

Brendan Rossmango
505 370 692

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

$$x^{(i)} \in \mathbb{R}^d \quad y^{(i)} \in \mathbb{R} \quad \theta \in \mathbb{R}^d \quad I \in \mathbb{R}^{d \times d} \quad \sigma \in \mathbb{R}$$

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

$$\delta^{(i)} \sim N(0, \sigma^2 I)$$

$$a) E_{\delta \sim N} [\tilde{L}(\theta)] = ?$$

$$= \frac{1}{N} \sum_{i=1}^N E_{\delta \sim N} [(y^{(i)} - (x^{(i)} + \delta^{(i)})^\top \theta)^2]$$

$$G = (y^{(i)} - (x^{(i)} + \delta^{(i)})^\top \theta)^2$$

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

$$= (y^{(i)} - x^{(i)\top} \theta)^2 - 2(y^{(i)} - x^{(i)\top} \theta)(\delta^{(i)\top} \theta) + (\delta^{(i)\top} \theta)^2$$

$E[\cdot]$ is linear so

$$E_{\delta \sim N} [(y^{(i)} - (x^{(i)} + \delta^{(i)})^\top \theta)^2]$$

$$= E[(y^{(i)} - x^{(i)\top} \theta)^2] + E_{\delta \sim N} [(-2(y^{(i)} - x^{(i)\top} \theta)(\delta^{(i)\top} \theta))]$$

$$+ E_{\delta \sim N} [(\delta^{(i)\top} \theta)^2]$$

$$E_{\delta \sim N} [(y^{(i)} - x^{(i)\top} \theta)^2] = \underline{(y^{(i)} - x^{(i)\top} \theta)^2} \quad \text{No } \delta \text{ dependence}$$

$$E_{\delta \sim N} [-2(y^{(i)} - x^{(i)\top} \theta)(\delta^{(i)\top} \theta)] = -2(y^{(i)} - x^{(i)\top} \theta) \cdot E_{\delta \sim N} [(\delta^{(i)\top} \theta)]$$

$$= 0 \quad E(\delta^{(i)}) = 0 \in \mathbb{R}^d \text{ since } \delta^{(i)} \sim N(0, \sigma^2 I)$$

$$E_{\delta \sim N} [(\delta^{(i)\top} \theta)^2] =$$

$$= E_{\delta \sim N} [\theta^\top \delta^{(i)} \delta^{(i)\top} \theta]$$

From the
Hint:

$$= \theta^\top E_{\delta \sim N} [\delta^{(i)} \delta^{(i)\top}] \theta$$

$$E_{\delta \sim N} [\delta \delta^\top] = \sigma^2 I$$

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

$$\text{All together: } (y^{(i)} - x^{(i)\top} \theta)^2 + \sigma^2 \|\theta\|_2^2$$

$$E_{\delta \sim N} (\tilde{L}(\theta)) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - x^{(i)\top} \theta)^2 + \sigma^2 \|\theta\|_2^2$$

$$= L(\theta) + R, \quad R = \sigma^2 \|\theta\|_2^2$$

- b) The noise addition would have a ridge regularization effect, where the regularization strength is the variance of the noise. This helps the model avoid overfitting by adding a penalty to shrink the coefficients and reduce model complexity.

- c) $\sigma \rightarrow 0$. There would be no regularization and the model might overfit the data.

- d) $\sigma \rightarrow \infty$. The objective of the cost function is to minimize the L^2 -norm of θ , so $\theta \rightarrow 0$, and the model underfits the data.

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

```
In [69]: 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
```

