

CS280 Hw1

1. First we have three formulae

$$(1) \ell(B) = \sum_{n=1}^N \log P(X_n|B)$$

$$P(X_n|B) = \prod_{k=1}^K \pi_k \cdot r(X_n|M_k, \bar{\Sigma}_k)$$

$$\mathcal{N}(X_n|M_k, \bar{\Sigma}_k) = \frac{1}{\sqrt{(2\pi)^d |\bar{\Sigma}_k|}} \cdot \exp^{-\frac{1}{2}(X_n - M_k)^T \bar{\Sigma}_k^{-1} (X_n - M_k)}$$

and then we can derive the likelihood w.r.t π_k .

$$\frac{\partial \ell(B)}{\partial \pi_k} = \sum_{n=1}^N \frac{1}{P(X_n|B)} \cdot \frac{\partial P(X_n|B)}{\partial \pi_k}$$

$$= \sum_{n=1}^N \left[\frac{1}{\prod_{k=1}^K \pi_k \cdot \mathcal{N}(X_n|M_k, \bar{\Sigma}_k)} \cdot \pi_k \cdot \mathcal{N}(X_n|M_k, \bar{\Sigma}_k) \cdot (X_n - M_k) \cdot \bar{\Sigma}_k^{-1} \right]$$

$$= \sum_{n=1}^N r_{nk} \cdot \bar{\Sigma}_k^{-1} \cdot (X_n - M_k)$$

(2) the process is same by cs (1) to derive π_k

$$\text{according to chain rule}$$

$$\text{without constraints}$$

$$= \sum_{n=1}^N \frac{1}{P(X_n|B)} \cdot \mathcal{N}(X_n|M_k, \bar{\Sigma}_k)$$

$$= \sum_{n=1}^N r_{nk}$$

\therefore finally, we get it

according to Lagrangian, we can have

$$\text{with constraints}$$

$$L(\pi_k, \lambda) = \ell(B) + \lambda \left(1 - \sum_{k=1}^K \pi_k \right)$$

$$\frac{\partial L}{\partial \pi_k} = \sum_{n=1}^N r_{nk} - \lambda = 0$$

$$\therefore \lambda = \frac{1}{K} \sum_{k=1}^K \sum_{n=1}^N r_{nk} = \frac{N}{K}$$

$$\therefore \frac{\partial \ell(B)}{\partial \pi_k} = \sum_{n=1}^N r_{nk} - \frac{N}{K}$$

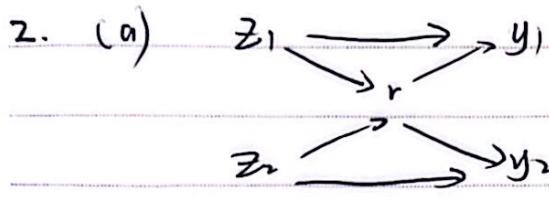
\therefore we get it.

REPLACEABLE



CS 扫描全能王

3亿人都在用的扫描App



on your definition,
 $\frac{\partial E}{\partial y_1}$ and $\frac{\partial E}{\partial y_2}$, depending
 and assuming we already have

(b) when given \bar{y}_i , we can compute error easily denoting it as $\bar{E}(y_i, \bar{y}_i)$

$$\frac{\partial \bar{E}}{\partial r} = \frac{\partial \bar{E}}{\partial y_1} \frac{\partial y_1}{\partial r} + \frac{\partial \bar{E}}{\partial y_2} \frac{\partial y_2}{\partial r}, \text{ for } \frac{\partial y_i}{\partial r} = -\frac{e^{z_i}}{r^2}$$

$$= \frac{\partial \bar{E}}{\partial y_1} \cdot \frac{-e^{z_1}}{r^2} + \frac{\partial \bar{E}}{\partial y_2} \cdot \frac{-e^{z_2}}{r^2}$$

$$\begin{aligned} \frac{\partial \bar{E}}{\partial z_1} &= \frac{\partial \bar{E}}{\partial y_1} \frac{\partial y_1}{\partial z_1} + \frac{\partial \bar{E}}{\partial r} \frac{\partial r}{\partial z_1} \\ &= \frac{\partial \bar{E}}{\partial y_1} \frac{e^{z_1}}{r} + \frac{\partial \bar{E}}{\partial r} \cdot e^{z_1} \end{aligned}$$

$$\begin{aligned} \frac{\partial \bar{E}}{\partial z_2} &= \frac{\partial \bar{E}}{\partial y_2} \frac{\partial y_2}{\partial z_2} + \frac{\partial \bar{E}}{\partial r} \frac{\partial r}{\partial z_2} \\ &= \frac{\partial \bar{E}}{\partial y_2} \frac{e^{z_2}}{r} + \frac{\partial \bar{E}}{\partial r} \cdot e^{z_2} \end{aligned}$$

and then we can use $\frac{\partial \bar{E}}{\partial z_1}, \frac{\partial \bar{E}}{\partial z_2}$ to update \bar{z}_j , if $\bar{E}(y_i, \bar{y}_i)$ is determined.

Also, for more general case, this formula is also useful if $\frac{\partial \bar{E}}{\partial y_j}$ is given

$$\begin{cases} \frac{\partial \bar{E}}{\partial z_j} = \frac{\partial \bar{E}}{\partial y_j} \frac{e^{z_j}}{r} + \frac{\partial \bar{E}}{\partial r} \cdot e^{z_j} \\ \frac{\partial \bar{E}}{\partial r} = \frac{\partial \bar{E}}{\partial y_j} \frac{-e^{z_j}}{r^2} \end{cases}$$

when given $\Rightarrow \frac{\partial \bar{E}}{\partial z_j} = \bar{z}_j = \bar{y}_j \frac{e^{z_j}}{r} + \bar{r} e^{z_j}$

$$\frac{\partial \bar{E}}{\partial r} = \bar{r} = -\bar{y}_j \frac{e^{z_j}}{r^2}$$

(c) def softmax_VJP(Z, Y-bar):

$$R = np.sum(np.exp(Z), axis=1, keepdims=True)$$

$$R-bar = -np.sum(I-bar * np.exp(Z), axis=1, keepdims=True)$$

$$1/R * x2 = 1 - R-bar / R$$

$$Z-bar = Y-bar * (np.exp(Z)/R) + R-bar * np.exp(Z)$$

return Z-bar



Perceptron Learning Algorithm

The perceptron is a simple supervised machine learning algorithm and one of the earliest neural network architectures. It was introduced by Rosenblatt in the late 1950s. A perceptron represents a binary linear classifier that maps a set of training examples (of d dimensional input vectors) onto binary output values using a $d - 1$ dimensional hyperplane. But Today, we will implement **Multi-Classes**

Perceptron Learning Algorithm Given:

- dataset $\{(x^i, y^i)\}, i \in (1, M)$
- x^i is d dimension vector, $x^i = (x_1^i, \dots, x_d^i)$
- y^i is multi-class target variable $y^i \in \{0, 1, 2\}$

A perceptron is trained using gradient descent. The training algorithm has different steps. In the beginning (step 0) the model parameters are initialized. The other steps (see below) are repeated for a specified number of training iterations or until the parameters have converged.

