# 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

```python
In [77]:  import numpy as np # for doing most of our calculations
          import matplotlib.pyplot as plt# for plotting
          from utils.data_utils import load_CIFAR10 # function to load the CIFAR-10 dataset.

          # Load matplotlib images inline
          %matplotlib inline

          # These are important for reloading any code you write in external .py files.
          # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
          %load_ext autoreload
          %autoreload 2
```
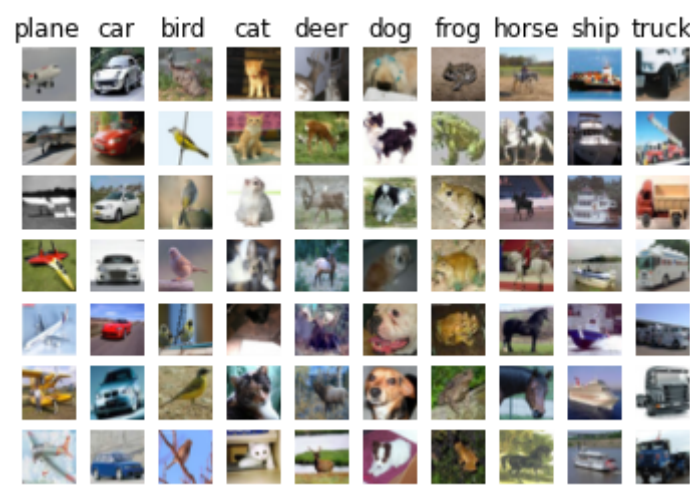
```
The autoreload extension is already loaded. To reload it, use:
  %reload_ext autoreload
```

```python
In [78]:  # 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,)
```

```python
In [79]:  # 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()
```



```python
In [80]:  # 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 [81]:
```python
# Import the KNN class

from nndl import KNN
```

In [82]:
```python
# 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're just storing the training data (X_train and y_train) in memory

(2) Pros - Easy to implement, training process is fast as there is no processing; Cons - Memory instensive, we need to store all the input 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 [83]:
```python
# 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: 22.104963064193726
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 [87]:
```python
# 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(dis
```

```
Time to run code: 0.21875691413879395
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 [88]:
```python
# 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, k=1)
error = np.mean(y_pred != 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.

In [89]:
```python
# Create the dataset folds for cross-valdiation.
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]
# ================================================================ #
fold_size = X_train.shape[0] // num_folds
for fold in range(num_folds):
    X_train_folds.append(X_train[fold*fold_size:(fold+1)*fold_size])
    y_train_folds.append(y_train[fold*fold_size:(fold+1)*fold_size])


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

## 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 [115...]:
```python
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.
# ================================================================ #
```
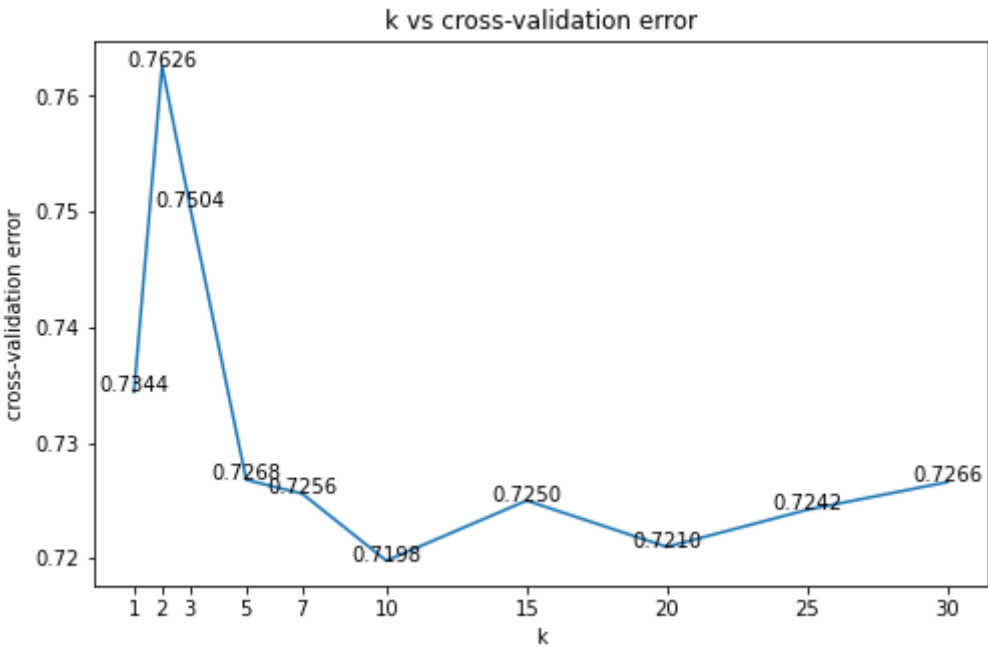
```
errors = [0] * len(ks)
for j in range(len(ks)):
#     print (ks[j])
    for i in range(len(X_train_folds)):
        X_train_ = np.concatenate(X_train_folds[:i] + X_train_folds[i+1:])
        y_train_ = np.concatenate(y_train_folds[:i] + y_train_folds[i+1:])
        knn.train(X=X_train_, y=y_train_)
        dists_L2_vectorized = knn.compute_L2_distances_vectorized(X=X_train_folds[i])
        y_pred = knn.predict_labels(dists_L2_vectorized, k=ks[j])
        errors[j] += np.mean(y_pred != y_train_folds[i])
    errors[j] /= len(X_train_folds)
# Plot
plt.figure(figsize=(8,5))
plt.title("k vs cross-validation error")
plt.xticks(ks)
plt.plot(ks, errors, label="cross-validation error")
plt.xlabel("k")
plt.ylabel("cross-validation error")

for i in range(len(ks)):
    plt.text(ks[i], errors[i], "{:.4f}".format(errors[i]), ha = 'center')

plt.show()


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

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