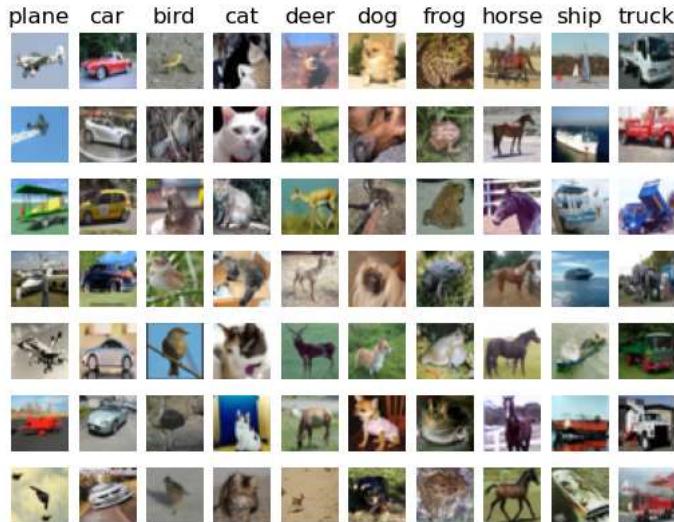
The autoreload extension is already loaded. To reload it, use:
%reload_ext autoreload

```
In [70]: # 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 [71]: # 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 [72]: # 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 [73]: # Import the KNN class
from nnndl import KNN
```

```
In [74]: # 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() simply remembers the data points (images) and their labels.
- (2) The pros are that there is no training time and it is simple and fast; the cons are that it is memory intensive as we must store all the input data. Additionally, predicting the test data is expensive since we must get the distance to all points.

KNN prediction

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

```
In [75]: # 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: 19.62835693359375
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 [76]: # 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)))

Time to run code: 0.11509299278259277
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 [79]: # 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 = 1 - np.mean(np.equal(y_test, y_pred))
# ===== #
# 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 [80]: # Create the dataset folds for cross-validation.
num_folds = 5

X_train_folds = []
y_train_folds = []

np.random.seed(123)
idx = np.random.permutation(num_training)
print(idx, idx.shape)
# ===== #
# 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_shuffle = X_train[idx]
y_train_shuffle = y_train[idx]
X_train_folds = np.array_split(X_train_shuffle, num_folds)
y_train_folds = np.array_split(y_train_shuffle, num_folds)

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

print(y_train_folds[0].shape)
```

[2648 2456 4557 ... 1346 3454 3582] (5000,)
(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 [81]: 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.
# ===== #

error_list = []
for k in ks:
    error = 0
    for i in np.arange(num_folds):
        X_fold_train = np.concatenate(X_train_folds[:i] + X_train_folds[i + 1:], axis=0)
        y_fold_train = np.concatenate(y_train_folds[:i] + y_train_folds[i + 1:], axis=0)
        X_fold_val = X_train_folds[i]
        y_fold_val = y_train_folds[i]

        knn.train(X_fold_train, y_fold_train)
        cur_pred = knn.predict_labels(knn.compute_L2_distances_vectorized(X_fold_val), k)

        error += (1 - np.mean(np.equal(y_fold_val, cur_pred)))
    error_list.append(error / num_folds)

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

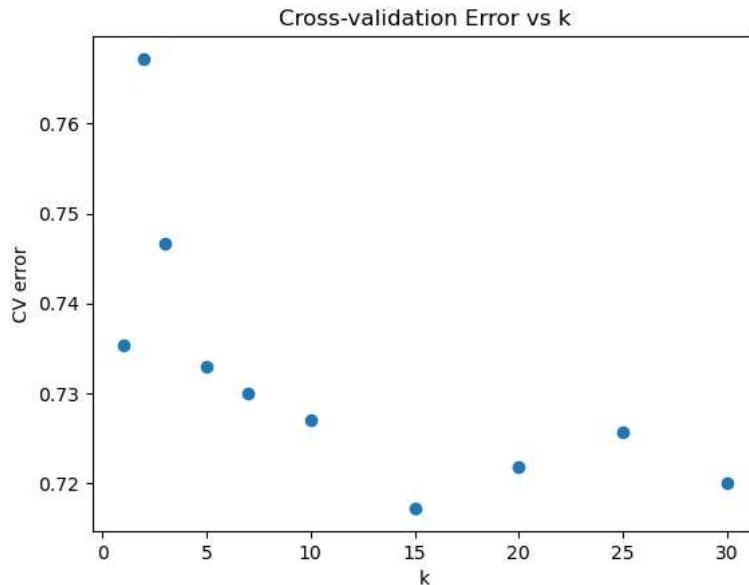
print('Computation time: %.2f' % (time.time() - time_start))
```