Step0: Initial the weight vector and bias with zeros

Step1: Compute the linear combination of the input features and weight.

$$y_{pred}^i = \arg \max_k W_k * x^i + b$$

Step2: Compute the gradients for parameters W_k, b . **Derive the parameter update equation Here (5 points)**

#####

TODO: Derive you answer hear

Firstly, our derivation of Eq. gives us

$$\frac{\partial L}{\partial w} = - \sum_{x^i} y^i x^i$$

$$\frac{\partial L}{\partial b} = - \sum_{x^i} y^i$$

And then we use it to update our parameters by multiplying with learning rate η

$$W_{y^i} = W_{y^i} + \eta y^i x^i$$

$$b_{y^i} = b_{y^i} + \eta y^i$$

Therefore, if the datapoint is classified correctly, the W_{y^i} and b_{y^i} will become larger. Otherwise, it will become smaller. Finally, we can get a trained perceptron model.

#####

In [1]:

```
from sklearn import datasets
import numpy as np
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
import random

np.random.seed(0)
random.seed(0)
```

```
In [2]: iris = datasets.load_iris()
X = iris.data
print(type(X))
y = iris.target
y = np.array(y)
print('X_Shape:', X.shape)
print('y_Shape:', y.shape)
print('Label Space:', np.unique(y))
```

```
<class 'numpy.ndarray'>
X_Shape: (150, 4)
y_Shape: (150,)
Label Space: [0 1 2]
```

```
In [3]: ## split the training set and test set
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.3, random_state=0)
print('X_train_Shape:', X_train.shape)
print('X_test_Shape:', X_test.shape)
print('y_train_Shape:', y_train.shape)
print('y_test_Shape:', y_test.shape)

print(type(y_train))
```

```
X_train_Shape: (105, 4)
X_test_Shape: (45, 4)
y_train_Shape: (105,)
y_test_Shape: (105,)
<class 'numpy.ndarray'>
```

```
In [19]: class MultiClsPLA(object):

    ## We recommend to absorb the bias into weight. w = [w, b]

    def __init__(self, X_train, y_train, X_test, y_test, lr, num_epoch, weight_d):
        super(MultiClsPLA, self).__init__()
        self.X_train = X_train
        self.y_train = y_train
        self.X_test = X_test
        self.y_test = y_test
        self.weight = self.initial_weight(weight_dimension, num_cls)
        self.sample_mean = np.mean(self.X_train, 0)
        self.sample_std = np.std(self.X_train, 0)
        self.num_epoch = num_epoch
        self.lr = lr
        self.total_acc_train = []
        self.total_acc_tst = []

    def initial_weight(self, weight_dimension, num_cls):
        #####
        ## Initialize the weight with      ##
        ## small std and zero mean gaussian  ##
        #####
        weight = np.random.normal(0, 0.01, (weight_dimension + 1, num_cls))
        return weight

    def data_preprocessing(self, data):
        #####
        ## Normalize the data      ##
        #####
        norm_data = (data - self.sample_mean) / self.sample_std
```

```

    return norm_data

def train_step(self, X_train, y_train, shuffle_idx):
    np.random.shuffle(shuffle_idx)
    X_train = X_train[shuffle_idx]
    y_train = y_train[shuffle_idx]

    # Append a column of ones to the end of X_train for the bias term
    X_train_bias = np.hstack((X_train, np.ones((X_train.shape[0], 1)))))

    for i in range(X_train.shape[0]):
        xi = X_train_bias[i]
        yi = y_train[i]
        yi_pred = np.argmax(np.dot(xi, self.weight))
        if yi_pred != yi:
            self.weight[:, yi] += self.lr * xi
            self.weight[:, yi_pred] -= self.lr * xi

    # Calculate training accuracy
    pred_train = np.argmax(np.dot(X_train_bias, self.weight), axis=1)
    train_acc = np.mean(pred_train == y_train)
    return train_acc

def test_step(self, X_test, y_test):
    X_test = self.data_preprocessing(data=X_test)

    # Append a column of ones to the end of X_test for the bias term
    X_test_bias = np.hstack((X_test, np.ones((X_test.shape[0], 1)))))

    pred_test = np.argmax(np.dot(X_test_bias, self.weight), axis=1)
    test_acc = np.mean(pred_test == y_test)

    return test_acc

def train(self):
    self.X_train = self.data_preprocessing(data=self.X_train)
    num_sample = self.X_train.shape[0]

    ##### In order to absorb the bias into weights #####
    #### we need to modify the input data. #####
    #####

    shuffle_index = np.array(range(0, num_sample))
    for epoch in range(self.num_epoch):
        training_acc = self.train_step(X_train=self.X_train, y_train=self.y_
        tst_acc = self.test_step(X_test=self.X_test, y_test=self.y_test)
        self.total_acc_train.append(training_acc)
        self.total_acc_tst.append(tst_acc)
        print('epoch:', epoch, 'training_acc:%.3f' % training_acc, 'test_acc'

def vis_acc_curve(self):
    train_acc = np.array(self.total_acc_train)
    test_acc = np.array(self.total_acc_tst)
    plt.plot(train_acc)
    plt.plot(test_acc)
    plt.legend(['train_acc', 'test_acc'])
    plt.xlabel('Epoch')
    plt.ylabel('Accuracy')

```

```
plt.title('Training and Test Accuracy over Epochs')
plt.show()
```

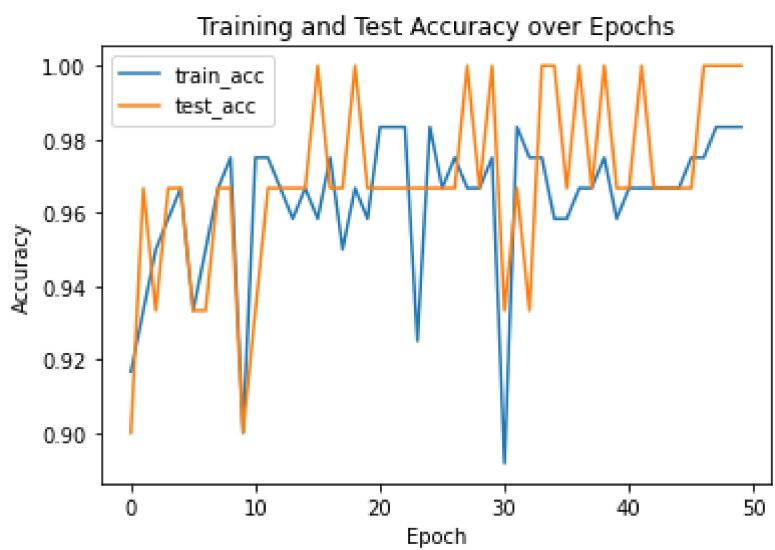
```
In [22]: np.random.seed(0)
random.seed(0)
#####
### TODO:
### 1. You need to import the model and pass some parameters.
### 2. Then training the model with some epoches.
### 3. Visualize the training acc and test acc verus epoches

# Parameters
lr = 0.001          # Learning rate
num_epoch = 50       # number of epochs
weight_dimension = X_train.shape[1] # number of features
num_cls = len(np.unique(y))         # number of classes

# Initialize and train the MultiClsPLA model
model = MultiClsPLA(X_train, y_train, X_test, y_test, lr, num_epoch, weight_dime
model.train()

# Visualize the training accuracy and test accuracy versus epochs
model.vis_acc_curve()
```

```
epoch: 0 training_acc:0.917 test_acc:0.900
epoch: 1 training_acc:0.933 test_acc:0.967
epoch: 2 training_acc:0.950 test_acc:0.933
epoch: 3 training_acc:0.958 test_acc:0.967
epoch: 4 training_acc:0.967 test_acc:0.967
epoch: 5 training_acc:0.933 test_acc:0.933
epoch: 6 training_acc:0.950 test_acc:0.933
epoch: 7 training_acc:0.967 test_acc:0.967
epoch: 8 training_acc:0.975 test_acc:0.967
epoch: 9 training_acc:0.900 test_acc:0.900
epoch: 10 training_acc:0.975 test_acc:0.933
epoch: 11 training_acc:0.975 test_acc:0.967
epoch: 12 training_acc:0.967 test_acc:0.967
epoch: 13 training_acc:0.958 test_acc:0.967
epoch: 14 training_acc:0.967 test_acc:0.967
epoch: 15 training_acc:0.958 test_acc:1.000
epoch: 16 training_acc:0.975 test_acc:0.967
epoch: 17 training_acc:0.950 test_acc:0.967
epoch: 18 training_acc:0.967 test_acc:1.000
epoch: 19 training_acc:0.958 test_acc:0.967
epoch: 20 training_acc:0.983 test_acc:0.967
epoch: 21 training_acc:0.983 test_acc:0.967
epoch: 22 training_acc:0.983 test_acc:0.967
epoch: 23 training_acc:0.925 test_acc:0.967
epoch: 24 training_acc:0.983 test_acc:0.967
epoch: 25 training_acc:0.967 test_acc:0.967
epoch: 26 training_acc:0.975 test_acc:0.967
epoch: 27 training_acc:0.967 test_acc:1.000
epoch: 28 training_acc:0.967 test_acc:0.967
epoch: 29 training_acc:0.975 test_acc:1.000
epoch: 30 training_acc:0.892 test_acc:0.933
epoch: 31 training_acc:0.983 test_acc:0.967
epoch: 32 training_acc:0.975 test_acc:0.933
epoch: 33 training_acc:0.975 test_acc:1.000
epoch: 34 training_acc:0.958 test_acc:1.000
epoch: 35 training_acc:0.958 test_acc:0.967
epoch: 36 training_acc:0.967 test_acc:1.000
epoch: 37 training_acc:0.967 test_acc:0.967
epoch: 38 training_acc:0.975 test_acc:1.000
epoch: 39 training_acc:0.958 test_acc:0.967
epoch: 40 training_acc:0.967 test_acc:0.967
epoch: 41 training_acc:0.967 test_acc:1.000
epoch: 42 training_acc:0.967 test_acc:0.967
epoch: 43 training_acc:0.967 test_acc:0.967
epoch: 44 training_acc:0.967 test_acc:0.967
epoch: 45 training_acc:0.975 test_acc:0.967
epoch: 46 training_acc:0.975 test_acc:1.000
epoch: 47 training_acc:0.983 test_acc:1.000
epoch: 48 training_acc:0.983 test_acc:1.000
epoch: 49 training_acc:0.983 test_acc:1.000
```



