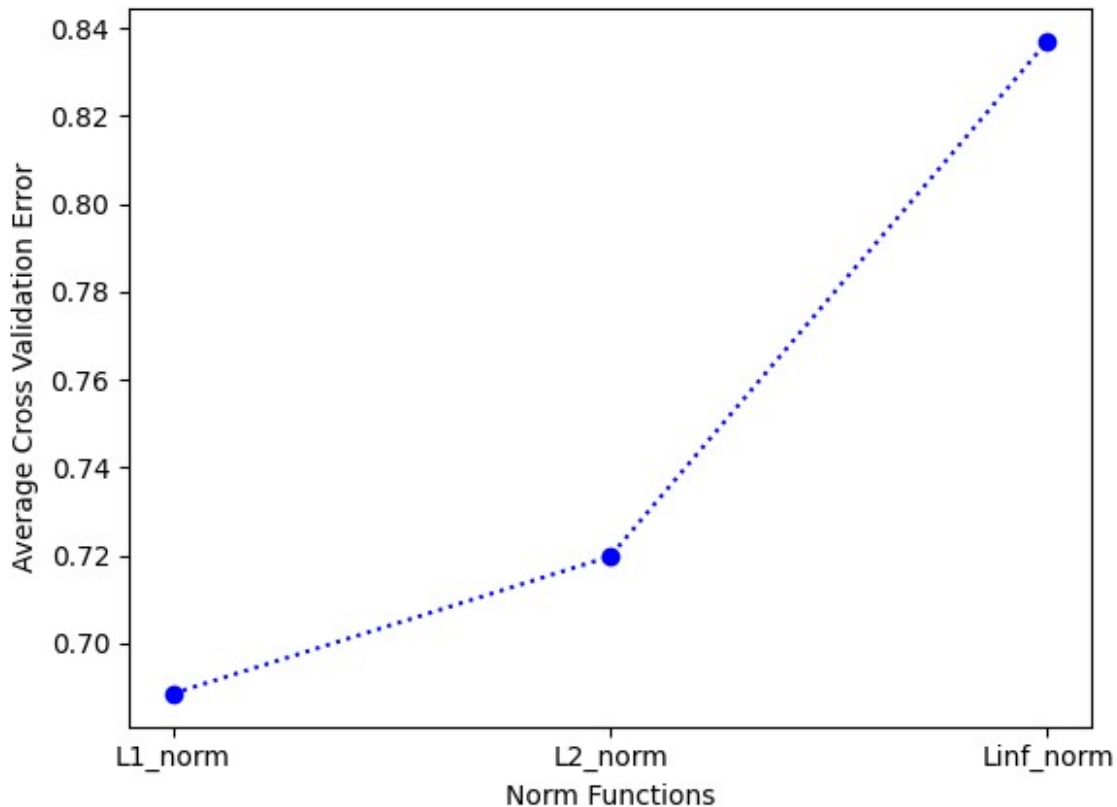
k = 10
avg_err = [CE_err(num_training,num_folds, kk = k, norm= norm) for
norm in norms]
norm_types = ["L1_norm", "L2_norm", "Linf_norm"]
[print("When k is = {}, Average cross validation error: {}".format(
    norm_types[k], avg_err[k])) for k in range(len(norms))]

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

fig, ax = plt.subplots(1, 1)
ax.plot(avg_err, "bo:")
ax.set_xlabel("Norm Functions")
ax.set_xticks(np.arange(3), norm_types)
ax.set_ylabel("Average Cross Validation Error")
Text(0, 0.5, 'Average Cross Validation Error')

```





## Questions:

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

## Answers:

- (1) The l1 norm has the best cross-validation error, out of all the norms
- (2) The cross validation error for the l1 norm with the optimal  $k = 10$  is 0.7344

## 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.
```

```

# ===== #
best_norm = lambda x: np.linalg.norm(x, ord=1)
knn.train(X=X_train, y=y_train)

dists_test = knn.compute_distances(X=X_test, norm=best_norm)
predictions_test = knn.predict_labels(dists=dists_test, k=10)

total_incorrect = num_test -
np.count_nonzero(predictions_test==y_test)
err = total_incorrect/num_test

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

print('Error rate achieved: {}'.format(error))
Error rate achieved: 0.722