Computation time: 19.94

```
In [82]: plt.scatter(ks, error_list)
plt.title('Cross-validation Error vs k')
plt.ylabel('CV error')
plt.xlabel('k')
plt.show()

for i in range(len(error_list)):
    print(str(ks[i]) + ': ' + str(error_list[i]))

print()
best_error = min(error_list)
k_idx = error_list.index(best_error)
best_k = ks[k_idx]
print("Best k, error: " + str(best_k) + ', ' + str(best_error))
```



```
1: 0.7354
2: 0.7672
3: 0.7465999999999999
5: 0.733
7: 0.7300000000000001
10: 0.727
15: 0.7172
20: 0.7218
25: 0.7255999999999999
30: 0.72

Best k, error: 15, 0.7172
```

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) $k=15$
- (2) 0.7172

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 [83]: 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.
# ===== #

norm_error_list = []
for norm in norms:
    norm_error = 0
    for i in np.arange(num_folds):
        X_fold_train = np.concatenate(X_train_folds[:i] + X_train_folds[i + 1:], axis=0)
        y_fold_train = np.concatenate(y_train_folds[:i] + y_train_folds[i + 1:], axis=0)
        X_fold_val = X_train_folds[i]
        y_fold_val = y_train_folds[i]

        knn.train(X_fold_train, y_fold_train)
        cur_pred = knn.predict_labels(knn.compute_distances(X_fold_val, norm), best_k)

        norm_error += (1 - np.mean(np.equal(y_fold_val, cur_pred)))
    norm_error_list.append(norm_error / num_folds)

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

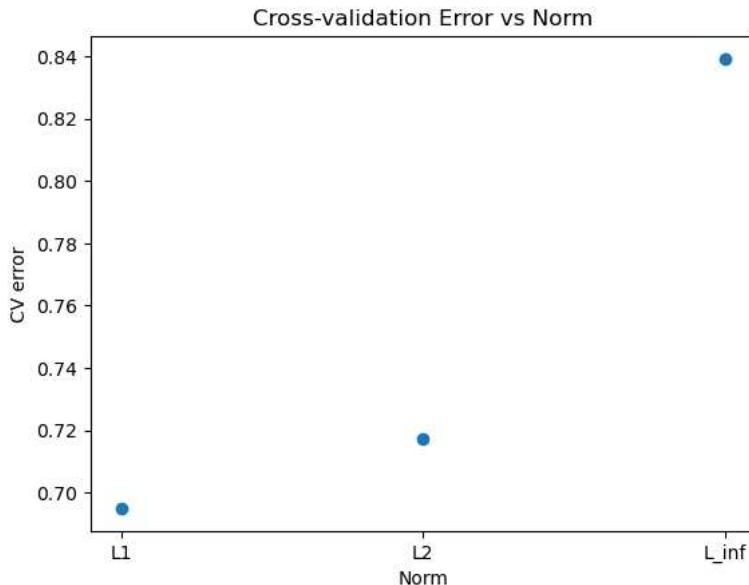
Computation time: 362.95

```
In [84]: ns = ["L1", "L2", "L_inf"]

plt.scatter(["L1", "L2", "L_inf"], norm_error_list)
plt.title('Cross-validation Error vs Norm')
plt.ylabel('CV error')
plt.xlabel('Norm')
plt.show()

for i in range(len(norm_error_list)):
    print(ns[i] + ' error: ' + str(norm_error_list[i]))

print()
best_error = min(norm_error_list)
norm_idx = norm_error_list.index(best_error)
print("Best norm, error: " + ns[norm_idx] + ', ' + str(best_error))
```



L1 error: 0.695
 L2 error: 0.7172
 L_inf error: 0.8392

Best norm, error: L1, 0.695

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) 0.695

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 [85]: 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 = KNN()
knn.train(X=X_train, y=y_train)
dists_L1 = knn.compute_distances(X=X_test, norm=L1_norm)
y_pred = knn.predict_labels(dists_L1, k=best_k)

error = 1 - np.mean(np.equal(y_test, y_pred))

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

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

Error rate achieved: 0.718

Question:

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

Answer:

Error with L2-norm and k=1: 0.726

Error with L1-norm and k=15: 0.718

Improvement: 0.008

knn.py related code sections

```
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.
    # ===== #
    # print(X.shape) # (500, 3072)
    # print(self.X_train.shape) # (5000, 3072)

    # X**2.sum(axis=1).reshape(-1,1): (N,1)
    # X_train**2.sum(axis=1): (M, )
    # X @ X_train.T: (500, 3072) x (3072, 5000) = (N, M)
    dists = np.sqrt((X**2).sum(axis=1).reshape(-1,1) - 2 * X @ self.X_train.T + (self.X_train**2).sum(axis=1))

    # ===== #
    # 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.
        # ===== #
        indices = np.argsort(dists[i])
        closest_y = self.y_train[indices][:k]
        closest_y = sorted(list(closest_y)) # sort by label number, then get most common label
        y_pred[i] = max(closest_y, key = closest_y.count)

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

    return y_pred

```

$\exists (x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)}) \in \mathbb{R}^n, y^{(j)} \in \{1, \dots, c\}$

$$\theta = \{w_i, b_i\}_{i=1, \dots, c}$$

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

$$\text{softmax}_i(x) = \frac{\exp(w_i^T x + b_i)}{\sum_{k=1}^c \exp(w_k^T x + b_k)}$$

$$L = ?$$

$$\nabla_{w_i} L = ?, \quad \tilde{w}_i = \begin{bmatrix} w_i \\ b_i \end{bmatrix}, \quad \tilde{x} = \begin{bmatrix} x \\ 1 \end{bmatrix}, \quad a_i(x) = \tilde{w}_i^T \tilde{x}$$

First get log likelihood L

$$\begin{aligned} \Pr(x^{(1)}, \dots, x^{(m)}, y^{(1)}, \dots, y^{(m)} | \theta) &= \prod_{i=1}^m p(x^{(i)}, y^{(i)} | \theta) \\ &= \prod_{i=1}^m p(x^{(i)} | \theta) p(y^{(i)} | x^{(i)}, \theta) \end{aligned}$$

↑ no dependence on θ

$$= \prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta) \underbrace{p(x^{(i)} | \theta)}_{= p(x^{(i)})}$$

$$L = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log (\text{softmax}_{y^{(i)}}(x^{(i)}))$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log \left[\frac{\exp(w_{y^{(i)}}^T x^{(i)} + b_{y^{(i)}})}{\sum_{k=1}^c \exp(w_k^T x^{(i)} + b_k)} \right] \rightarrow$$

$$= \arg \max_{\theta} \frac{1}{m} \sum_{i=1}^m \left[w_{y^{(i)}}^T x^{(i)} + b_{y^{(i)}} - \log \sum_{k=1}^c \exp(w_k^T x^{(i)} + b_k) \right]$$

$a_{y^{(i)}}(x^{(i)})$
normalize

$L(\theta) = \frac{1}{m} \sum_{i=1}^m \left[a_{y^{(i)}}(x^{(i)}) - \log \left(\sum_{k=1}^c e^{a_k(x^{(i)})} \right) \right]$

log-likelihood
(not negative)
log-likelihood

$$\frac{\partial L}{\partial \tilde{w}_j} = \frac{1}{m} \sum_{i=1}^m \left[\frac{\partial a_{y^{(i)}}(x^{(i)})}{\partial \tilde{w}_j} - \frac{\partial \left(\sum_{k=1}^c e^{a_k(x^{(i)})} \right)}{\partial \tilde{w}_j} \right]$$

Do j instead of i to not get confused with the other i's

nonzero only when $y^{(i)} = j$

nonzero only when $k = j$

$$= \frac{1}{m} \sum_{i=1}^m \left[1_{y^{(i)}=j} \cdot \frac{x^{(i)}}{\sum_{k=1}^c e^{a_k(x^{(i)})}} - \frac{e^{a_j(x^{(i)})}}{\sum_{k=1}^c e^{a_k(x^{(i)})}} \right]$$

$$\nabla_{\tilde{w}_j} L = \frac{1}{m} \sum_{i=1}^m \frac{x^{(i)}}{\sum_{k=1}^c e^{a_k(x^{(i)})}} \cdot \left(1_{y^{(i)}=j} - \frac{e^{a_j(x^{(i)})}}{\sum_{k=1}^c e^{a_k(x^{(i)})}} \right)$$

normalization term doesn't have to be present softmax;($x^{(i)}$)

$$4. \quad \mathcal{D} = \{(x^{(1)}, y^{(1)}), \dots, x^{(K)}, y^{(K)}\} \quad w \in \mathbb{R}^d, b \in \mathbb{R}$$