k-Nearest Neighbor (kNN) exercise

Complete and hand in this completed worksheet (including its outputs and any supporting code outside of the worksheet) with your assignment submission. For more details see the [assignments page](#) on the course website.

The kNN classifier consists of two stages:

- During training, the classifier takes the training data and simply remembers it
- During testing, kNN classifies every test image by comparing to all training images and transferring the labels of the k most similar training examples
- The value of k is cross-validated

In this exercise you will implement these steps and understand the basic Image Classification pipeline, cross-validation, and gain proficiency in writing efficient, vectorized code.

In [1]: *# Run some setup code for this notebook.*

```
import random
import numpy as np
from cs231n.data_utils import load_CIFAR10
import matplotlib.pyplot as plt

# This is a bit of magic to make matplotlib figures appear inline in the notebook
# rather than in a new window.
%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# Some more magic so that the notebook will reload external python modules;
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2
```

In [2]: *# Load the raw CIFAR-10 data.*

```
cifar10_dir = 'cs231n/datasets/cifar-10-batches-py'

# Cleaning up variables to prevent loading data multiple times (which may cause
try:
    del X_train, y_train
    del X_test, y_test
    print('Clear previously loaded data.')
except:
    pass

X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

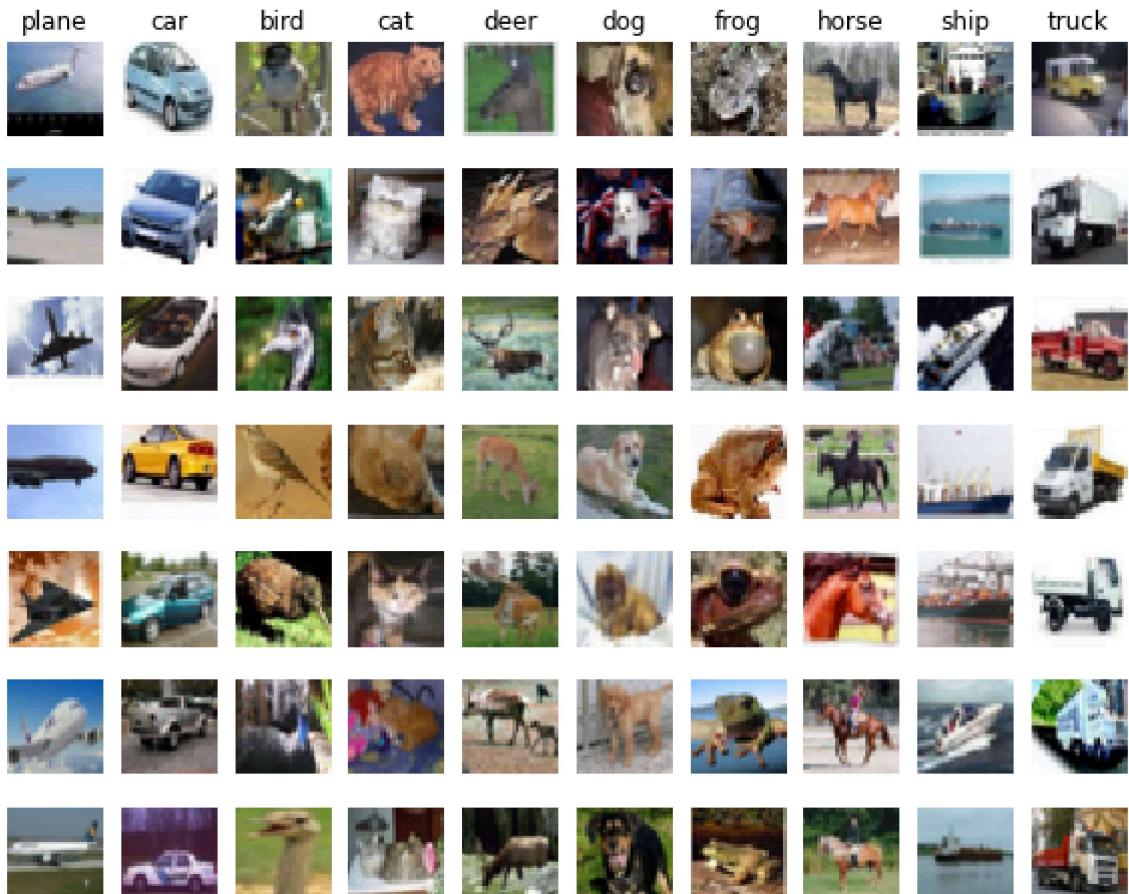
# As a sanity check, we print out the size of the training and test data.
print('Training data shape: ', X_train.shape)
print('Training labels shape: ', y_train.shape)
```

```
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']
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)
```

In [6]:

```
from cs231n.classifiers import KNearestNeighbor

# Create a kNN classifier instance.
# Remember that training a kNN classifier is a noop:
# the Classifier simply remembers the data and does no further processing
classifier = KNearestNeighbor()
classifier.train(X_train, y_train)
```

We would now like to classify the test data with the kNN classifier. Recall that we can break down this process into two steps:

1. First we must compute the distances between all test examples and all train examples.
2. Given these distances, for each test example we find the k nearest examples and have them vote for the label

Lets begin with computing the distance matrix between all training and test examples. For example, if there are **Ntr** training examples and **Nte** test examples, this stage should result in a **Nte x Ntr** matrix where each element (i,j) is the distance between the i-th test and j-th train example.