```

## Question:

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

## Answer:

There wasn't that great of an improvement, however the new error rate is 0.004% lower than naively choosing the  $k=1$  and l2 norm for the computation.(0.726 vs 0.722)

```
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].
                # ===== #
                distance = X[i,:] - self.X_train[j,:]
                dists[i,j] = norm(distance)
                # ===== #
                # 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.
# ===== #
dists = np.sqrt(np.sum(X**2, axis=1).reshape((num_test,1)) + np.sum(self.X_train**2,
axis=1) - 2*np.dot(X,self.X_train.T))
# ===== #
# END YOUR CODE HERE
# ===== #
#np.sqrt((X- Xtrain))

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.
        # ===== #

        sorted = np.argsort(dists[i,:])
        closest_y = self.y_train[sorted[:k]]
        y_pred[i] = np.argmax(np.bincount(closest_y))

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

    return y_pred

```



# 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`

*#Downloading the CIFAR-10 data*

```
import shutil
import urllib.request
urllib.request.urlretrieve(
    "http://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz",
    "utils/datasets/cifar-10-python.tar.gz")
```

*#Unzipping the downloaded file*

```
shutil.unpack_archive(
    "utils/datasets/cifar-10-python.tar.gz", "utils/datasets/")
```

```
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 = 'utils/datasets/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,)
```

## 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 nnrl import Softmax
import subprocess

# List of pip commands

# 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
```

## Question:

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

## Answer:

Initially, no training is done on the Softmax. Thus, as the weights are initialized to zero, the value of the loss according to the softmax is:  $-1/m * \log((1/10)^m) = \log(10)$ , which is close to 2.3

### 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: -0.659919 analytic: -0.659919, relative error: 4.811308e-08
numerical: -0.887144 analytic: -0.887144, relative error: 1.872054e-08
numerical: 0.593711 analytic: 0.593711, relative error: 1.212518e-08
numerical: 2.179519 analytic: 2.179519, relative error: 3.471660e-10
numerical: -0.557324 analytic: -0.557324, relative error: 1.079517e-07
numerical: 1.998294 analytic: 1.998293, relative error: 1.228015e-08
numerical: -0.522391 analytic: -0.522392, relative error: 1.038401e-07
numerical: -1.286610 analytic: -1.286610, relative error: 2.583998e-09
numerical: 0.518052 analytic: 0.518052, relative error: 9.485669e-09
numerical: -4.808219 analytic: -4.808220, relative error: 1.449431e-08
```