$x^{(i)} \in \mathbb{R}^d$
 $y^{(i)} \in \{-1, 1\}$

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

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

Not differentiable at $x=1$ if $y^{(i)}(w^T x^{(i)} + b) > 1$:

$$\nabla_w \text{hinge}_{y^{(i)}}(x^{(i)}) = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^d, \quad \nabla_b \text{hinge}_{y^{(i)}}(x^{(i)}) = 0 \in \mathbb{R}$$

$$\text{If } y^{(i)}(w^T x^{(i)} + b) \leq 1: \quad \nabla_w \text{hinge}_{y^{(i)}}(x^{(i)}) = -y^{(i)} x^{(i)} \in \mathbb{R}^d, \quad \nabla_b \text{hinge}_{y^{(i)}}(x^{(i)}) = -y^{(i)}$$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{1}{K} \sum_{i=1}^K \frac{\partial}{\partial w} \max(0, 1 - y^{(i)}(w^T x^{(i)} + b))$$

$$\nabla_w \mathcal{L} = \left\{ \begin{array}{ll} \frac{1}{K} \sum_{i=1}^K -y^{(i)} x^{(i)}, & y^{(i)}(w^T x^{(i)} + b) < 1 \\ \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}, & y^{(i)}(w^T x^{(i)} + b) \geq 1 \end{array} \right. \in \mathbb{R}^d$$

vector $\in \mathbb{R}^d$

$$\nabla_b \mathcal{L} = \left\{ \begin{array}{ll} \frac{1}{K} \sum_{i=1}^K -y^{(i)}, & y^{(i)}(w^T x^{(i)} + b) < 1 \\ 0, & y^{(i)}(w^T x^{(i)} + b) \geq 1 \end{array} \right. \in \mathbb{R}$$

Scalar

$$\nabla_w \mathcal{L} = \frac{1}{K} \sum_{i=1}^K \mathbf{1}_{\{y^{(i)}(w^T x^{(i)} + b) < 1\}} \cdot f(y^{(i)} x^{(i)}) \in \mathbb{R}^d$$

$\mathbf{1}_{\{y^{(i)}(w^T x^{(i)} + b) < 1\}}$ vector of 1s in \mathbb{R}^d , if < 1 , zero vector if ≥ 1

$$\nabla_b \mathcal{L} = \frac{1}{K} \sum_{i=1}^K \mathbf{1}_{\{y^{(i)}(w^T x^{(i)} + b) < 1\}} \cdot -y^{(i)} \in \mathbb{R}$$

$\mathbf{1}_{\{y^{(i)}(w^T x^{(i)} + b) < 1\}}$, 1 if < 1 , 0 if ≥ 1

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

Please follow the notebook linearly to implement a softmax classifier.

Please print out the workbook entirely when completed.

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

```
In [159]: 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

```
In [160]: 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)
```

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 [161]: from nnndl import Softmax
```

```
In [162]: # 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 [163]: ## Implement the Loss function of the softmax using a for Loop over
# the number of examples

loss = softmax.loss(X_train, y_train)

In [164]: 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:

There are 10 classes and $\ln(10) \approx 2.3$. It makes sense for the softmax loss to be about 2.3 because the weights W are initialized to a random value, so we have a 10% chance of choosing the right class.

Softmax gradient

```
In [165]: ## 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.
print(loss)
softmax.grad_check_sparse(X_dev, y_dev, grad)

