

# This is the svm workbook for ECE C147/C247 Assignment #2 ¶

Please follow the notebook linearly to implement a linear support vector machine.

Please print out the workbook entirely when completed.

We thank Serena Yeung & Justin Johnson for permission to use code written for the CS 231n class (cs231n.stanford.edu). These are the functions in the cs231n folders and includes code to preprocess and show the images. The classifiers used are based off of code prepared for CS 231n as well.

The goal of this workbook is to give you experience with training an SVM classifier via gradient descent.

## Importing libraries and data setup

```
In [2]: import numpy as np # for doing most of our calculations
import matplotlib.pyplot as plt # for plotting
from cs231n.data_utils import load_CIFAR10 # function to load the CIFAR-10 data set.
import pdb

# 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
```

```
In [3]: # 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)

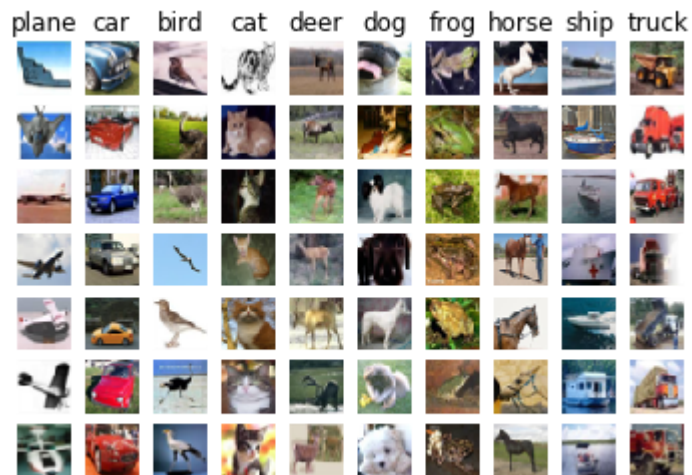
# As a sanity check, we print out the size of the training and test data.
print('Training data shape: ', X_train.shape)
print('Training labels shape: ', y_train.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)
```

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

```

In [4]: # 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()

```



```
In [5]: # Split the data into train, val, and test sets. In addition we will
# create a small development set as a subset of the training data;
# we can use this for development so our code runs faster.
num_training = 49000
num_validation = 1000
num_test = 1000
num_dev = 500

# Our validation set will be num_validation points from the original
# training set.
mask = range(num_training, num_training + num_validation)
X_val = X_train[mask]
y_val = y_train[mask]

# Our training set will be the first num_train points from the original
# training set.
mask = range(num_training)
X_train = X_train[mask]
y_train = y_train[mask]

# We will also make a development set, which is a small subset of
# the training set.
mask = np.random.choice(num_training, num_dev, replace=False)
X_dev = X_train[mask]
y_dev = y_train[mask]

# We use the first num_test points of the original test set as our
# test set.
mask = range(num_test)
X_test = X_test[mask]
y_test = y_test[mask]

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, 32, 32, 3)
Train labels shape: (49000,)
Validation data shape: (1000, 32, 32, 3)
Validation labels shape: (1000,)
Test data shape: (1000, 32, 32, 3)
Test labels shape: (1000,)
Dev data shape: (500, 32, 32, 3)
Dev labels shape: (500,)
```

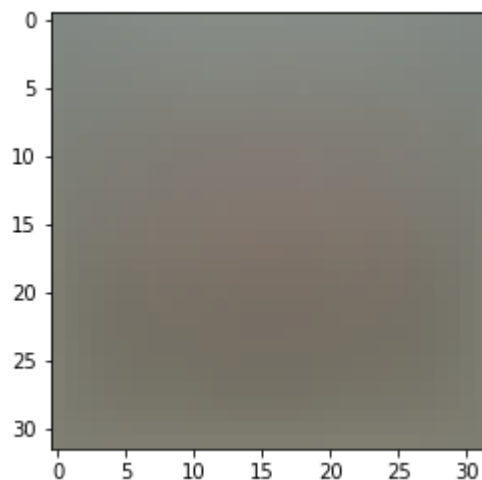
```
In [6]: # Preprocessing: reshape the image data into rows
X_train = np.reshape(X_train, (X_train.shape[0], -1))
X_val = np.reshape(X_val, (X_val.shape[0], -1))
X_test = np.reshape(X_test, (X_test.shape[0], -1))
X_dev = np.reshape(X_dev, (X_dev.shape[0], -1))

# As a sanity check, print out the shapes of the data
print('Training data shape: ', X_train.shape)
print('Validation data shape: ', X_val.shape)
print('Test data shape: ', X_test.shape)
print('dev data shape: ', X_dev.shape)
```

```
Training data shape: (49000, 3072)
Validation data shape: (1000, 3072)
Test data shape: (1000, 3072)
dev data shape: (500, 3072)
```

```
In [7]: # Preprocessing: subtract the mean image
# first: compute the image mean based on the training data
mean_image = np.mean(X_train, axis=0)
print(mean_image[:10]) # print a few of the elements
plt.figure(figsize=(4,4))
plt.imshow(mean_image.reshape((32,32,3)).astype('uint8')) # visualize the mean image
plt.show()
```

```
[130.64189796 135.98173469 132.47391837 130.05569388 135.34804082
 131.75402041 130.96055102 136.14328571 132.47636735 131.48467347]
```



```
In [8]: # second: subtract the mean image from train and test data
X_train -= mean_image
X_val -= mean_image
X_test -= mean_image
X_dev -= mean_image
```

```
In [9]: # third: append the bias dimension of ones (i.e. bias trick) so that our SVM
# only has to worry about optimizing a single weight matrix W.
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))])

print(X_train.shape, X_val.shape, X_test.shape, X_dev.shape)

(49000, 3073) (1000, 3073) (1000, 3073) (500, 3073)
```