```
Computation time: 26.01
```

## 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) k = 10

(2) 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.

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

errors = [0] * len(norms)
```
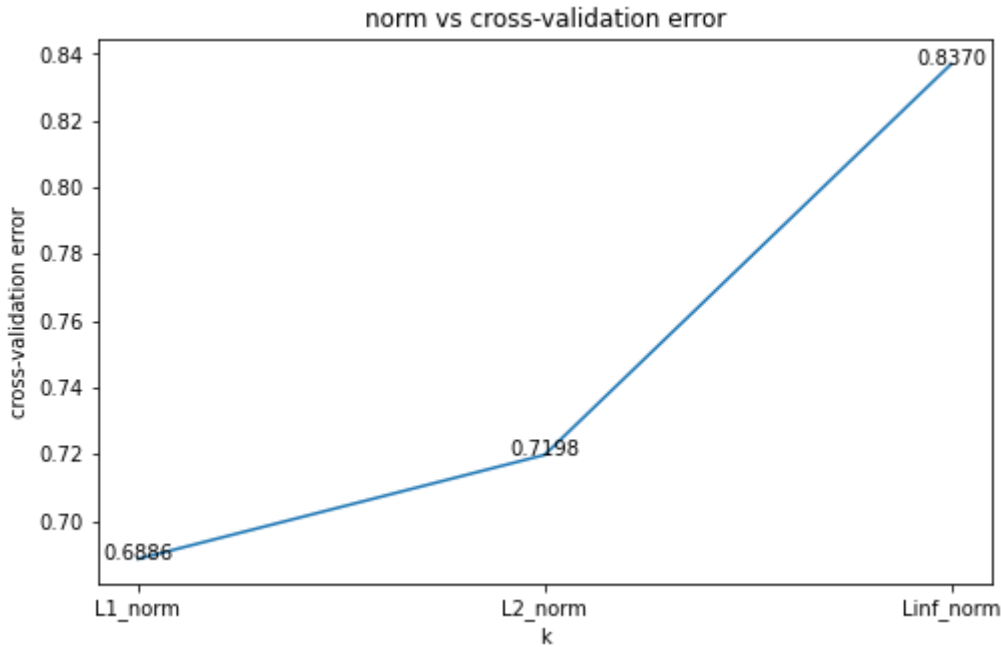
```
for j in range(len(norms)):
    for i in range(len(X_train_folds)):
        X_train_ = np.concatenate(X_train_folds[:i] + X_train_folds[i+1:])
        y_train_ = np.concatenate(y_train_folds[:i] + y_train_folds[i+1:])
        knn.train(X=X_train_, y=y_train_)
        dists = knn.compute_distances(X=X_train_folds[i], norm=norms[j])
        y_pred = knn.predict_labels(dists, k=10)
        errors[j] += np.mean(y_pred != y_train_folds[i])
    errors[j] /= len(X_train_folds)

# Plot
plt.figure(figsize=(8,5))
plt.title("norm vs cross-validation error")
plt.xticks(np.arange(3), ['L1_norm', 'L2_norm', 'Linf_norm'])
plt.plot(range(len(norms)), errors, label="cross-validation error")
plt.xlabel("k")
plt.ylabel("cross-validation error")

for i in range(len(norms)):
    plt.text(i, errors[i], "{:.4f}".format(errors[i]), ha = 'center')

plt.show()

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


norm vs cross-validation error

Computation time: 461.55

## 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) L1 norm

(2) 0.6886

# 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 [117...

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

knn.train(X=X_train, y=y_train)
dists = knn.compute_distances(X=X_test, norm=L1_norm)
y_pred = knn.predict_labels(dists, k=10)
error = np.mean(y_pred != y_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:

0.004

In [ ]:

```python
import numpy as np
import pdb


class KNN(object):

  def __init__(self):
    pass

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

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

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

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

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

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

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


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

    return dists

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

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

    Returns:
    - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
      is the Euclidean distance between the ith test point and the jth training
      point.
    """
    num_test = X.shape[0]
    num_train = self.X_train.shape[0]
    dists = np.zeros((num_test, num_train))