## 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: 2.329781476362285 / 400.2127530552219
computed in 0.07850503921508789s
Vectorized loss / grad: 2.329781476362286 / 400.212753055222 computed
in 0.024672985076904297s
difference in loss / grad: -1.3322676295501878e-15
/3.705124159971763e-13

```

## 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

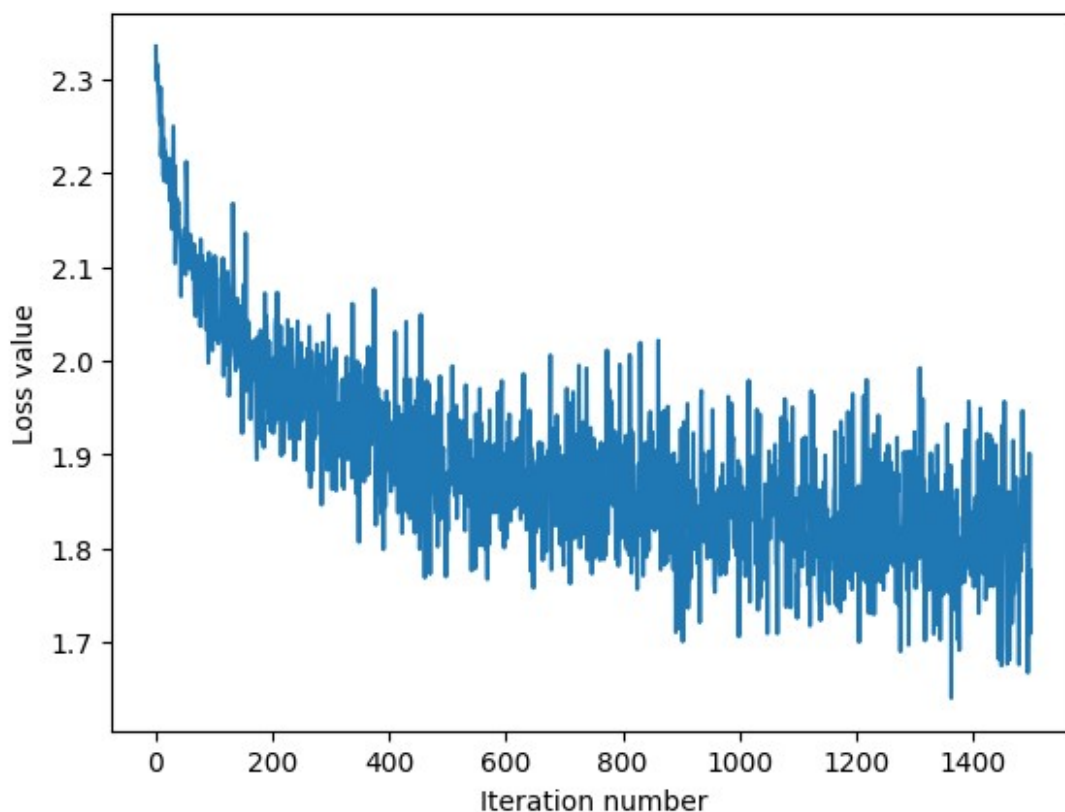
tic = time.time()
loss_hist = softmax.train(X_train, y_train, learning_rate=1e-7,
                           num_iters=1500, verbose=True)
toc = time.time()
print('That took {}s'.format(toc - tic))

plt.plot(loss_hist)
plt.xlabel('Iteration number')
plt.ylabel('Loss value')
plt.show()

iteration 0 / 1500: loss 2.3353835450891545
iteration 100 / 1500: loss 2.0225093946317187
iteration 200 / 1500: loss 1.982172871654982
iteration 300 / 1500: loss 1.9356442081331486

```

```
iteration 400 / 1500: loss 1.882893396815689
iteration 500 / 1500: loss 1.81818696973945
iteration 600 / 1500: loss 1.874513153185746
iteration 700 / 1500: loss 1.8361832500173585
iteration 800 / 1500: loss 1.8584086819212184
iteration 900 / 1500: loss 1.9275087067564147
iteration 1000 / 1500: loss 1.824667969507725
iteration 1100 / 1500: loss 1.7731817984393607
iteration 1200 / 1500: loss 1.8636308568113118
iteration 1300 / 1500: loss 1.924074621260815
iteration 1400 / 1500: loss 1.7846918635831293
That took 9.59458303451538s
```