Note: For the three distance computations that we require you to implement in this notebook, you may not use the np.linalg.norm() function that numpy provides.

First, open `cs231n/classifiers/k_nearest_neighbor.py` and implement the function `compute_distances_two_loops` that uses a (very inefficient) double loop over all pairs of (test, train) examples and computes the distance matrix one element at a time.

In [7]:

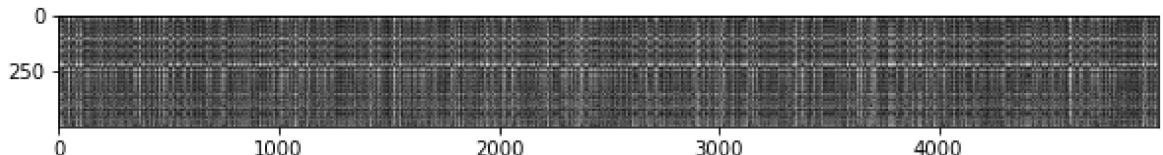
```
# Open cs231n/classifiers/k_nearest_neighbor.py and implement
# compute_distances_two_loops.

# Test your implementation:
dists = classifier.compute_distances_two_loops(X_test)
print(dists.shape)

(500, 5000)
```

In [8]:

```
# We can visualize the distance matrix: each row is a single test example and
# its distances to training examples
plt.imshow(dists, interpolation='none')
plt.show()
```



Inline Question 1

Notice the structured patterns in the distance matrix, where some rows or columns are visible brighter. (Note that with the default color scheme black indicates low distances while white indicates high distances.)

- What in the data is the cause behind the distinctly bright rows?
- What causes the columns?

Your Answer :

1. Because white scheme color indicates high distances, therefore, the distinctly bright rows mean that the test data is distinct to almost all data in the training set.
2. Because white scheme color indicates high distances, therefore, the distinctly bright columns mean that the training data is distinct to almost all data in the test set.

```
In [9]: # Now implement the function predict_labels and run the code below:  
# We use k = 1 (which is Nearest Neighbor).  
y_test_pred = classifier.predict_labels(dists, k=1)  
  
# Compute and print the fraction of correctly predicted examples  
num_correct = np.sum(y_test_pred == y_test)  
accuracy = float(num_correct) / num_test  
print('Got %d / %d correct => accuracy: %f' % (num_correct, num_test, accuracy))
```

Got 137 / 500 correct => accuracy: 0.274000

You should expect to see approximately 27% accuracy. Now lets try out a larger k , say $k = 5$:

```
In [10]: y_test_pred = classifier.predict_labels(dists, k=5)  
num_correct = np.sum(y_test_pred == y_test)  
accuracy = float(num_correct) / num_test  
print('Got %d / %d correct => accuracy: %f' % (num_correct, num_test, accuracy))
```

Got 139 / 500 correct => accuracy: 0.278000

You should expect to see a slightly better performance than with $k = 1$.

Inline Question 2

We can also use other distance metrics such as L1 distance. For pixel values $p_{ij}^{(k)}$ at location (i, j) of some image I_k ,

the mean μ across all pixels over all images is

$$\mu = \frac{1}{nhw} \sum_{k=1}^n \sum_{i=1}^h \sum_{j=1}^w p_{ij}^{(k)}$$

And the pixel-wise mean μ_{ij} across all images is

$$\mu_{ij} = \frac{1}{n} \sum_{k=1}^n p_{ij}^{(k)}.$$

The general standard deviation σ and pixel-wise standard deviation σ_{ij} is defined similarly.

Which of the following preprocessing steps will not change the performance of a Nearest Neighbor classifier that uses L1 distance? Select all that apply.

1. Subtracting the mean μ ($\tilde{p}_{ij}^{(k)} = p_{ij}^{(k)} - \mu$)
2. Subtracting the per pixel mean μ_{ij} ($\tilde{p}_{ij}^{(k)} = p_{ij}^{(k)} - \mu_{ij}$)
3. Subtracting the mean μ and dividing by the standard deviation σ .
4. Subtracting the pixel-wise mean μ_{ij} and dividing by the pixel-wise standard deviation σ_{ij} .
5. Rotating the coordinate axes of the data.

Your Answer :

1, 2, 5

Your Explanation :

1. Because L1 means the abs distance, that is $\sum \|x_1 - x_2\|$, then subtracting one constant will not have a influence on performance.
2. Same as 1.
3. Because it will divide std, then it will influence the distance by changing scaling.
4. Same as 3.
5. Just rotating the axis will not have a influence on the L1 distance, because it doesnot change the scaling.

```
In [11]: # Now Lets speed up distance matrix computation by using partial vectorization
# with one loop. Implement the function compute_distances_one_loop and run the
# code below:
dists_one = classifier.compute_distances_one_loop(X_test)

# To ensure that our vectorized implementation is correct, we make sure that it
# agrees with the naive implementation. There are many ways to decide whether
# two matrices are similar; one of the simplest is the Frobenius norm. In case
# you haven't seen it before, the Frobenius norm of two matrices is the square
# root of the squared sum of differences of all elements; in other words, reshape
# the matrices into vectors and compute the Euclidean distance between them.
difference = np.linalg.norm(dists - dists_one, ord='fro')
print('One loop difference was: %f' % (difference, ))
if difference < 0.001:
    print('Good! The distance matrices are the same')
```

```
else:
    print('Uh-oh! The distance matrices are different')
```

One loop difference was: 0.000000
Good! The distance matrices are the same

```
In [13]: # Now implement the fully vectorized version inside compute_distances_no_loops
# and run the code
dists_two = classifier.compute_distances_no_loops(X_test)

# check that the distance matrix agrees with the one we computed before:
difference = np.linalg.norm(dists - dists_two, ord='fro')
print('No loop difference was: %f' % (difference, ))
if difference < 0.001:
    print('Good! The distance matrices are the same')
else:
    print('Uh-oh! The distance matrices are different')
```

No loop difference was: 0.000000
Good! The distance matrices are the same

```
In [14]: # Let's compare how fast the implementations are
def time_function(f, *args):
    """
    Call a function f with args and return the time (in seconds) that it took to
    """
    import time
    tic = time.time()
    f(*args)
    toc = time.time()
    return toc - tic

two_loop_time = time_function(classifier.compute_distances_two_loops, X_test)
print('Two loop version took %f seconds' % two_loop_time)

one_loop_time = time_function(classifier.compute_distances_one_loop, X_test)
print('One loop version took %f seconds' % one_loop_time)

no_loop_time = time_function(classifier.compute_distances_no_loops, X_test)
print('No loop version took %f seconds' % no_loop_time)

# You should see significantly faster performance with the fully vectorized impl

# NOTE: depending on what machine you're using,
# you might not see a speedup when you go from two loops to one loop,
# and might even see a slow-down.
```

Two loop version took 41.052701 seconds
One loop version took 68.560830 seconds
No loop version took 0.177495 seconds