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

    X_norm = np.sum(np.square(X), axis=1)
    X_norm = X_norm.reshape(X_norm.shape[0], 1)

    X_train_norm = np.sum(np.square(self.X_train), axis=1)

    X_dot_X_train = X @ (self.X_train).T
    dists = np.sqrt(X_norm + X_train_norm - 2*X_dot_X_train)

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

```python
    # 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 betwen the ith test point and the jth training point.

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

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

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

    return y_pred
```

Samples: $(x^{(1)}, y^{(1)}), \ldots\ldots, (x^{(m)}, y^{(m)})$

$\qquad x^{(j)} \in \mathbb{R}^n$

$\qquad y^{(j)} \in \{1, \ldots, c\}$

Parameters: $\theta = \{w_i, b_i\}_{i=1,\ldots,c}$

Model: $\Pr(y^{(j)} = i \mid x^{(j)}, \theta) = \text{Softmax}(x^{(j)})$

$\qquad = \dfrac{e^{w_i^T x^{(j)} + b_i}}{\sum\limits_{k=1}^{c} e^{w_k^T x^{(j)} + b_k}}$

Likelihood: $P(x^{(1)}, \ldots, x^{(m)}, y^{(1)}, \ldots, y^{(m)} \mid \theta) = \prod\limits_{j=1}^{m} P(x^{(j)}, y^{(j)} \mid \theta)$

$\qquad\qquad\qquad\qquad = \prod\limits_{j=1}^{m} P(x^{(j)} \mid \theta) \, P(y^{(j)} \mid x^{(j)}, \theta)$

$\underset{\theta}{\arg\max} \, (\text{Likelihood}) = \underset{\theta}{\arg\max} \prod\limits_{j=1}^{m} P(x^{(j)} \mid \theta) \, P(y^{(j)} \mid x^{(j)}, \theta)$

$\qquad\qquad\qquad = \underset{\theta}{\arg\max} \prod\limits_{j=1}^{m} P(y^{(j)} \mid x^{(j)}, \theta)$

$\qquad\qquad\qquad = \underset{\theta}{\arg\max} \sum\limits_{j=1}^{m} \log\left[ P(y^{(j)} \mid x^{(j)}, \theta) \right]$

$\qquad\qquad\qquad = \underset{\theta}{\arg\max} \sum\limits_{j=1}^{m} \log\left[ \dfrac{e^{a_{y^{(j)}}(x^{(j)})}}{\sum\limits_{k=1}^{c} e^{a_k(x^{(j)})}} \right]$

Let:

$\boxed{a_i(x^{(j)}) = w_i^T x^{(j)} + b_i}$

$\qquad\qquad\qquad = \underset{\theta}{\arg\max} \dfrac{1}{m} \sum\limits_{j=1}^{m} \log\left[ \dfrac{e^{a_{y^{(j)}}(x^{(j)})}}{\sum\limits_{k=1}^{c} e^{a_k(x^{(j)})}} \right]$ ⟶ Log Likelihood $\mathcal{L}$

Hence to get the optimum parameters, we need to maximise the log likelihood

$\qquad \underset{\theta}{\max} \, f(\theta) = \underset{\theta}{\min} \, -f(\theta)$

⟹ Our loss function is negative of log likelihood

$\mathcal{L} = \dfrac{1}{m} \sum\limits_{j=1}^{m} \log\left[ \dfrac{e^{a_{y^{(j)}}(x^{(j)})}}{\sum\limits_{k=1}^{c} e^{a_k(x^{(j)})}} \right]$

Let,

$\mathcal{L}_j = \log\left[ \dfrac{e^{a_{y^{(j)}}(x^{(j)})}}{\sum\limits_{k=1}^{c} e^{a_k(x^{(j)})}} \right] = \left[ a_{y^{(j)}}(x^{(j)}) - \log \sum\limits_{k=1}^{c} e^{a_k(x^{(j)})} \right]$

Let :

$$\sigma_{j^{(i)}}\left(x^{(i)}\right) = \frac{e^{a_{j^{(i)}}\left(x^{(i)}\right)}}{\sum\limits_{k=1}^{c} e^{a_k\left(x^{(i)}\right)}}$$

using : $\left[\dfrac{f(x)}{g(x)}\right]' = \dfrac{g(x)f'(x) - f(x)g'(x)}{[g(x)]^2}$

we get,

$$\frac{\partial \sigma_{y^{(i)}}\left(x^{(i)}\right)}{\partial a_i\left(x^{(i)}\right)} = \begin{cases} \dfrac{\sum\limits_{k=1}^{c} e^{a_k\left(x^{(i)}\right)} \, e^{a_{y^{(i)}}\left(x^{(i)}\right)} - e^{2 a_{y^{(i)}}\left(x^{(i)}\right)}}{\left(\sum\limits_{k=1}^{c} e^{a_k\left(x^{(i)}\right)}\right)^2} & i = y^{(i)} \\[20pt] \dfrac{e^{a_{y^{(i)}}\left(x^{(i)}\right)} \, e^{a_i\left(x^{(i)}\right)}}{\left(\sum\limits_{k=1}^{c} e^{a_k\left(x^{(i)}\right)}\right)^2} & i \neq y^{(i)} \end{cases}$$