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: 0.37881632653061226
validation accuracy: 0.39
```

## 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.
# ===== #

learning_rates = np.linspace(2, -8, 12)
learning_rates = 10**learning_rates
print(learning_rates)

num_iters = 1500
num_rates = len(learning_rates)
loss_histories = np.zeros((num_rates, num_iters), dtype=float)
validation_accuracies = np.zeros(num_rates, dtype=float)
for i in range(num_rates):
    loss_histories[i] = softmax.train(X_train, y_train,
                                     learning_rate=learning_rates[i],
                                     num_iters=num_iters,
                                     verbose=True)
    y_pred = softmax.predict(X_val)
    accuracy = np.mean(y_pred == y_val)
    validation_accuracies[i] = accuracy

best_lr = learning_rates[np.argmax(validation_accuracies)]
best_accuracy =
validation_accuracies[np.argmax(validation_accuracies)]

print("Best learning rate:", best_lr)
print("Best validation accuracy:", best_accuracy)
# ===== #
```

```
# END YOUR CODE HERE
```

```
# ===== #
```

```
[1.00000000e+02 1.23284674e+01 1.51991108e+00 1.87381742e-01  
2.31012970e-02 2.84803587e-03 3.51119173e-04 4.32876128e-05  
5.33669923e-06 6.57933225e-07 8.11130831e-08 1.00000000e-08]
```

```
iteration 0 / 1500: loss 2.376850104198011  
iteration 100 / 1500: loss 21930450.818899762  
iteration 200 / 1500: loss 18002661.772628788  
iteration 300 / 1500: loss 30139112.424423013  
iteration 400 / 1500: loss 35643494.85798577  
iteration 500 / 1500: loss 22668817.58273412  
iteration 600 / 1500: loss 41879530.25533271  
iteration 700 / 1500: loss 23980976.744694404  
iteration 800 / 1500: loss 28425963.735318568  
iteration 900 / 1500: loss 41653045.50052511  
iteration 1000 / 1500: loss 42389915.23338503  
iteration 1100 / 1500: loss 20859953.49689545  
iteration 1200 / 1500: loss 22777862.354245503  
iteration 1300 / 1500: loss 21338356.136914648  
iteration 1400 / 1500: loss 38965792.72549838  
iteration 0 / 1500: loss 2.431368148995201  
iteration 100 / 1500: loss 3911024.0104619255  
iteration 200 / 1500: loss 3483064.3656539526  
iteration 300 / 1500: loss 3829290.333834709  
iteration 400 / 1500: loss 4001500.0194998938  
iteration 500 / 1500: loss 2473336.2546379687  
iteration 600 / 1500: loss 3306166.6987504843  
iteration 700 / 1500: loss 3536105.661674255  
iteration 800 / 1500: loss 2981740.0236274083  
iteration 900 / 1500: loss 3172550.679629564  
iteration 1000 / 1500: loss 4485552.570294583  
iteration 1100 / 1500: loss 2403146.447038425  
iteration 1200 / 1500: loss 2428078.177848376  
iteration 1300 / 1500: loss 3443958.519265981  
iteration 1400 / 1500: loss 3076566.152687497  
iteration 0 / 1500: loss 2.363535760198207  
iteration 100 / 1500: loss 372500.94385254284  
iteration 200 / 1500: loss 367745.3050547632  
iteration 300 / 1500: loss 266145.98215907655  
iteration 400 / 1500: loss 386647.6504217968  
iteration 500 / 1500: loss 389080.4494084273  
iteration 600 / 1500: loss 314276.8420210867  
iteration 700 / 1500: loss 266822.1915628476  
iteration 800 / 1500: loss 424908.9585048904  
iteration 900 / 1500: loss 347603.45680869895  
iteration 1000 / 1500: loss 415194.3586194626  
iteration 1100 / 1500: loss 284708.4698438062  
iteration 1200 / 1500: loss 299317.73308567784  
iteration 1300 / 1500: loss 336885.17984189413
```



```
iteration 1400 / 1500: loss 550591.4480033843
iteration 0 / 1500: loss 2.3577965172090742
iteration 100 / 1500: loss 49415.40619365274
iteration 200 / 1500: loss 48094.422779138775
iteration 300 / 1500: loss 62282.17075980758
iteration 400 / 1500: loss 55219.77597413651
iteration 500 / 1500: loss 79037.59562536365
iteration 600 / 1500: loss 46623.959253629255
iteration 700 / 1500: loss 55592.77071718095
iteration 800 / 1500: loss 50293.75325395955
iteration 900 / 1500: loss 45938.17900252937
iteration 1000 / 1500: loss 60126.71120511817
iteration 1100 / 1500: loss 52526.98692114738
iteration 1200 / 1500: loss 71199.01830421321
iteration 1300 / 1500: loss 47138.68953778818
iteration 1400 / 1500: loss 83020.0362768881
iteration 0 / 1500: loss 2.3623163529832842
iteration 100 / 1500: loss 7386.785167363726
iteration 200 / 1500: loss 6561.724714169041
iteration 300 / 1500: loss 7762.457975035249
iteration 400 / 1500: loss 8185.593338401697
iteration 500 / 1500: loss 9241.982105076282
iteration 600 / 1500: loss 6448.578867873589
iteration 700 / 1500: loss 5490.309019938016
iteration 800 / 1500: loss 5302.935330988444
iteration 900 / 1500: loss 6600.164089494958
iteration 1000 / 1500: loss 4467.617866491232
iteration 1100 / 1500: loss 7304.733435829171
iteration 1200 / 1500: loss 9017.892633791746
iteration 1300 / 1500: loss 6068.499578394603
iteration 1400 / 1500: loss 5582.4341722882955
iteration 0 / 1500: loss 2.3303211245338518