## Question:

(1) For the SVM, we perform mean-subtraction on the data. However, for the KNN notebook, we did not. Why?

## Answer:

(1) KNN's algorithm basically works by the closest k-nearest points, and declaring the label of a point based off of the number of the nearest k-points. Evidently, this only cares about relative distance from one point to k-others. Mean subtraction on SVM is done in order to focus on the information-carrying pixels of the image.

## Training an SVM

The following cells will take you through building an SVM. 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.

```
In [10]: from nndl.svm import SVM
```

```
In [11]: # Declare an instance of the SVM class.
# Weights are initialized to a random value.
# Note, to keep people's initial 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]

svm = SVM(dims=[num_classes, num_features])
```

## SVM loss

```
In [12]: ## Implement the loss function for in the SVM class(nndl/svm.py), svm.loss()

loss = svm.loss(X_train, y_train)
print('The training set loss is {}'.format(loss))

# If you implemented the loss correctly, it should be 15569.98
```

The training set loss is 15569.977915410243.

## SVM gradient

```
In [13]: ## Calculate the gradient of the SVM class.
# For convenience, we'll write one function that computes the loss
# and gradient together. Please modify svm.loss_and_grad(X, y).
# You may copy and paste your loss code from svm.loss() here, and then
# use the appropriate intermediate values to calculate the gradient.

loss, grad = svm.loss_and_grad(X_dev, y_dev)

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

numerical: -5.683140 analytic: -5.683140, relative error: 1.294764e-09
numerical: 2.267528 analytic: 2.267529, relative error: 1.264653e-07
numerical: -9.292961 analytic: -9.292962, relative error: 4.265034e-08
numerical: 13.249561 analytic: 13.249561, relative error: 3.131214e-09
numerical: 4.672054 analytic: 4.672055, relative error: 4.055359e-08
numerical: 7.825805 analytic: 7.825805, relative error: 2.219825e-08
numerical: 7.475661 analytic: 7.475661, relative error: 2.117599e-08
numerical: -3.105096 analytic: -3.105096, relative error: 8.219805e-08
numerical: -5.497770 analytic: -5.497770, relative error: 4.429351e-08
numerical: -20.060421 analytic: -20.060421, relative error: 1.443535e-09
```

## A vectorized version of SVM

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

```
In [14]: import time
```

```
In [15]: ## Implement svm.fast_loss_and_grad which calculates the loss and gradient
# WITHOUT using any for loops.

# Standard Loss and gradient
tic = time.time()
loss, grad = svm.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 = svm.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 fa
ster.
print('difference in loss / grad: {} / {}'.format(loss - loss_vectorized, np.l
    inalg.norm(grad - grad_vectorized)))

# You should notice a speedup with the same output, i.e., differences on the o
rder of 1e-12
```

```
Normal loss / grad_norm: 15510.045561497369 / 2295.4427712149827 computed in
0.13462305068969727s
Vectorized loss / grad: 15510.045561497373 / 2295.442771214983 computed in 0.
056557655334472656s
difference in loss / grad: -3.637978807091713e-12 / 2.746320501330392e-12
```

## Stochastic gradient descent

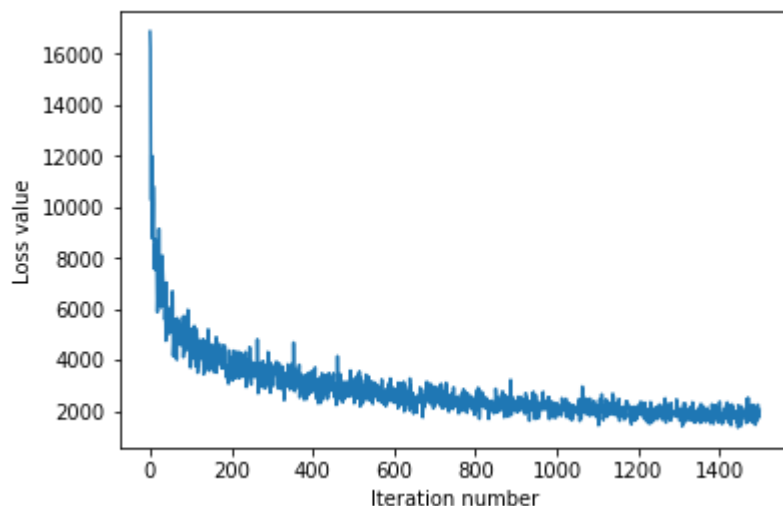
We now implement stochastic gradient descent. This uses the same principles of gradient descent we discussed in class, however, it calculates the gradient by only using examples from a subset of the training set (so each gradient calculation is faster).