$$= \begin{cases} \dfrac{e^{a_{y^{(i)}}\left(x^{(i)}\right)}}{\sum\limits_{k=1}^{c} e^{a_k\left(x^{(i)}\right)}} - \left(\dfrac{e^{a_{y^{(i)}}\left(x^{(i)}\right)}}{\sum\limits_{k=1}^{c} e^{a_k\left(x^{(i)}\right)}}\right)^2 & i = y^{(i)} \\[20pt] \dfrac{e^{a_{y^{(i)}}\left(x^{(i)}\right)} \, e^{a_i\left(x^{(i)}\right)}}{\left(\sum\limits_{k=1}^{c} e^{a_k\left(x^{(i)}\right)}\right)^2} & i \neq y^{(i)} \end{cases}$$

$$\frac{\partial \sigma_{y^{(i)}}\left(x^{(i)}\right)}{\partial a_i\left(x^{(i)}\right)} = \begin{cases} \sigma_{y^{(i)}}\left(x^{(i)}\right)\left[1 - \sigma_{y^{(i)}}\left(x^{(i)}\right)\right] & i = y^{(i)} \\[8pt] -\sigma_{y^{(i)}}\left(x^{(i)}\right)\,\sigma_i\left(x^{(i)}\right) & i \neq y^{(i)} \end{cases}$$

using chain Rule :

$$\frac{\partial L_i}{\partial a_i\left(x^{(i)}\right)} \quad \frac{\partial \log\left[\sigma_{y^{(i)}}\left(x^{(i)}\right)\right]}{\partial a_i\left(x^{(i)}\right)} = \frac{\partial \log\left[\sigma_{y^{(i)}}\left(x^{(i)}\right)\right]}{\partial \sigma_{y^{(i)}}\left(x^{(i)}\right)} \frac{\partial \sigma_{y^{(i)}}\left(x^{(i)}\right)}{\partial a_i\left(x^{(i)}\right)}$$

$$= \frac{1}{\sigma_{y^{(i)}}\left(x^{(i)}\right)} \frac{\partial \sigma_{y^{(i)}}\left(x^{(i)}\right)}{\partial a_i\left(x^{(i)}\right)}$$

$$= \begin{cases} 1 - \sigma_{y^{(i)}}\left(x^{(i)}\right) & i = y^{(i)} \\[8pt] -\sigma_i\left(x^{(i)}\right) & i \neq y^{(i)} \end{cases}$$

Now :

$$\frac{\partial L_i}{\partial w_i} = \frac{\partial \log\left[\sigma_{y^{(i)}}\left(x^{(i)}\right)\right]}{\partial w_i} = \frac{\partial \log\left[\sigma_{y^{(i)}}\left(x^{(i)}\right)\right]}{\partial \sigma_{y^{(i)}}\left(x^{(i)}\right)} \frac{\partial \sigma_{y^{(i)}}\left(x^{(i)}\right)}{\partial a_i\left(x^{(i)}\right)} \frac{\partial a_i\left(x^{(i)}\right)}{\partial w_i}$$

$$= \frac{\partial \log\left[\sigma_{y^{(n)}}(x^{(n)})\right]}{\partial \sigma_{y^{(n)}}(x^{(n)})} \cdot \frac{\partial \sigma_{y^{(n)}}(x^{(n)})}{\partial a_i(x^{(n)})} \, x^{(n)}$$

$$= \begin{cases} \left[1 - \sigma_{y^{(n)}}(x^{(n)})\right] x^{(n)} & i = y^{(n)} \\[2mm] \left[-\sigma_i(x^{(n)})\right] x^{(n)} & i \neq y^{(n)} \end{cases}$$

$$\frac{\partial \ell_j}{\partial b_v} = \frac{\partial \log\left[\sigma_{y^{(n)}}(x^{(n)})\right]}{\partial b_i} = \frac{\partial \log\left[\sigma_{y^{(n)}}(x^{(n)})\right]}{\partial \sigma_{y^{(n)}}(x^{(n)})} \cdot \frac{\partial \sigma_{y^{(n)}}(x^{(n)})}{\partial a_i(x^{(n)})} \cdot \frac{\partial a_i(x^{(n)})}{\partial b_i}$$

$$= \frac{\partial \log\left[\sigma_{y^{(n)}}(x^{(n)})\right]}{\partial \sigma_{y^{(n)}}(x^{(n)})} \cdot \frac{\partial \sigma_{y^{(n)}}(x^{(n)})}{\partial a_i(x^{(n)})} \quad (1)$$

$$= \begin{cases} \left[1 - \sigma_{y^{(n)}}(x^{(n)})\right] & i = y^{(n)} \\[2mm] \left[-\sigma_i(x^{(n)})\right] & i \neq y^{(n)} \end{cases}$$

# 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.

```python
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
```

```python
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.