Cross-validation

We have implemented the k-Nearest Neighbor classifier but we set the value $k = 5$ arbitrarily. We will now determine the best value of this hyperparameter with cross-validation.

```
In [15]: num_folds = 5
k_choices = [1, 3, 5, 8, 10, 12, 15, 20, 50, 100]
```

```

X_train_folds = []
y_train_folds = []
#####
# TODO:
# Split up the training data into folds. After splitting, X_train_folds and
# y_train_folds should each be lists of length num_folds, where
# y_train_folds[i] is the label vector for the points in X_train_folds[i].
# Hint: Look up the numpy array_split function.
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

X_train_folds = np.array_split(X_train,num_folds)
y_train_folds = np.array_split(y_train,num_folds)

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****


# A dictionary holding the accuracies for different values of k that we find
# when running cross-validation. After running cross-validation,
# k_to_accuracies[k] should be a list of length num_folds giving the different
# accuracy values that we found when using that value of k.
k_to_accuracies = {}

#####
# TODO:
# Perform k-fold cross validation to find the best value of k. For each
# possible value of k, run the k-nearest-neighbor algorithm num_folds times,
# where in each case you use all but one of the folds as training data and the
# last fold as a validation set. Store the accuracies for all fold and all
# values of k in the k_to_accuracies dictionary.
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****


for k in k_choices:
    accuracy_set = []
    for i in range(num_folds):
        #cross validation
        X_train_fold = np.concatenate([fold_content for index, fold_content in enumerate(X_train_folds) if index != i])
        y_train_fold = np.concatenate([fold_content for index, fold_content in enumerate(y_train_folds) if index != i])
        X_test_fold = X_train_folds[i]
        y_test_fold = y_train_folds[i]

        #train model
        model = KNearestNeighbor()
        model.train(X_train_fold, y_train_fold)

        #predict
        y_test_pred = model.predict(X_test_fold, k)

        #metric
        acc_rate = np.mean(y_test_pred == y_test_fold)
        accuracy_set.append(acc_rate)
    k_to_accuracies[k] = accuracy_set

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****


# Print out the computed accuracies
for k in sorted(k_to_accuracies):

```

```

for accuracy in k_to_accuracies[k]:
    print('k = %d, accuracy = %f' % (k, accuracy))

k = 1, accuracy = 0.263000
k = 1, accuracy = 0.257000
k = 1, accuracy = 0.264000
k = 1, accuracy = 0.278000
k = 1, accuracy = 0.266000
k = 3, accuracy = 0.239000
k = 3, accuracy = 0.249000
k = 3, accuracy = 0.240000
k = 3, accuracy = 0.266000
k = 3, accuracy = 0.254000
k = 5, accuracy = 0.248000
k = 5, accuracy = 0.266000
k = 5, accuracy = 0.280000
k = 5, accuracy = 0.292000
k = 5, accuracy = 0.280000
k = 8, accuracy = 0.262000
k = 8, accuracy = 0.282000
k = 8, accuracy = 0.273000
k = 8, accuracy = 0.290000
k = 8, accuracy = 0.273000
k = 10, accuracy = 0.265000
k = 10, accuracy = 0.296000
k = 10, accuracy = 0.276000
k = 10, accuracy = 0.284000
k = 10, accuracy = 0.280000
k = 12, accuracy = 0.260000
k = 12, accuracy = 0.295000
k = 12, accuracy = 0.279000
k = 12, accuracy = 0.283000
k = 12, accuracy = 0.280000
k = 15, accuracy = 0.252000
k = 15, accuracy = 0.289000
k = 15, accuracy = 0.278000
k = 15, accuracy = 0.282000
k = 15, accuracy = 0.274000
k = 20, accuracy = 0.270000
k = 20, accuracy = 0.279000
k = 20, accuracy = 0.279000
k = 20, accuracy = 0.282000
k = 20, accuracy = 0.285000
k = 50, accuracy = 0.271000
k = 50, accuracy = 0.288000
k = 50, accuracy = 0.278000
k = 50, accuracy = 0.269000
k = 50, accuracy = 0.266000
k = 100, accuracy = 0.256000
k = 100, accuracy = 0.270000
k = 100, accuracy = 0.263000
k = 100, accuracy = 0.256000
k = 100, accuracy = 0.263000

```

```

In [16]: # plot the raw observations
for k in k_choices:
    accuracies = k_to_accuracies[k]
    plt.scatter([k] * len(accuracies), accuracies)

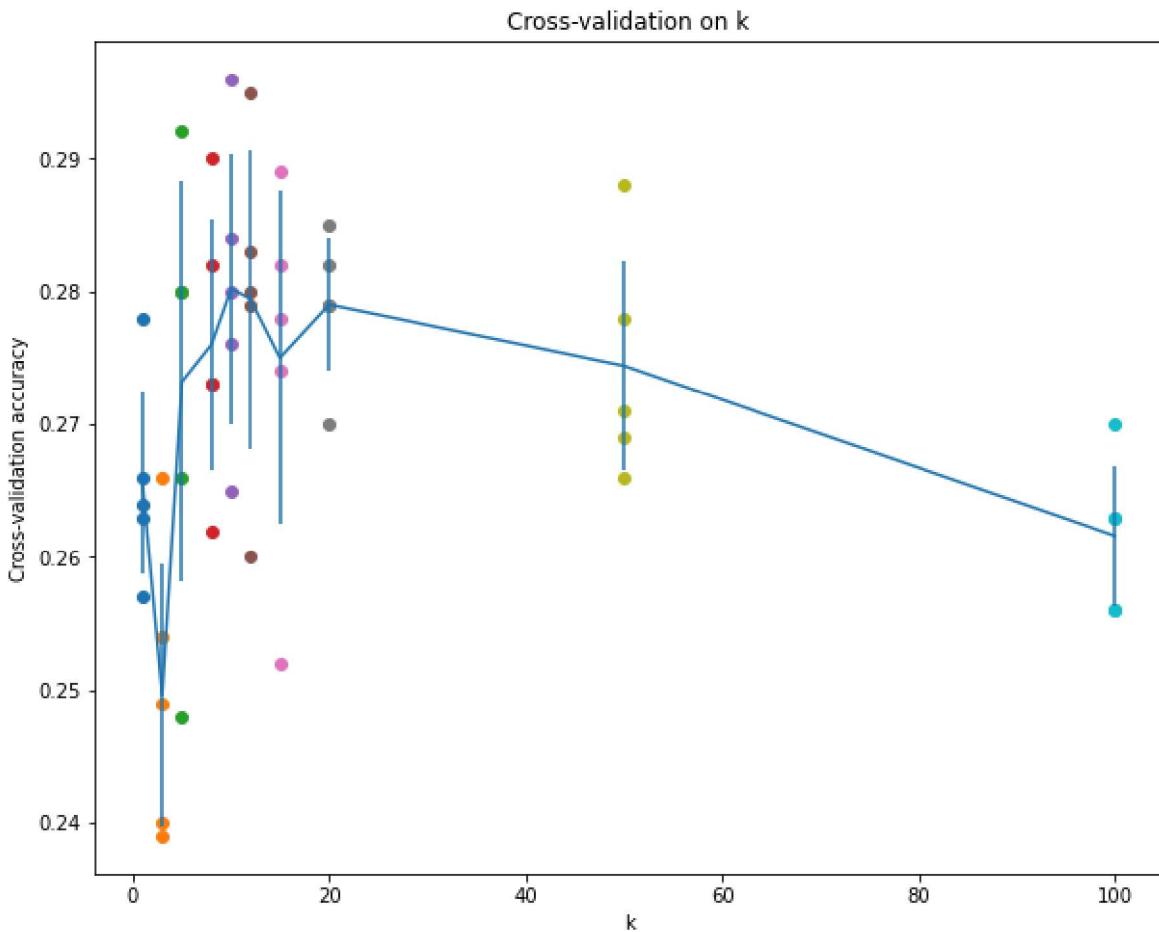
# plot the trend line with error bars that correspond to standard deviation