In [16]: *# Implement svm.train() by filling in the code to extract a batch of data  
# and perform the gradient step.*

```
tic = time.time()
loss_hist = svm.train(X_train, y_train, learning_rate=5e-4,
                      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 16878.859625643894
iteration 100 / 1500: loss 3698.3744294458324
iteration 200 / 1500: loss 3749.3722795434505
iteration 300 / 1500: loss 3232.201604880448
iteration 400 / 1500: loss 2786.9285054561765
iteration 500 / 1500: loss 2911.8902158067995
iteration 600 / 1500: loss 2696.123480933879
iteration 700 / 1500: loss 2959.5756376002882
iteration 800 / 1500: loss 2512.8602634753493
iteration 900 / 1500: loss 2105.9669921625746
iteration 1000 / 1500: loss 2313.428837770445
iteration 1100 / 1500: loss 1732.4754464411596
iteration 1200 / 1500: loss 2114.4851353256367
iteration 1300 / 1500: loss 2050.994800231626
iteration 1400 / 1500: loss 1814.3598110234714
That took 18.10395121574402s
```



**Evaluate the performance of the trained SVM on the validation data.**



```
In [17]: ## Implement svm.predict() and use it to compute the training and testing error.

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

```
training accuracy: 0.29612244897959183
validation accuracy: 0.303
```

## Optimize the SVM

Note, to make things faster and simpler, we won't do k-fold cross-validation, but will only optimize the hyperparameters on the validation dataset (X\_val, y\_val).

```

In [20]: # ===== #
# YOUR CODE HERE:
#   Train the SVM 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.
#   Note: You do not need to modify SVM class for this section
# ===== #
my_time = time.time()
learning_rates = np.linspace(0, 0.01, 100)
y_val_accs = []
best_learning_rate = -1
best_val_acc = -1

for rate in learning_rates:
    loss_hist = svm.train(X_train, y_train, learning_rate=rate, num_iters=1500
, verbose=False)
    y_train_pred = svm.predict(X_train)
    train_acc = np.mean(np.equal(y_train, y_train_pred))
    y_val_pred = svm.predict(X_val)
    y_val_acc = np.mean(np.equal(y_val, y_val_pred))
    y_val_accs.append(y_val_acc)
    if y_val_acc > best_val_acc:
        best_val_acc = y_val_acc
        best_learning_rate = rate
    #print(rate)
    print('.')

print("Best learning rate: ", best_learning_rate, " Best Accuracy: ", best_val
_acc, "Err: ", 1 - best_val_acc)

loss_hist = svm.train(X_train, y_train, learning_rate=best_learning_rate, num_
iters=1500, verbose=False)
y_test_pred = svm.predict(X_test)
test_acc = np.mean(np.equal(y_test, y_test_pred))
print("Test Acc: ", test_acc, " Test Error: ", 1-test_acc)

plt.plot(learning_rates, y_val_accs)
plt.xlabel('Learning Rate')
plt.ylabel('Accuracy')

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

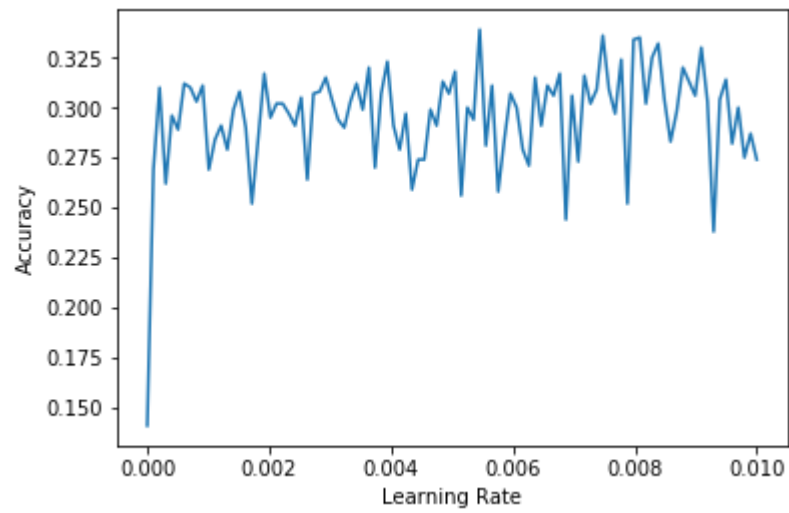
print("Time took {}s".format(time.time() - my_time))

```

.....

•