```python
from nndl import Softmax
```

```python
# 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

```python
## Implement the loss function of the softmax using a for loop over
#   the number of examples

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

```python
print(loss)
```

```
2.327760702804897
```

## Question:

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

## Answer:

Initially all weights are close to 0 (0.0001) so e^w becomes 1; Sofmax becomes 1/10; and -log(softmax) becomes 2.3. In loss we're taking a mean of for all examples, which are 2.3. Hence the average loss is 2.3

### Softmax gradient

```python
## 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 correct
softmax.grad_check_sparse(X_dev, y_dev, grad)
```

```
numerical: 0.394800 analytic: 0.394800, relative error: 3.192154e-08
numerical: -0.559981 analytic: -0.559981, relative error: 8.947427e-10
numerical: -1.673186 analytic: -1.673186, relative error: 1.066596e-08
numerical: 1.881785 analytic: 1.881785, relative error: 3.070772e-10
numerical: 0.921359 analytic: 0.921359, relative error: 3.323017e-08
numerical: 1.654556 analytic: 1.654556, relative error: 5.940917e-09
numerical: -0.069415 analytic: -0.069415, relative error: 2.468808e-07
numerical: -0.658161 analytic: -0.658161, relative error: 4.087038e-09
numerical: 1.061699 analytic: 1.061699, relative error: 2.079294e-08
numerical: -3.228319 analytic: -3.228319, relative error: 1.398503e-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.

```python
import time
```

```python
## 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,

# 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: 2.3312721708461672 / 384.9360658044393 computed in 0.056385040283203125s
```

```
Vectorized loss / grad: 2.3312721708461686 / 384.93606580443924 computed in 0.002574920654296875s
difference in loss / grad: -1.3322676295501878e-15 /2.607030649150689e-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:

In SVM, we have a different loss function compared to softmax. So gradient calculation will change but gradient decent will remain same

In [212...
```python
# 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 2.3365926606637544
iteration 100 / 1500: loss 2.0557222613850827
iteration 200 / 1500: loss 2.0357745120662813
iteration 300 / 1500: loss 1.9813348165609888
iteration 400 / 1500: loss 1.9583142443981612
iteration 500 / 1500: loss 1.8622653073541355
iteration 600 / 1500: loss 1.853261145435938
iteration 700 / 1500: loss 1.8353062223725827
iteration 800 / 1500: loss 1.829389246882764
iteration 900 / 1500: loss 1.8992158530357484
iteration 1000 / 1500: loss 1.97835035402523
iteration 1100 / 1500: loss 1.8470797913532633
iteration 1200 / 1500: loss 1.8411450268664082
iteration 1300 / 1500: loss 1.7910402495792102
iteration 1400 / 1500: loss 1.8705803029382257
That took 3.5640828609466553s
```



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

In [213...
```python
## 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: 0.3811428571428571
validation accuracy: 0.398
```

## Optimize the softmax classifier

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

In [214...
```python
np.finfo(float).eps
```

Out[214... 
```
2.220446049250313e-16
```

In [230...
```python
# =============================================================== #
```

```python
# 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 = [1e-9, 1e-8, 5e-7, 2e-7, 1e-7, 8e-6, 5e-6, 1e-6, 1e-5]
validation_accuracies = []

for learning_rate in learning_rates:
    tic = time.time()
    loss_hist = softmax.train(X_train, y_train, learning_rate=learning_rate,
                              num_iters=1500, verbose=False)
    toc = time.time()
    print('That took {}s'.format(toc - tic))

    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)), ))
    validation_accuracies.append(np.mean(np.equal(y_val, y_val_pred)))

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


best_learning_rate = learning_rates[np.argmax(validation_accuracies)]
best_validation_accuracy = np.max(validation_accuracies)
print('best validation accuracy: {}'.format(best_validation_accuracy))
print('best validation error: {}'.format(1-best_validation_accuracy))
print('best learning rate: {}'.format(best_learning_rate))

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)
print('test accuracy: {}'.format(np.mean(np.equal(y_test, y_test_pred)), ))
print('test error: {}'.format(1-np.mean(np.equal(y_test, y_test_pred)), ))