```

```

accuracies_mean = np.array([np.mean(v) for k,v in sorted(k_to_accuracies.items())])
accuracies_std = np.array([np.std(v) for k,v in sorted(k_to_accuracies.items())])
plt.errorbar(k_choices, accuracies_mean, yerr=accuracies_std)
plt.title('Cross-validation on k')
plt.xlabel('k')
plt.ylabel('Cross-validation accuracy')
plt.show()

```



In [20]:

```

# Based on the cross-validation results above, choose the best value for k,
# retrain the classifier using all the training data, and test it on the test
# data. You should be able to get above 28% accuracy on the test data.
best_k = 10

classifier = KNearestNeighbor()
classifier.train(X_train, y_train)
y_test_pred = classifier.predict(X_test, k=best_k)

# Compute and display the accuracy
num_correct = np.sum(y_test_pred == y_test)
accuracy = float(num_correct) / num_test
print('Got %d / %d correct => accuracy: %f' % (num_correct, num_test, accuracy))

```

Got 141 / 500 correct => accuracy: 0.282000

Inline Question 3

Which of the following statements about k -Nearest Neighbor (k -NN) are true in a classification setting, and for all k ? Select all that apply.

1. The decision boundary of the k -NN classifier is linear.
2. The training error of a 1-NN will always be lower than that of 5-NN.

3. The test error of a 1-NN will always be lower than that of a 5-NN.
4. The time needed to classify a test example with the k-NN classifier grows with the size of the training set.
5. None of the above.

Your Answer :

2,4

Your Explanation :

1. False. Because we use L_2 distance, therefore, the decision boundary of the K-NN should be nonlinear.
2. True. Because when $k = 1$, during the training process, it will choose its self label as its result, and then the training error will decrease too much. But for test process, it will not be a brilliant k value choice.
3. False. The reason is same as 2, that is it is not a general model without robustness.
4. True. Because when set becomes bigger, we need more time to calculate the distance matrix.
5. False.

Multiclass Support Vector Machine exercise

Complete and hand in this completed worksheet (including its outputs and any supporting code outside of the worksheet) with your assignment submission. For more details see the [assignments page](#) on the course website.

In this exercise you will:

- implement a fully-vectorized **loss function** for the SVM
- implement the fully-vectorized expression for its **analytic gradient**
- **check your implementation** using numerical gradient
- use a validation set to **tune the learning rate and regularization strength**
- **optimize** the loss function with **SGD**
- **visualize** the final learned weights

```
In [1]: # Run some setup code for this notebook.
import random
import numpy as np
from cs231n.data_utils import load_CIFAR10
import matplotlib.pyplot as plt

# This is a bit of magic to make matplotlib figures appear inline in the
# notebook rather than in a new window.
%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# Some more magic so that the notebook will reload external python modules;
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2
```

CIFAR-10 Data Loading and Preprocessing

```
In [2]: # Load the raw CIFAR-10 data.
cifar10_dir = 'cs231n/datasets/cifar-10-batches-py'

# Cleaning up variables to prevent Loading data multiple times (which may cause
try:
    del X_train, y_train
    del X_test, y_test
    print('Clear previously loaded data.')
except:
    pass

X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

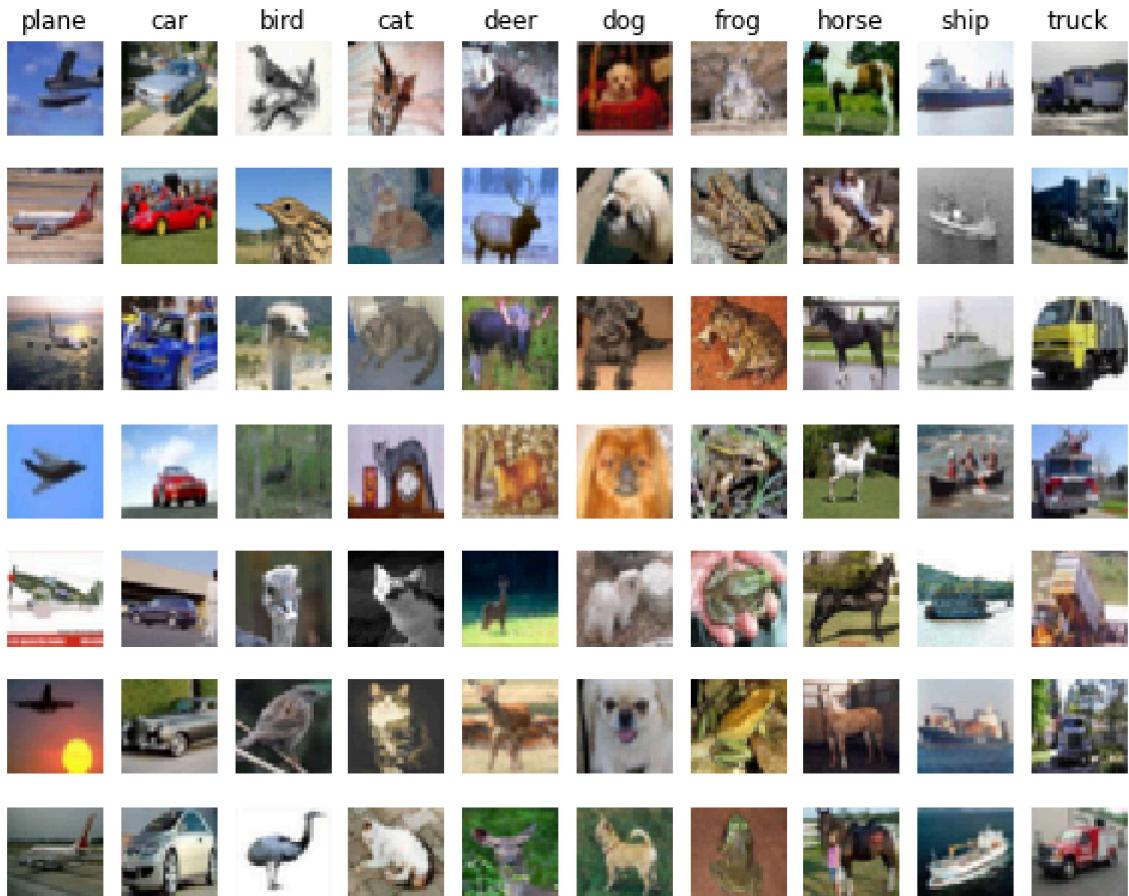
# As a sanity check, we print out the size of the training and test data.
print('Training data shape: ', X_train.shape)
```

```
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']
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]:

```
# 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)

```

```

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,)

```

In [5]:

```

# 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 [6]:

```

# 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()

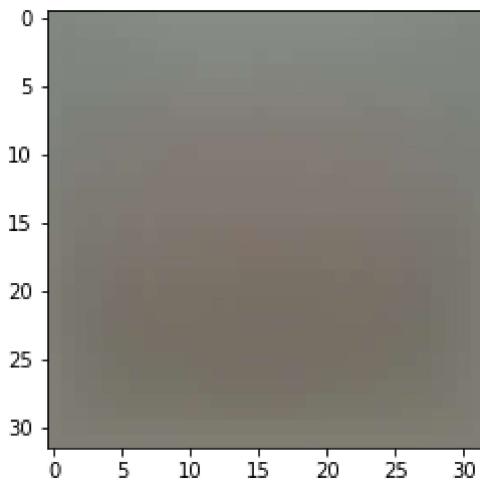
# second: subtract the mean image from train and test data
X_train -= mean_image
X_val -= mean_image
X_test -= mean_image
X_dev -= mean_image

# 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)

```

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



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

SVM Classifier

Your code for this section will all be written inside
`cs231n/classifiers/linear_svm.py`.