```
Best learning rate: 0.005454545454545455 Best Accuracy: 0.339 Err: 0.661
Test Acc: 0.315 Test Error: 0.685
Time took 1877.8136940002441s
```



In [ ]:

```

import numpy as np
import pdb

"""
This code was based off of code from cs231n at Stanford University, and modified for ECE C147/C247 at UCLA.
"""

class SVM(object):

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

    def init_weights(self, dims):
        """
        Initializes the weight matrix of the SVM. 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)

    def loss(self, X, y):
        """
        Calculates the SVM 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
        """

        # compute the loss and the gradient
        num_classes = self.W.shape[0]
        num_train = X.shape[0]
        loss = 0.0

        for i in np.arange(num_train):
            L = 0
            for j in range(num_classes):
                zj = 0
                if y[i] != j:
                    zj = 1 + np.dot(self.W[j].T, X[i]) - np.dot(self.W[y[i]].T, X[i])
                    L = L + max(0, zj)
            loss = loss + L
        loss = loss / num_train #normalize

        # ===== #
        # 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.

```

```

"""

# compute the loss and the gradient
num_classes = self.W.shape[0]
num_train = X.shape[0]
loss = 0.0
grad = np.zeros_like(self.W)

for i in np.arange(num_train):
    L = 0
    for j in range(num_classes):
        zj = 0
        if y[i] != j:
            zj = 1 + np.dot(self.W[j].T, X[i]) - np.dot(self.W[y[i]].T, X[i])
            L = L + max(0, zj)
        #Hinge
        grad[j,:] += X[i].T * (zj > 0)
        grad[y[i],:] -= X[i].T * (zj > 0)
    loss = loss + L

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

loss /= num_train
grad /= num_train

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 ouptuts as loss_and_grad.
    """
    loss = 0.0
    grad = np.zeros(self.W.shape) # initialize the gradient as zero

    # ===== #
    # YOUR CODE HERE:
    #   Calculate the SVM Loss WITHOUT any for loops.
    # ===== #
    num_classes = self.W.shape[1]

```

```

num_train = X.shape[0]

aj = np.dot(X, self.W.T)
ayi = np.resize(aj[np.arange(num_train), y], (num_train, 1))
Losses = np.maximum(0, 1 + aj - ayi)
loss = (np.sum(Losses) - num_train)/num_train

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

# ===== #
# YOUR CODE HERE:
#   Calculate the SVM grad WITHOUT any for loops.
# ===== #
m = np.maximum(0, X.dot(self.W.T) - X.dot(self.W.T)[np.arange(num_train), y].reshape(-1, 1) + 1)
m[np.arange(num_train), y] = 0
m[m > 0] = 1
m[np.arange(num_train), y] = -np.sum(m, axis=1)
grad = ((X.T.dot(m)).T) / num_train
# ===== #
# 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: (dim, batch_size)

```



```

# - 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.
# ===== #
a = list(range(len(X)))
idxs = np.random.choice(a, size = batch_size, replace=False)
X_batch = np.vstack([X[i] for i in idxs])
y_batch = [y[i] for i in idxs]
# ===== #
# 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 - 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 with the parameter self.W.
    # ===== #
    y_pred = np.argmax(np.dot(X, self.W.T), axis=1)
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    return y_pred

```

# 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.

We thank Serena Yeung & Justin Johnson for permission to use code written for the CS 231n class (cs231n.stanford.edu). These are the functions in the cs231n folders and code in the jupyter notebook to preprocess and show the images. The classifiers used are based off of code prepared for CS 231n as well.

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

```
In [1]: import numpy as np # for doing most of our calculations
import matplotlib.pyplot as plt # for plotting
from cs231n.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
```

```
In [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,)
```

```
In [3]: # 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()
```



```
In [4]: # 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.

```
In [5]: # Import the KNN class

from nndl import KNN
print('Loaded Kernel') #Debug to see if my file Loaded
```

Loaded Kernel

```
In [6]: # 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) The function `knn.train` assigns `X_train` and `y_train` to the member variables of the KNN class representing the training data.
- (2) Pros: Simple and fast implementation. Cons: This may be more memory intensive since we have to store the training data

## 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.

```
In [7]: # 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: 67.74084901809692

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.

```
In [8]: # 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: 0.6519057750701904

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.