# ================================================================ #
# END YOUR CODE HERE
# ================================================================ #
```
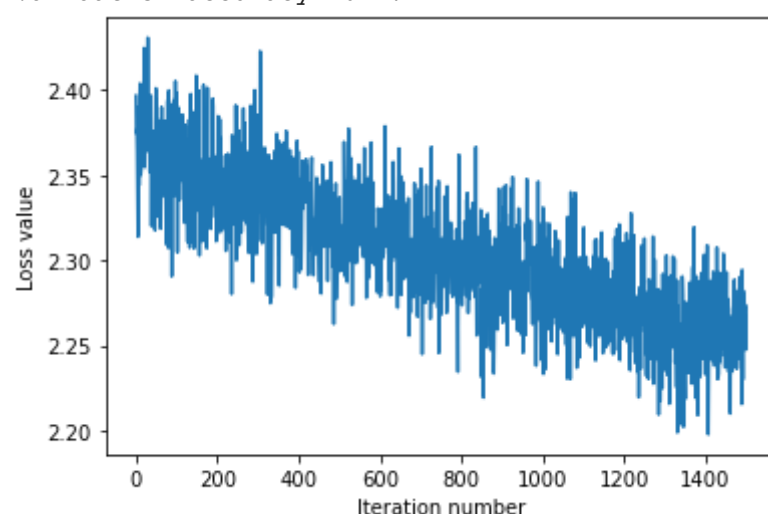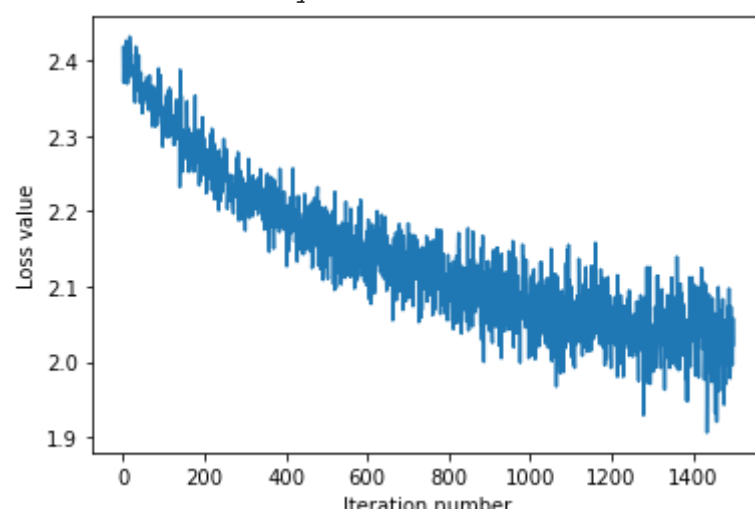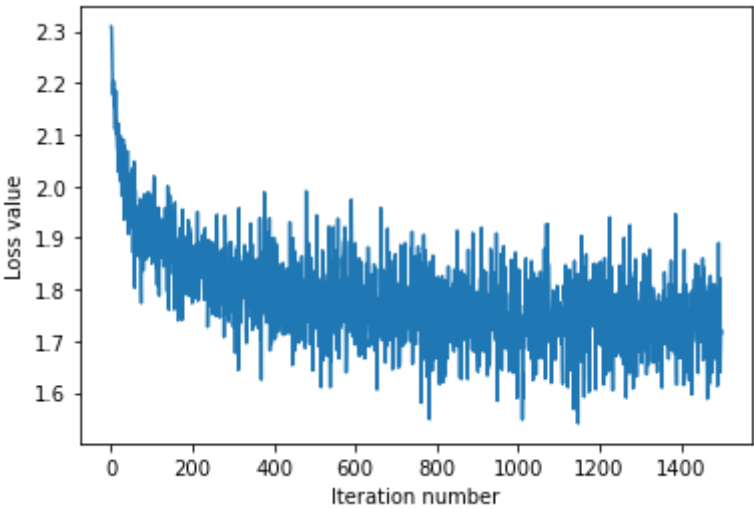
```
That took 5.370471000671387s
training accuracy: 0.15177551020408164
validation accuracy: 0.174
```
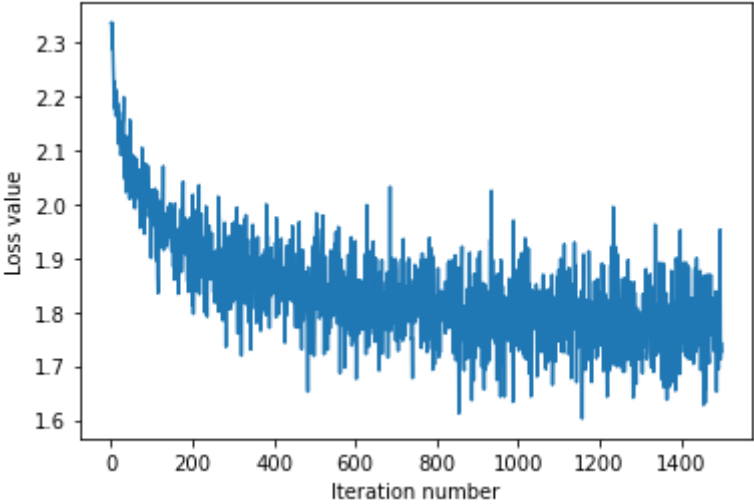