2.336973608923815
numerical: 1.478083 analytic: 1.478083, relative error: 2.428878e-08
numerical: -0.086784 analytic: -0.086784, relative error: 1.494163e-07
numerical: 0.816771 analytic: 0.816771, relative error: 2.611513e-08
numerical: 1.503617 analytic: 1.503617, relative error: 2.402709e-09
numerical: 1.234142 analytic: 1.234142, relative error: 1.813902e-08
numerical: 1.110206 analytic: 1.110206, relative error: 1.141227e-08
numerical: 0.311683 analytic: 0.311683, relative error: 1.086696e-07
numerical: -1.192795 analytic: -1.192795, relative error: 1.390972e-08
numerical: 0.240884 analytic: 0.240884, relative error: 4.498970e-08
numerical: -3.824836 analytic: -3.824836, relative error: 1.157722e-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 [166]: import time
```

```
In [167]: ## Implement softmax.fast_loss_and_grad which calculates the Loss and gradient
# WITHOUT using any for Loops.

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

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

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

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

Normal loss / grad_norm: 2.336973608923815 / 360.8904953247931 computed in 0.056603193283081055s
Vectorized loss / grad: 2.336973608923816 / 360.89049532479316 computed in 0.0020036697387695312s
difference in loss / grad: -8.881784197001252e-16 / 2.7062638080135757e-13

Stochastic gradient descent

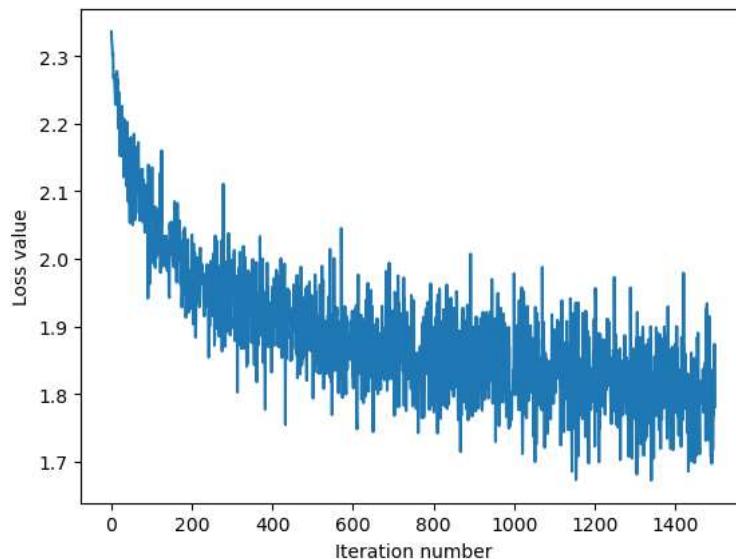
We now implement stochastic gradient descent. This uses the same principles of gradient descent we discussed in class, however, it calculates the gradient by only using examples from a subset of the training set (so each gradient calculation is faster).

```
In [168]: # 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 {}'.format(toc - tic))

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

```
iteration 0 / 1500: loss 2.336592660663754
iteration 100 / 1500: loss 2.0557222613850827
iteration 200 / 1500: loss 2.0357745120662813
iteration 300 / 1500: loss 1.9813348165609888
iteration 400 / 1500: loss 1.9583142443981614
iteration 500 / 1500: loss 1.8622653073541355
iteration 600 / 1500: loss 1.8532611454359387
iteration 700 / 1500: loss 1.8353062223725827
iteration 800 / 1500: loss 1.829389246882764
iteration 900 / 1500: loss 1.8992158530357484
iteration 1000 / 1500: loss 1.9783503540252303
iteration 1100 / 1500: loss 1.8470797913532635
iteration 1200 / 1500: loss 1.8411450268664082
iteration 1300 / 1500: loss 1.7910402495792102
iteration 1400 / 1500: loss 1.8705803029382257
That took 2.3441543579101562s
```



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