```
In [9]: # 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  
  
y_pred = knn.predict_labels(dists_L2_vectorized, k=1)  
count_correct = np.sum(y_pred == y_test)  
count_incorrect = num_test - count_correct  
err = count_incorrect / num_test  
  
error = err  
  
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.

```

In [10]: # 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]
# ===== #
X_train_folds = np.split(X_train, num_folds)
y_train_folds = np.split(y_train, num_folds)
print(len(y_train_folds[0]))
# ===== #
# END YOUR CODE HERE
# ===== #

```

1000

## 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.

```

In [47]: 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.
# ===== #

cross_val_err = np.empty(len(ks))
m_num_test = len(y_train_folds[0])

for i in range(0, len(ks)):
    ks_err = np.empty(5)
    for j in range(0, num_folds):
        #Make training folds
        train_idx = list(range(0, num_folds))
        train_idx.remove(j)

        m_x_train = X_train_folds[train_idx[0]]
        m_y_train = y_train_folds[train_idx[0]]

        for l in range(1, len(train_idx)):
            m_x_train = np.concatenate((m_x_train, X_train_folds[train_idx[l]]), axis=0)
            m_y_train = np.concatenate((m_y_train, y_train_folds[train_idx[l]]), axis=0)

        #Testing folds
        m_x_test = X_train_folds[j]
        m_y_test = y_train_folds[j]

        knn2 = KNN()
        knn2.train(X = m_x_train, y = m_y_train)
        #Do test

        dists_l2 = knn2.compute_L2_distances_vectorized(X=m_x_test)

        y_pred = knn2.predict_labels(dists_l2, ks[i])
        count_correct = np.sum(y_pred == m_y_test)
        count_incorrect = m_num_test - count_correct
        ks_err[j] = count_incorrect / m_num_test
        #print("len_crossvalerr:", len(cross_val_err))
        #print(cross_val_err)
        cross_val_err[i] = np.average(ks_err)

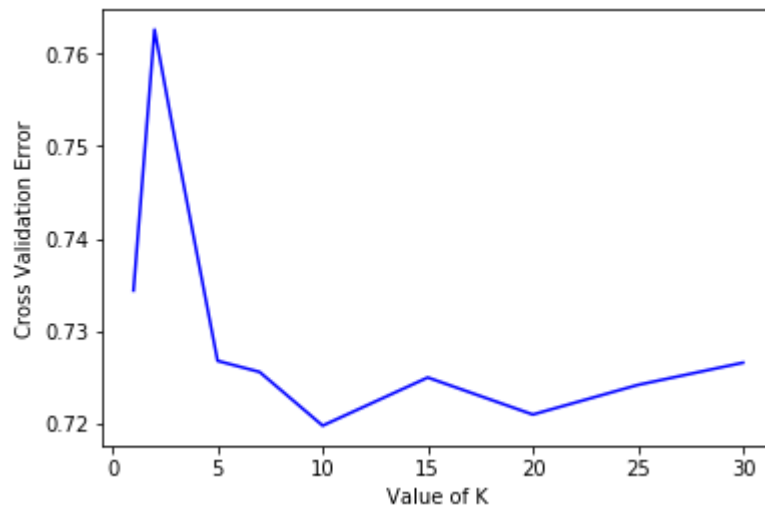
print("KS: " + str(ks))
print("ERROR: ", cross_val_err.round(decimals=4))
plt.plot(ks, cross_val_err, 'b')
plt.xlabel("Value of K")
plt.ylabel("Cross Validation Error")
plt.show()
# ===== #

```



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

```
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
len_crossvalerr: 10  
KS: [1, 2, 3, 5, 7, 10, 15, 20, 25, 30]  
ERROR: [0.7344 0.7626 0.7504 0.7268 0.7256 0.7198 0.725 0.721 0.7242 0.726  
6]
```



Computation time: 49.74

## 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) The lowest value for cross-validation occurred at  $k = 10$ . Lower error means better accuracy.
- (2) At which the CV error is 0.7198

## 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.

```

In [13]: 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]

norm_name_list = ["L1", "L2", "Linf"]

k = 10
cross_val_error = np.empty(len(norms))
m_num_test = len(y_train_folds[0])

for i in range(0, len(norms)):
    norm_err = np.empty(5)
    for j in range(0, num_folds):
        train_idx = list(range(0, num_folds))
        train_idx.remove(j)

        m_x_train = X_train_folds[train_idx[0]]
        m_y_train = y_train_folds[train_idx[0]]

        for l in range(1, len(train_idx)):
            m_x_train = np.concatenate((m_x_train, X_train_folds[train_idx[l]]), axis=0)
            m_y_train = np.concatenate((m_y_train, y_train_folds[train_idx[l]]), axis=0)

        m_x_test = X_train_folds[j]
        m_y_test = y_train_folds[j]

        knn3 = KNN()
        knn3.train(X = m_x_train, y = m_y_train)
        #Do test
        if i!= 1:
            dists = knn3.compute_distances(X=m_x_test, norm = norms[i])
        else:
            dists = knn3.compute_L2_distances_vectorized(X=m_x_test)

        y_pred = knn3.predict_labels(dists, k)
        count_correct = np.sum(y_pred == m_y_test)
        count_incorrect = m_num_test - count_correct
        norm_err[j] = count_incorrect / m_num_test
    cross_val_error[i] = np.average(norm_err)

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