iteration 100 / 1500: loss 702.8104745719538
iteration 200 / 1500: loss 1003.0492399680344
iteration 300 / 1500: loss 700.1187822421518
iteration 400 / 1500: loss 835.1644800620383
iteration 500 / 1500: loss 446.30354572286905
iteration 600 / 1500: loss 982.7065184736753
iteration 700 / 1500: loss 879.9026852273258
iteration 800 / 1500: loss 822.1161610487652
iteration 900 / 1500: loss 579.046723437298
iteration 1000 / 1500: loss 536.6872540150157
iteration 1100 / 1500: loss 703.7149140578282
iteration 1200 / 1500: loss 756.385763502355
iteration 1300 / 1500: loss 648.743959299979
iteration 1400 / 1500: loss 568.723610575774
iteration 0 / 1500: loss 2.3819157186171527
iteration 100 / 1500: loss 101.90366538590519
iteration 200 / 1500: loss 118.06337378241717
```

```
iteration 300 / 1500: loss 77.94405186044338
iteration 400 / 1500: loss 103.57086415363067
iteration 500 / 1500: loss 88.07708962144505
iteration 600 / 1500: loss 87.00524664953944
iteration 700 / 1500: loss 99.1464793458803
iteration 800 / 1500: loss 80.50596718681074
iteration 900 / 1500: loss 54.873063256993994
iteration 1000 / 1500: loss 138.13711973359645
iteration 1100 / 1500: loss 120.9919651195003
iteration 1200 / 1500: loss 92.57520058567
iteration 1300 / 1500: loss 133.3349194177065
iteration 1400 / 1500: loss 81.83726018936659
iteration 0 / 1500: loss 2.3484318532800312
iteration 100 / 1500: loss 12.182852376071834
iteration 200 / 1500: loss 12.087503451124547
iteration 300 / 1500: loss 10.615130409607861
iteration 400 / 1500: loss 13.84461942104537
iteration 500 / 1500: loss 12.571806312922758
iteration 600 / 1500: loss 6.984199121128224
iteration 700 / 1500: loss 8.64487613420592
iteration 800 / 1500: loss 14.502744343807787
iteration 900 / 1500: loss 10.463870556146347
iteration 1000 / 1500: loss 13.839418943398323
iteration 1100 / 1500: loss 8.439131732079757
iteration 1200 / 1500: loss 25.471756400015657
iteration 1300 / 1500: loss 11.112198046857
iteration 1400 / 1500: loss 12.513600477833668
iteration 0 / 1500: loss 2.3884800630134353
iteration 100 / 1500: loss 1.9244324754685642
iteration 200 / 1500: loss 1.8307907909717744
iteration 300 / 1500: loss 1.9401241351781335
iteration 400 / 1500: loss 1.9390054291197396
iteration 500 / 1500: loss 1.7642694060698612
iteration 600 / 1500: loss 1.914377770811749
iteration 700 / 1500: loss 2.140990947900194
iteration 800 / 1500: loss 1.823280960376048
iteration 900 / 1500: loss 1.7823081429383212
iteration 1000 / 1500: loss 1.8361834433051791
iteration 1100 / 1500: loss 1.7918254575079688
iteration 1200 / 1500: loss 1.6720707159338417
iteration 1300 / 1500: loss 1.67405276230121
iteration 1400 / 1500: loss 1.8196460558278313
iteration 0 / 1500: loss 2.3738551189467474
iteration 100 / 1500: loss 1.8800757751651815
iteration 200 / 1500: loss 1.9032012169003765
iteration 300 / 1500: loss 1.7552064783210315
iteration 400 / 1500: loss 1.8057427708923495
iteration 500 / 1500: loss 1.703297149344325
iteration 600 / 1500: loss 1.8106431097810396
```

```
iteration 700 / 1500: loss 1.781909808759197
iteration 800 / 1500: loss 1.8347189009757465
iteration 900 / 1500: loss 1.772411828877675
iteration 1000 / 1500: loss 1.7426963229988734
iteration 1100 / 1500: loss 1.7939978026703687
iteration 1200 / 1500: loss 1.7390890428035715
iteration 1300 / 1500: loss 1.6237130713164825
iteration 1400 / 1500: loss 1.6319798366216436
iteration 0 / 1500: loss 2.381247622289375
iteration 100 / 1500: loss 2.0597882886583454
iteration 200 / 1500: loss 2.0907885620978175
iteration 300 / 1500: loss 1.908021469908436
iteration 400 / 1500: loss 1.891647808828347
iteration 500 / 1500: loss 1.8553358364243135
iteration 600 / 1500: loss 2.018174773065773
iteration 700 / 1500: loss 1.9029340789229914
iteration 800 / 1500: loss 1.8551084631237507
iteration 900 / 1500: loss 1.9528490745867333
iteration 1000 / 1500: loss 1.8609197826141204
iteration 1100 / 1500: loss 1.8487548977230108
iteration 1200 / 1500: loss 1.8788884653758466
iteration 1300 / 1500: loss 1.876206304543427
iteration 1400 / 1500: loss 1.8559469827187676
iteration 0 / 1500: loss 2.472523652438493
iteration 100 / 1500: loss 2.310952064158016
iteration 200 / 1500: loss 2.261801262801085
iteration 300 / 1500: loss 2.172240024995832
iteration 400 / 1500: loss 2.2023515733714847
iteration 500 / 1500: loss 2.170122715607117
iteration 600 / 1500: loss 2.147718560680545
iteration 700 / 1500: loss 2.145091568186487
iteration 800 / 1500: loss 2.1314891893644052
iteration 900 / 1500: loss 2.130507158028264
iteration 1000 / 1500: loss 2.0524036308637954
iteration 1100 / 1500: loss 2.054309380022291
iteration 1200 / 1500: loss 2.1269059271129307
iteration 1300 / 1500: loss 2.0956204583447713
iteration 1400 / 1500: loss 2.0905403988265276
Best learning rate: 6.579332246575682e-07
Best validation accuracy: 0.395
```