As you can see, we have prefilled the function `svm_loss_naive` which uses for loops to evaluate the multiclass SVM loss function.

In [7]:

```

# Evaluate the naive implementation of the Loss we provided for you:
from cs231n.classifiers.linear_svm import svm_loss_naive
import time

# generate a random SVM weight matrix of small numbers
W = np.random.randn(3073, 10) * 0.0001

loss, grad = svm_loss_naive(W, X_dev, y_dev, 0.000005)
print('loss: %f' % (loss, ))

```

loss: 8.971985

The `grad` returned from the function above is right now all zero. Derive and implement the gradient for the SVM cost function and implement it inline inside the function `svm_loss_naive`. You will find it helpful to interleave your new code inside the existing function.

To check that you have correctly implemented the gradient correctly, you can numerically estimate the gradient of the loss function and compare the numeric estimate to the gradient that you computed. We have provided code that does this for you:

```
In [8]: # Once you've implemented the gradient, recompute it with the code below
# and gradient check it with the function we provided for you

# Compute the loss and its gradient at W.
loss, grad = svm_loss_naive(W, X_dev, y_dev, 0.0)

# Numerically compute the gradient along several randomly chosen dimensions, and
# compare them with your analytically computed gradient. The numbers should match
# almost exactly along all dimensions.
from cs231n.gradient_check import grad_check_sparse
f = lambda w: svm_loss_naive(w, X_dev, y_dev, 0.0)[0]
grad_numerical = grad_check_sparse(f, W, grad)

# do the gradient check once again with regularization turned on
# you didn't forget the regularization gradient did you?
loss, grad = svm_loss_naive(W, X_dev, y_dev, 5e1)
f = lambda w: svm_loss_naive(w, X_dev, y_dev, 5e1)[0]
grad_numerical = grad_check_sparse(f, W, grad)

numerical: 7.896355 analytic: 7.896355, relative error: 2.626071e-11
numerical: -27.043911 analytic: -27.043911, relative error: 1.211571e-11
numerical: 8.610863 analytic: 8.610863, relative error: 3.757815e-11
numerical: -6.985576 analytic: -6.985576, relative error: 2.018770e-11
numerical: -23.808947 analytic: -23.808947, relative error: 3.948297e-12
numerical: -24.416677 analytic: -24.416677, relative error: 7.454803e-12
numerical: 22.122328 analytic: 22.122328, relative error: 2.767677e-12
numerical: 14.112074 analytic: 14.112074, relative error: 6.778871e-12
numerical: -11.918589 analytic: -11.918589, relative error: 8.883506e-12
numerical: -18.474127 analytic: -18.474127, relative error: 3.061632e-12
numerical: -4.286388 analytic: -4.286388, relative error: 3.723350e-11
numerical: 15.038756 analytic: 15.038756, relative error: 1.188120e-11
numerical: -0.393470 analytic: -0.393470, relative error: 4.276086e-10
numerical: -3.736506 analytic: -3.736506, relative error: 6.732757e-12
numerical: 38.159890 analytic: 38.159890, relative error: 2.547050e-12
numerical: 3.342397 analytic: 3.342397, relative error: 1.884225e-11
numerical: 3.198887 analytic: 3.198887, relative error: 1.044514e-11
numerical: -26.398113 analytic: -26.398113, relative error: 1.217336e-11
numerical: 5.376527 analytic: 5.376527, relative error: 2.958570e-12
numerical: 7.654437 analytic: 7.654437, relative error: 5.646457e-11
```

Inline Question 1

It is possible that once in a while a dimension in the gradcheck will not match exactly. What could such a discrepancy be caused by? Is it a reason for concern? What is a simple example in one dimension where a gradient check could fail?

How would change the margin affect of the frequency of this happening? Hint:
the SVM loss function is not strictly speaking differentiable

Your Answer :

The discrepancy would be caused when the loss function is not differentiable at some points. For example, the svm hinge loss function

$L = \max(0, 1 - y * (w * x))$ is not differentiable at some points, especially when $1 - y * (w * x) = 0$. And then for some of those points, we can use

numerical gradient formula $\frac{\partial f}{\partial x} \approx \frac{f(x+h) - f(x-h)}{2h}$ to calculate its numerical

gradient. And actually, sometimes, the numerical one would mismatch with the true gradient, that is discrepancy.

```
In [9]: # Next implement the function svm_loss_vectorized; for now only compute the Loss
# we will implement the gradient in a moment.
tic = time.time()
loss_naive, grad_naive = svm_loss_naive(W, X_dev, y_dev, 0.000005)
toc = time.time()
print('Naive loss: %e computed in %fs' % (loss_naive, toc - tic))

from cs231n.classifiers.linear_svm import svm_loss_vectorized
tic = time.time()
loss_vectorized, _ = svm_loss_vectorized(W, X_dev, y_dev, 0.000005)
toc = time.time()
print('Vectorized loss: %e computed in %fs' % (loss_vectorized, toc - tic))

# The losses should match but your vectorized implementation should be much fast
print('difference: %f' % (loss_naive - loss_vectorized))
```

Naive loss: 8.971985e+00 computed in 0.086998s
 Vectorized loss: 8.971985e+00 computed in 0.002002s
 difference: -0.000000

```
In [10]: # Complete the implementation of svm_loss_vectorized, and compute the gradient
# of the Loss function in a vectorized way.

# The naive implementation and the vectorized implementation should match, but
# the vectorized version should still be much faster.
tic = time.time()
_, grad_naive = svm_loss_naive(W, X_dev, y_dev, 0.000005)
toc = time.time()
print('Naive loss and gradient: computed in %fs' % (toc - tic))

tic = time.time()
_, grad_vectorized = svm_loss_vectorized(W, X_dev, y_dev, 0.000005)
toc = time.time()
print('Vectorized loss and gradient: computed in %fs' % (toc - tic))

# The Loss is a single number, so it is easy to compare the values computed
# by the two implementations. The gradient on the other hand is a matrix, so
# we use the Frobenius norm to compare them.
difference = np.linalg.norm(grad_naive - grad_vectorized, ord='fro')
print('difference: %f' % difference)
```

Naive loss and gradient: computed in 0.086991s
 Vectorized loss and gradient: computed in 0.002000s
 difference: 0.000000

Stochastic Gradient Descent

We now have vectorized and efficient expressions for the loss, the gradient and our gradient matches the numerical gradient. We are therefore ready to do SGD to minimize the loss. Your code for this part will be written inside

```
cs231n/classifiers/linear_classifier.py .
```