```

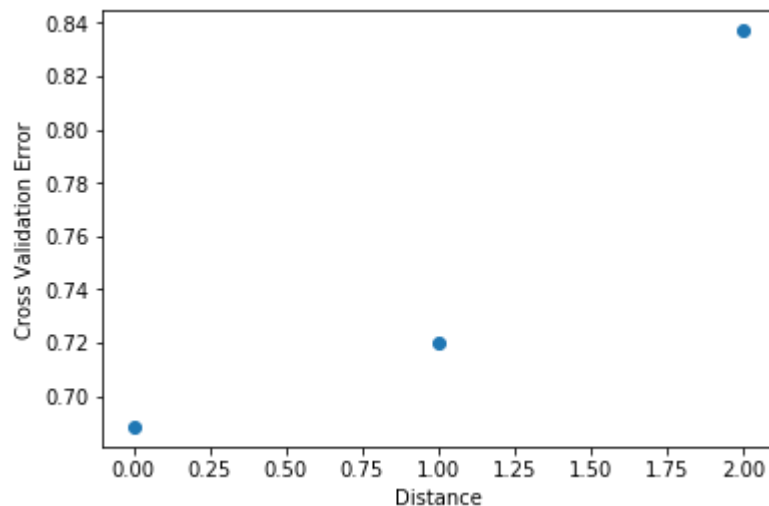
```

print("L1, L2, LInf Distance measures -> errors")
print(cross_val_error)
plt.plot(cross_val_error, 'o')
plt.xlabel("Distance")
plt.ylabel("Cross Validation Error")

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

```

L1, L2, LInf Distance measures -> errors  
 [0.6886 0.7198 0.837 ]  
 Computation time: 960.58



## 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) Note that in the plot above, L1 is at 0.00, L2 at 1.00, LInf at 2.00. L1 norm had the best CV Error of 0.688
- (2) It had a CV error of 0.6886 when k was equal to 10.

## 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.

```

In [55]: error = 1

test_num = len(y_test)
norm_l1 = lambda x: np.linalg.norm(x, ord=1)
k = 10

knn4 = KNN()
knn4.train(X=X_train, y=y_train)
dists = knn4.compute_distances(X=X_test, norm=norm_l1)

y_pred = knn4.predict_labels(dists,k)
count_correct = np.sum(y_pred == y_test)
count_incorrect = test_num - count_correct
error = count_incorrect / test_num

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

# ===== #
# 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:

The error did improve from 0.726 to 0.722 when  $k$  was equal to 10 when using the L1 norm, signifying a 0.004 improvement. The L2 norm resulted in an error of 0.7198 which was also about 0.004 better than using the L1 norm. This was slightly unexpected. However, in general L2 did not perform better than L1, whereas L1 performed significantly better than L2 norm when a 10-point knn was used during k-fold cross validation. This can be seen in the graph above

In [ ]:

```

import numpy as np
import pdb

"""
This code was based off of code from cs231n at Stanford University, and modified for ECE C147/C247 at UCLA.
"""

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):
                # This is my code
                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))

        # sqrt(x^2 + y^2 - 2xy)
        dists = -2*(np.dot(X, (self.X_train).T)) + np.square(X).sum(axis=1).reshape(num_test, 1) + np.square(self.X_train).sum(axis = 1)
        dists = np.sqrt(dists)

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

    pass

# ===== #
# 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.
        # ===== #
        sort_indexes = np.argsort(dists[i,:])
        k_near_indices = sort_indexes[0:k]
        k_near_classes = self.y_train[k_near_indices]
        bins = np.bincount(k_near_classes)
        max_indices = np.argmax(bins)
        y_pred[i] = np.amin(max_indices)
        pass

        # ===== #
        # 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.

We thank Serena Yeung & Justin Johnson for permission to use code written for the CS 231n class (cs231n.stanford.edu). These are the functions in the cs231n folders and code in the jupyter notebook to preprocess and show the images. The classifiers used are based off of code prepared for CS 231n as well.

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

```
In [1]: import random
import numpy as np
from cs231n.data_utils import load_CIFAR10
import matplotlib.pyplot as plt

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

```

In [9]: 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 = '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.

```
In [10]: from nndl import Softmax
```

```
In [11]: # 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

```
In [12]: ## Implement the loss function of the softmax using a for loop over
# the number of examples

loss = softmax.loss(X_train, y_train)
```

```
In [13]: 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:

The SVM loss is relatively large, but in comparison, this softmax loss is a lot smaller. This makes sense since the loss function can be comprised of an exponential divided by more exponentials (divided by a sum of exponentials), usually creating a smaller loss than SVM. When this loss is logged, the loss is even then smaller.