```
That took 2.529613971710205s
training accuracy: 0.2861020408163265
validation accuracy: 0.303
```
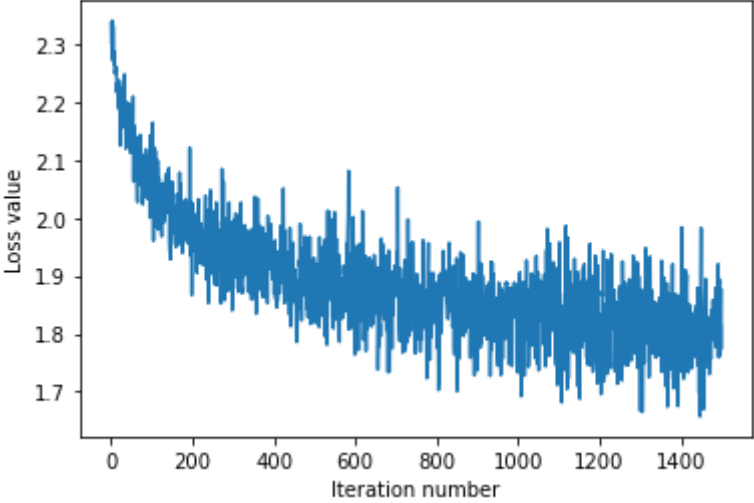


```
That took 2.3880419731140137s
training accuracy: 0.4126326530612245
validation accuracy: 0.42
```
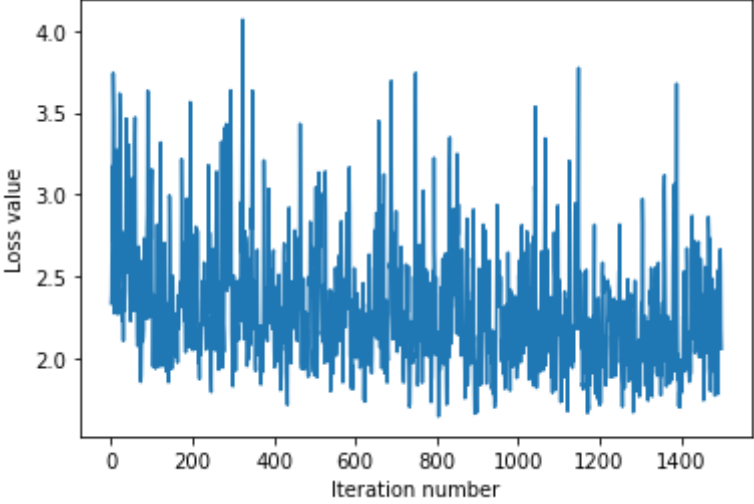
That took 2.314687967300415s
training accuracy: 0.39516326530612245
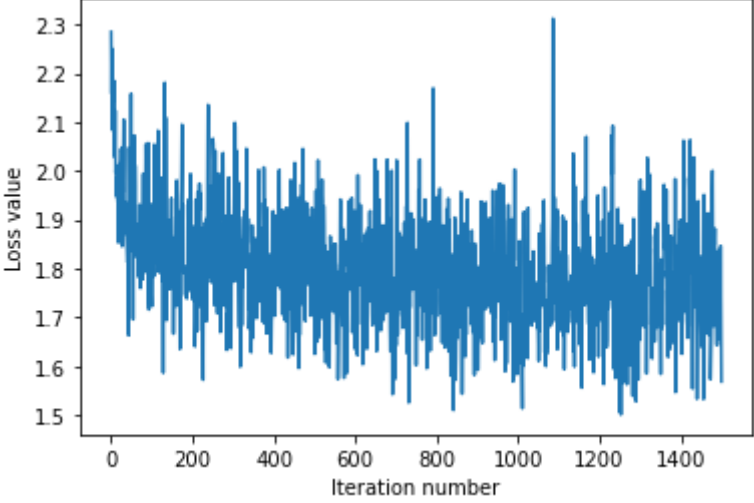validation accuracy: 0.399



That took 2.2825448513031006s
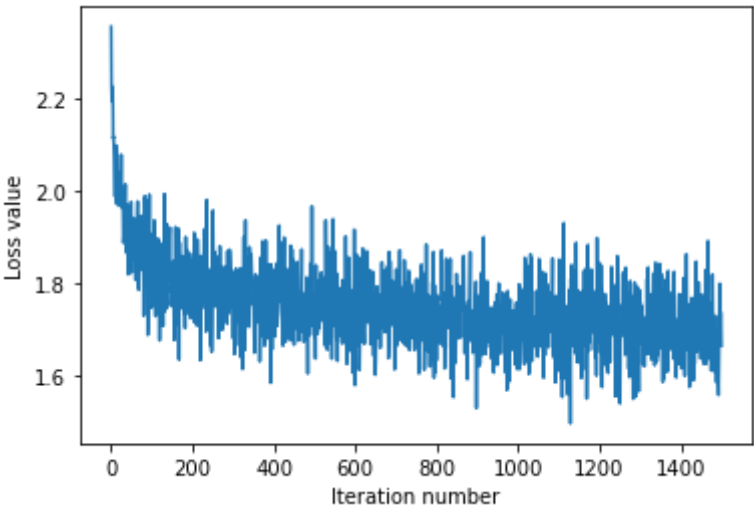training accuracy: 0.3826530612244898
validation accuracy: 0.389



That took 6.018434047698975s
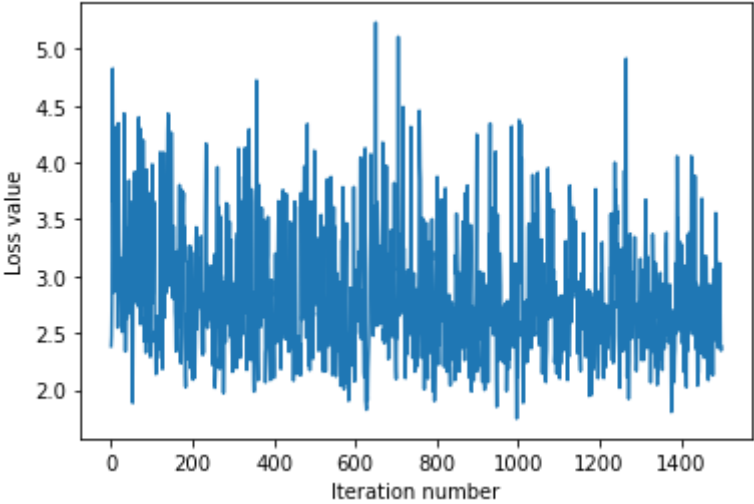training accuracy: 0.37620408163265306
validation accuracy: 0.354



That took 2.9429879188537598s
training accuracy: 0.41459183673469385
validation accuracy: 0.39



That took 3.5817999839782715s
training accuracy: 0.42136734693877553
validation accuracy: 0.407