According to the output above, the best validation accuracy is when the learning rate is around  $6.579332246575682 \times 10^{-7}$  out of 12 different learning rates. The best validation accuracy that was achieved for this learning rate = 0.395

```
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.)
    # ===== #
```

```
N, D = X.shape
scores = np.dot(X, self.W.T)
for i in range(N):
    score_ith = scores[i,:]
    score_ith -= np.max(score_ith)
    current_score = score_ith[y[i]]
    loss += current_score - np.log(np.sum(np.exp(score_ith)))
```

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

```
    return -loss/N
```

```
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.
# ===== #
N = X.shape[0]
C = self.W.shape[0]
scores = np.dot(X, self.W.T)
for i in range(N):
    score_ith = scores[i,:]
    score_ith -= np.max(score_ith)
    grad[y[i]] += X[i]
    for c in range(C):
        grad[c] -= np.exp(score_ith[c])/np.sum(np.exp(score_ith))*X[i]

loss = self.loss(X, y)
# ===== #
# END YOUR CODE HERE
# ===== #

return loss, -grad/N

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 ouputs 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.
    # ===== #
    X = np.array(X)
    N = X.shape[0]
    C = self.W.shape[0]
    scores = np.dot(X, self.W.T)
    grad = np.sum(np.exp)
    scores = np.dot(X, self.W.T)
    scores = (scores.T - np.max(scores, axis=1)).T

```

```
loss = -1/N * np.sum((scores[range(N), y] - np.log(np.sum(np.exp(scores), axis=1))),
axis=0)
```

```
probs = -np.exp(scores)/np.sum(np.exp(scores), axis=1, keepdims=True)
probs[range(N), y] += 1
grad = -1/N * np.dot(probs.T, X)
```

```
# ===== #
# 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  $0 \leq c < C$  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.
```

```
# ===== #
```

```
indexes = np.random.choice(num_train, batch_size)
```

```
X_batch = X[indexes]
```

```
y_batch = y[indexes]
```

```
# ===== #
```

```
# 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 -= learning_rate* grad

# ===== #
# 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.
    # ===== #
    # predicted = X.dot(self.W.T)
    # #normalize
    # predicted = predicted - np.argmax(predicted)

    # y_predict = np.exp(predicted) / np.sum(np.exp(predicted))
    # y_pred = np.argmax(y_predict)

    scores = np.dot(X, self.W.T) # scores.shape = num_examples * num_classes
    scores = (scores.T - np.max(scores, axis=1)).T # to control the overflow

    probs = np.exp(scores) / np.sum(np.exp(scores), axis=1, keepdims=True)
    y_pred = np.argmax(probs, axis=1)

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

    return y_pred

```