Additionally, the  $\ln(10) = 2.3$ , at which we have 10 classes. So if we "randomly" picked a class, we would have a log of 2.3

### Softmax gradient

```
In [15]: ## 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: 3.460603 analytic: 3.460602, relative error: 2.697676e-08  
numerical: 2.390041 analytic: 2.390041, relative error: 4.646907e-09  
numerical: 2.354785 analytic: 2.354785, relative error: 1.076003e-08  
numerical: 2.031387 analytic: 2.031387, relative error: 2.763128e-08  
numerical: -0.066246 analytic: -0.066246, relative error: 1.003596e-06  
numerical: 0.623177 analytic: 0.623176, relative error: 6.279616e-08  
numerical: 0.137967 analytic: 0.137967, relative error: 6.507170e-08  
numerical: 0.477074 analytic: 0.477073, relative error: 1.080586e-07  
numerical: 1.669293 analytic: 1.669293, relative error: 2.096844e-08  
numerical: 2.807520 analytic: 2.807520, relative error: 1.475514e-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.

```
In [16]: import time
```

```
In [17]: ## 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 fa
ster.
print('difference in loss / grad: {} / {} '.format(loss - loss_vectorized, np.l
    inalnorm(grad - grad_vectorized)))

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

```
Normal loss / grad_norm: 2.306647487014123 / 333.1309299317958 computed in 0.
11967897415161133s
Vectorized loss / grad: 2.3066474870141245 / 333.1309299317958 computed in 0.
03989458084106445s
difference in loss / grad: -1.3322676295501878e-15 / 2.1708887438091485e-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).

### Question:

How should the softmax gradient descent training step differ from the svm training step, if at all?

### Answer:

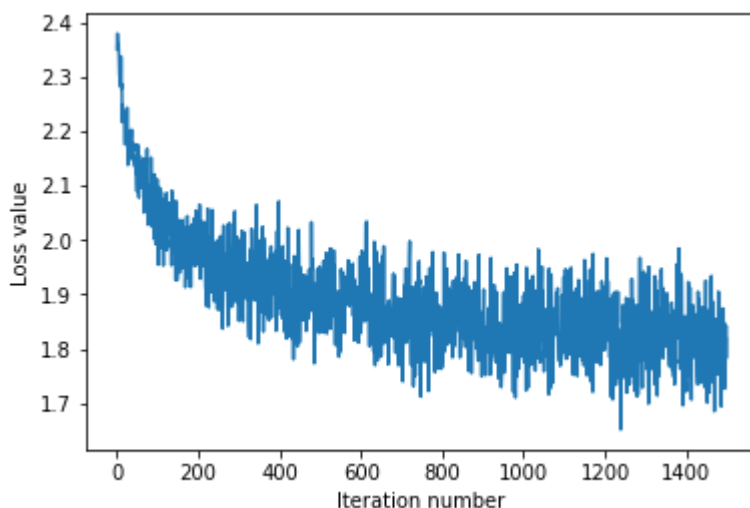
The softmax GD (gradient descent) should have a smaller learning rate than in SVM as vectors are exponentiated in softmax (whereas they aren't in SVM), so having a smaller learning rate allows us to be more accurate and also possibly overflow due to exponentiation.

```
In [22]: # 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 2.3506768483824065
iteration 100 / 1500: loss 2.039956147048289
iteration 200 / 1500: loss 1.9906507989783668
iteration 300 / 1500: loss 1.850921781073332
iteration 400 / 1500: loss 1.8903552634188256
iteration 500 / 1500: loss 1.9115070935052418
iteration 600 / 1500: loss 1.8287083973003708
iteration 700 / 1500: loss 1.8915737345409425
iteration 800 / 1500: loss 1.7766508935283596
iteration 900 / 1500: loss 1.8428792352925394
iteration 1000 / 1500: loss 1.8977364578449698
iteration 1100 / 1500: loss 1.8573034819689855
iteration 1200 / 1500: loss 1.8238967675441229
iteration 1300 / 1500: loss 1.7784640840991994
iteration 1400 / 1500: loss 1.7210576127921302
That took 37.584728479385376s
```



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

```
In [23]: ## Implement softmax.predict() and use it to compute the training and testing error.

y_train_pred = softmax.predict(X_train)
print(len(y_train_pred))
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)),
))

49000
training accuracy: 0.3824489795918367
validation accuracy: 0.38
```

## Optimize the softmax classifier

You may copy and paste your optimization code from the SVM here.