```
In [169]: ## 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

```
In [170]: np.finfo(float).eps
```

```
Out[170]: 2.220446049250313e-16
```

```
In [171]: # ===== #
# 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-10, 1e-9, 1e-8, 1e-7, 1e-6, 1e-5, 1e-4, 1e-3, 1e-2, 1e-1]
val_accs = {}
for r in learning_rates:
    softmax.train(X_train, y_train, learning_rate=r, num_iters=1500, verbose=False)
    y_val_pred = softmax.predict(X_val)
    acc = np.mean(np.equal(y_val, y_val_pred))
    val_accs[r] = acc

print(val_accs)

val_accs = sorted(val_accs.items(), key=lambda x : x[1])

best_learning_rate = val_accs[-1][0]
best_learning_rate_acc = val_accs[-1][1]

print("Best learning rate, validation accuracy:", max(val_accs, key = lambda x: x[1]))
print("Best validation error:", 1 - best_learning_rate_acc)

softmax.train(X_train, y_train, learning_rate=best_learning_rate, num_iters=1500, verbose=False)
y_test_pred = softmax.predict(X_test)
acc = np.mean(np.equal(y_test, y_test_pred))
error = 1 - acc
print("Error rate on test set with best learning rate:", error)
# ===== #
# END YOUR CODE HERE
# ===== #
```

{1e-10: 0.131, 1e-09: 0.122, 1e-08: 0.323, 1e-07: 0.388, 1e-06: 0.408, 1e-05: 0.292, 0.0001: 0.245, 0.001: 0.087, 0.01: 0.087, 0.1: 0.087}
Best learning rate, validation accuracy: (1e-06, 0.408)
Best validation error: 0.5920000000000001
Error rate on test set with best learning rate: 0.598

softmax.py related code sections

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

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

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

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

# Initialize the loss to zero.
loss = 0.0

# ===== #
# YOUR CODE HERE:
#   Calculate the normalized softmax loss. Store it as the variable loss.
#   (That is, calculate the sum of the losses of all the training
#   set margins, and then normalize the loss by the number of
#   training examples.)
# ===== #
N = X.shape[0]
# print(self.W.shape[0]) # 10

# print((X[0] @ self.W.T).shape)
for i in range(N):
    a = X[i] @ self.W.T
    a -= np.max(a) # normalize to avoid overflow
    sm = np.exp(a[y[i]]) / np.sum(np.exp(a))
    loss -= np.log(sm)

loss /= N
# ===== #
# END YOUR CODE HERE
# ===== #

return loss
```

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

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

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

    # ===== #
    # YOUR CODE HERE:
    # Calculate the softmax loss and the gradient. Store the gradient
    # as the variable grad.
    # ===== #
    N = X.shape[0]
    C = self.W.shape[0]
    # print(X.shape[0])

    for i in range(N):
        a = X[i] @ self.W.T
        a -= np.max(a) # normalize to avoid overflow
        sm = np.exp(a[y[i]]) / np.sum(np.exp(a))
        loss -= np.log(sm)
        for j in range(C):
            grad[j] += np.exp(a[j]) / np.sum(np.exp(a)) * X[i]
            if y[i] == j:
                grad[j] -= X[i]

    loss /= N
    grad /= N
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    return loss, grad
```

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

    scores = X @ self.W.T
    a_y = scores[range(N), y]
    log_sum = np.log(np.sum(np.exp(scores), axis=1))

    loss = np.sum(log_sum - a_y)

    sm = np.exp(scores) / np.matrix(np.sum(np.exp(scores), axis=1)).T
    sm[range(N), y] -= 1
    grad = sm.T @ X

    loss /= N
    grad /= N

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

    return loss, grad
```

```

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

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

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

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

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

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

        # ===== #
        # YOUR CODE HERE:
        #   Sample batch_size elements from the training data for use in
        #   gradient descent. After sampling,
        #   - X_batch should have shape: (batch_size, dim)
        #   - y_batch should have shape: (batch_size,)
        #   The indices should be randomly generated to reduce correlations
        #   in the dataset. Use np.random.choice. It's okay to sample with
        #   replacement.
        # ===== #
        rand = np.random.choice(num_train, batch_size)
        X_batch, y_batch = X[rand], y[rand]
        # ===== #
        # END YOUR CODE HERE
        # ===== #
        # evaluate loss and gradient
        loss, grad = self.fast_loss_and_grad(X_batch, y_batch)
        loss_history.append(loss)

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

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

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

    return loss_history

```

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

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

    return y_pred
```