```
That took 4.134116172790527s
training accuracy: 0.3256734693877551
validation accuracy: 0.306
```



```
best validation accuracy: 0.42
best validation error: 0.5800000000000001
best learning rate: 5e-07
test accuracy: 0.393
test error: 0.607
```

In [ ]:

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

    m = X.shape[0]
    for j in range(m):
      z = self.W @ X[j]
      e_x = np.exp(z - np.max(z))
      softmax = e_x / np.sum(e_x)
      loss += -np.log(softmax[y[j]])
    loss /= m

    # ================================================================ #
    # 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.
    # ================================================================ #
    m = X.shape[0]
    for j in range(m):
      z = self.W @ X[j]
      e_x = np.exp(z - np.max(z))
      softmax = e_x / np.sum(e_x)
      loss += -np.log(softmax[y[j]])

      for i in range(len(softmax)):
        if (y[j] == i):
          grad[i] += -(1-softmax[i]) * X[j]
        else:
          grad[i] += softmax[i] * X[j]

    loss /= m
    grad /= m
```

```python
 91          # ================================================================ #
 92          # END YOUR CODE HERE
 93          # ================================================================ #

 95          return loss, grad

 97      def grad_check_sparse(self, X, y, your_grad, num_checks=10, h=1e-5):
 98          """
 99          sample a few random elements and only return numerical
100          in these dimensions.
101          """

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

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

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

118      def fast_loss_and_grad(self, X, y):
119          """
120          A vectorized implementation of loss_and_grad. It shares the same
121          inputs and ouptuts as loss_and_grad.
122          """
123          loss = 0.0
124          grad = np.zeros(self.W.shape) # initialize the gradient as zero

126          # ================================================================ #
127          # YOUR CODE HERE:
128          #   Calculate the softmax loss and gradient WITHOUT any for loops.
129          # ================================================================ #

131          z = X @ self.W.T
132          max_z = np.max(z, axis=1, keepdims=True)
133          e_x = np.exp(z - max_z)
134          e_x_sum = np.sum(e_x, axis=1, keepdims=True)
135          softmax = e_x / e_x_sum
136          loss = -np.log(np.choose(y, softmax.T))
137          loss = np.mean(loss)


140          softmax[np.arange(X.shape[0]),y] -= 1
141          grad += (X.T @ softmax).T
142          grad /= X.shape[0]



146          # ================================================================ #
147          # END YOUR CODE HERE
148          # ================================================================ #

150          return loss, grad

152      def train(self, X, y, learning_rate=1e-3, num_iters=100,
153                batch_size=200, verbose=False):
154          """
155          Train this linear classifier using stochastic gradient descent.

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

167          Outputs:
168          A list containing the value of the loss function at each training iteration.
169          """
170          num_train, dim = X.shape
171          num_classes = np.max(y) + 1 # assume y takes values 0...K-1 where K is number of classes

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

175          # Run stochastic gradient descent to optimize W
176          loss_history = []

178          for it in np.arange(num_iters):
179              X_batch = None
180              y_batch = None
181
```

```python
        # ========================================================= #
        # 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.
        # ========================================================= #
        idx = np.random.randint(X.shape[0], size=batch_size)
        X_batch = X[idx, :]
        y_batch = y[idx]

        # ========================================================= #
        # 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.
    # ========================================================= #
    z = X @ self.W.T
    max_z = np.max(z, axis=1, keepdims=True)
    e_x = np.exp(z - max_z)
    e_x_sum = np.sum(e_x, axis=1, keepdims=True)
    softmax = e_x / e_x_sum
    y_pred = np.argmax(softmax, axis=1)
    # ========================================================= #
    # END YOUR CODE HERE
    # ========================================================= #

    return y_pred
```