```
In [24]: np.finfo(float).eps
```

```
Out[24]: 2.220446049250313e-16
```

```

In [26]: # ===== #
# 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.
# ===== #

my_time = time.time()
learning_rates = np.linspace(0, 1e-5, 20)
y_val_accs = []
best_learning_rate = -1
best_val_acc = -1

for rate in learning_rates:
    loss_hist = softmax.train(X_train, y_train, learning_rate=rate, num_iters=
1500, verbose=False)
    y_train_pred = softmax.predict(X_train)
    train_acc = np.mean(np.equal(y_train, y_train_pred))
    y_val_pred = softmax.predict(X_val)
    y_val_acc = np.mean(np.equal(y_val, y_val_pred))
    y_val_accs.append(y_val_acc)
    if y_val_acc > best_val_acc:
        best_val_acc = y_val_acc
        best_learning_rate = rate

print("Best learning rate: ", best_learning_rate, " Best Accuracy: ", best_val
_acc, "Err: ", 1 - best_val_acc)

loss_hist = softmax.train(X_train, y_train, learning_rate=best_learning_rate,
num_iters=1500, verbose=False)
y_test_pred = softmax.predict(X_test)
test_acc = np.mean(np.equal(y_test, y_test_pred))

plt.plot(learning_rates, y_val_accs)
plt.xlabel('Learning Rate')
plt.ylabel('Accuracy')

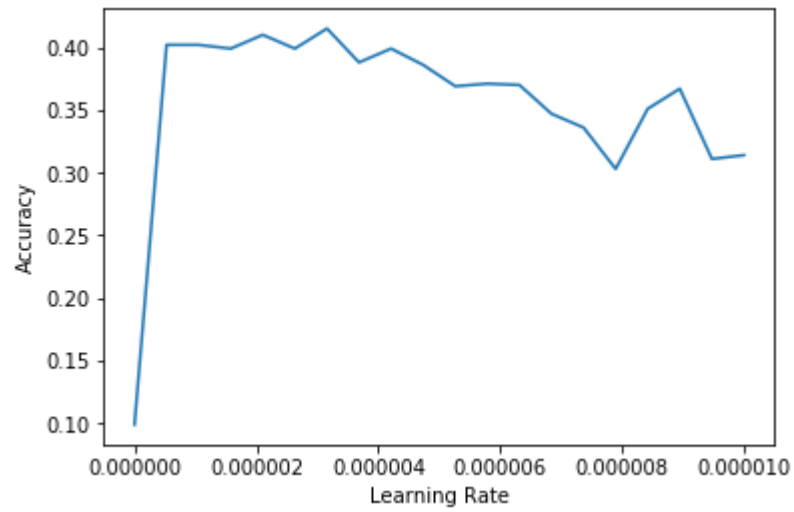
print("Test Acc: ", test_acc, " Test Error: ", 1-test_acc)
# ===== #
# END YOUR CODE HERE
# ===== #

print("Time took {}s".format(time.time() - my_time))

```



Best learning rate: 3.1578947368421056e-06 Best Accuracy: 0.415 Err: 0.58  
5  
Test Acc: 0.378 Test Error: 0.622  
Time took 419.63290762901306s



In [ ]:

```

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.)
        # ===== #
        num_train = X.shape[0]
        num_classes = max(y) + 1
        for i in range(num_train):
            ayi = np.dot(self.W[y[i]].T, X[i])
            maj = 0
            for j in range(num_classes):
                maj += np.exp(np.dot(self.W[j].T, X[i]))
            aj = np.log(maj)
            p = -ayi + aj
            loss += p
        loss /= num_train

        # ===== #
        # 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.*

*# ===== #*

num\_train = X.shape[0]

num\_classes = max(y) + 1

for i in range(num\_train):

    ayi = np.dot(self.W[y[i]].T, X[i])

    maj = 0

    for j in range(num\_classes):

        maj += np.exp(np.dot(self.W[j].T, X[i]))

    aj = np.log(maj)

    p = -ayi + aj

    loss += p

    for j in range(num\_classes):

        a = np.exp(np.dot(self.W[j].T, X[i]))

        grad[j] += -X[i] \* ((j == y[i]) - a / maj)

loss /= num\_train

grad /= num\_train

*# ===== #*

*# 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.
# ===== #
num_train = X.shape[0]
num_classes = max(y) + 1

score = np.dot(X, self.W.T)
score = np.exp(score - score.max())
denominator = np.sum(score, axis = 1)
num = score[range(num_train), y]
loss = -np.sum(np.log(num / denominator))/num_train
s = np.divide(score, denominator.reshape(num_train, 1))
s[range(num_train), y] = -(denominator - num) / denominator
grad = np.dot(s.T, X) / num_train

# ===== #
# 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: (dim, batch_size)
        #   - 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.
# ===== #
a = list(range(len(X)))
indx = np.random.choice(a, size = batch_size, replace=False)
X_batch = np.vstack([X[i] for i in indx])
y_batch = [y[i] for i in indx]

# ===== #
# 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 - 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.
    # ===== #
    y_pred = np.argmax(np.dot(X, self.W.T), axis=1)
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    return y_pred

```