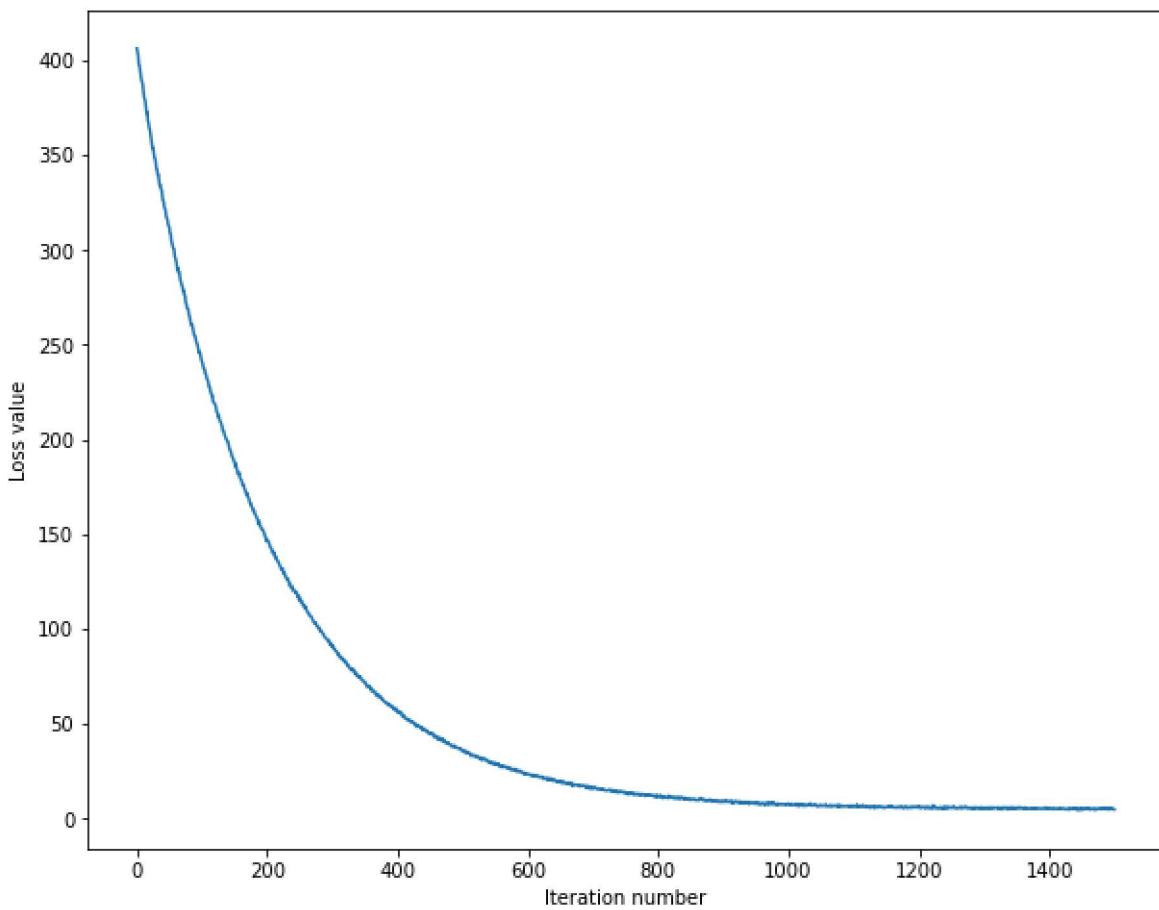
In [11]:

```
# In the file Linear_classifier.py, implement SGD in the function
# LinearClassifier.train() and then run it with the code below.
from cs231n.classifiers import LinearSVM
svm = LinearSVM()
tic = time.time()
loss_hist = svm.train(X_train, y_train, learning_rate=1e-7, reg=2.5e4,
                      num_iters=1500, verbose=True)
toc = time.time()
print('That took %fs' % (toc - tic))
```

```
iteration 0 / 1500: loss 406.435177
iteration 100 / 1500: loss 240.962028
iteration 200 / 1500: loss 147.149135
iteration 300 / 1500: loss 91.117906
iteration 400 / 1500: loss 55.958883
iteration 500 / 1500: loss 35.796272
iteration 600 / 1500: loss 23.572038
iteration 700 / 1500: loss 15.919244
iteration 800 / 1500: loss 11.825303
iteration 900 / 1500: loss 8.800333
iteration 1000 / 1500: loss 7.771851
iteration 1100 / 1500: loss 6.350576
iteration 1200 / 1500: loss 6.481001
iteration 1300 / 1500: loss 5.467648
iteration 1400 / 1500: loss 5.335129
That took 3.092999s
```

In [12]:

```
# A useful debugging strategy is to plot the Loss as a function of
# iteration number:
plt.plot(loss_hist)
plt.xlabel('Iteration number')
plt.ylabel('Loss value')
plt.show()
```



```
In [13]: # Write the LinearSVM.predict function and evaluate the performance on both the
# training and validation set
y_train_pred = svm.predict(X_train)
print('training accuracy: %f' % (np.mean(y_train == y_train_pred), ))
y_val_pred = svm.predict(X_val)
print('validation accuracy: %f' % (np.mean(y_val == y_val_pred), ))
```

training accuracy: 0.379714
validation accuracy: 0.383000

```
In [18]: # Use the validation set to tune hyperparameters (regularization strength and
# Learning rate). You should experiment with different ranges for the Learning
# rates and regularization strengths; if you are careful you should be able to
# get a classification accuracy of about 0.39 on the validation set.

# Note: you may see runtime/overflow warnings during hyper-parameter search.
# This may be caused by extreme values, and is not a bug.

# results is dictionary mapping tuples of the form
# (learning_rate, regularization_strength) to tuples of the form
# (training_accuracy, validation_accuracy). The accuracy is simply the fraction
# of data points that are correctly classified.
results = {}
best_val = -1 # The highest validation accuracy that we have seen so far.
best_svm = None # The LinearSVM object that achieved the highest validation rate

#####
# TODO:
# Write code that chooses the best hyperparameters by tuning on the validation #
# set. For each combination of hyperparameters, train a Linear SVM on the      #
# training set, compute its accuracy on the training and validation sets, and    #
# store these numbers in the results dictionary. In addition, store the best     #
#
```

```

# validation accuracy in best_val and the LinearSVM object that achieves this #
# accuracy in best_svm.
#
# Hint: You should use a small value for num_iters as you develop your #
# validation code so that the SVMs don't take much time to train; once you are #
# confident that your validation code works, you should rerun the validation   #
# code with a larger value for num_iters.
#####
# Provided as a reference. You may or may not want to change these hyperparameters
learning_rates = [1e-7, 15e-8]
regularization_strengths = [1e3, 1e4]

# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

for rs in regularization_strengths:
    for lr in learning_rates:
        svm = LinearSVM()
        loss_hist = svm.train(X_train, y_train, lr, rs, num_iters=5000)
        y_train_pred = svm.predict(X_train)
        train_accuracy = np.mean(y_train_pred == y_train)

        y_val_pred = svm.predict(X_val)
        test_accuracy = np.mean(y_val == y_val_pred)
        if test_accuracy > best_val:
            best_val = test_accuracy
            best_svm = svm
        results[(lr,rs)] = train_accuracy, test_accuracy

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****


# Print out results.
for lr, reg in sorted(results):
    train_accuracy, val_accuracy = results[(lr, reg)]
    print('lr %e reg %e train accuracy: %f val accuracy: %f' % (
        lr, reg, train_accuracy, val_accuracy))

print('best validation accuracy achieved during cross-validation: %f' % best_val)

```

```

lr 1.000000e-07 reg 1.000000e+03 train accuracy: 0.373837 val accuracy: 0.371000
lr 1.000000e-07 reg 1.000000e+04 train accuracy: 0.394184 val accuracy: 0.391000
lr 1.500000e-07 reg 1.000000e+03 train accuracy: 0.392306 val accuracy: 0.379000
lr 1.500000e-07 reg 1.000000e+04 train accuracy: 0.392020 val accuracy: 0.404000
best validation accuracy achieved during cross-validation: 0.404000

```

In [19]: # Visualize the cross-validation results

```

import math
import pdb

# pdb.set_trace()

x_scatter = [math.log10(x[0]) for x in results]
y_scatter = [math.log10(x[1]) for x in results]

# plot training accuracy
marker_size = 100
colors = [results[x][0] for x in results]
plt.subplot(2, 1, 1)
plt.tight_layout(pad=3)
plt.scatter(x_scatter, y_scatter, marker_size, c=colors, cmap=plt.cm.coolwarm)

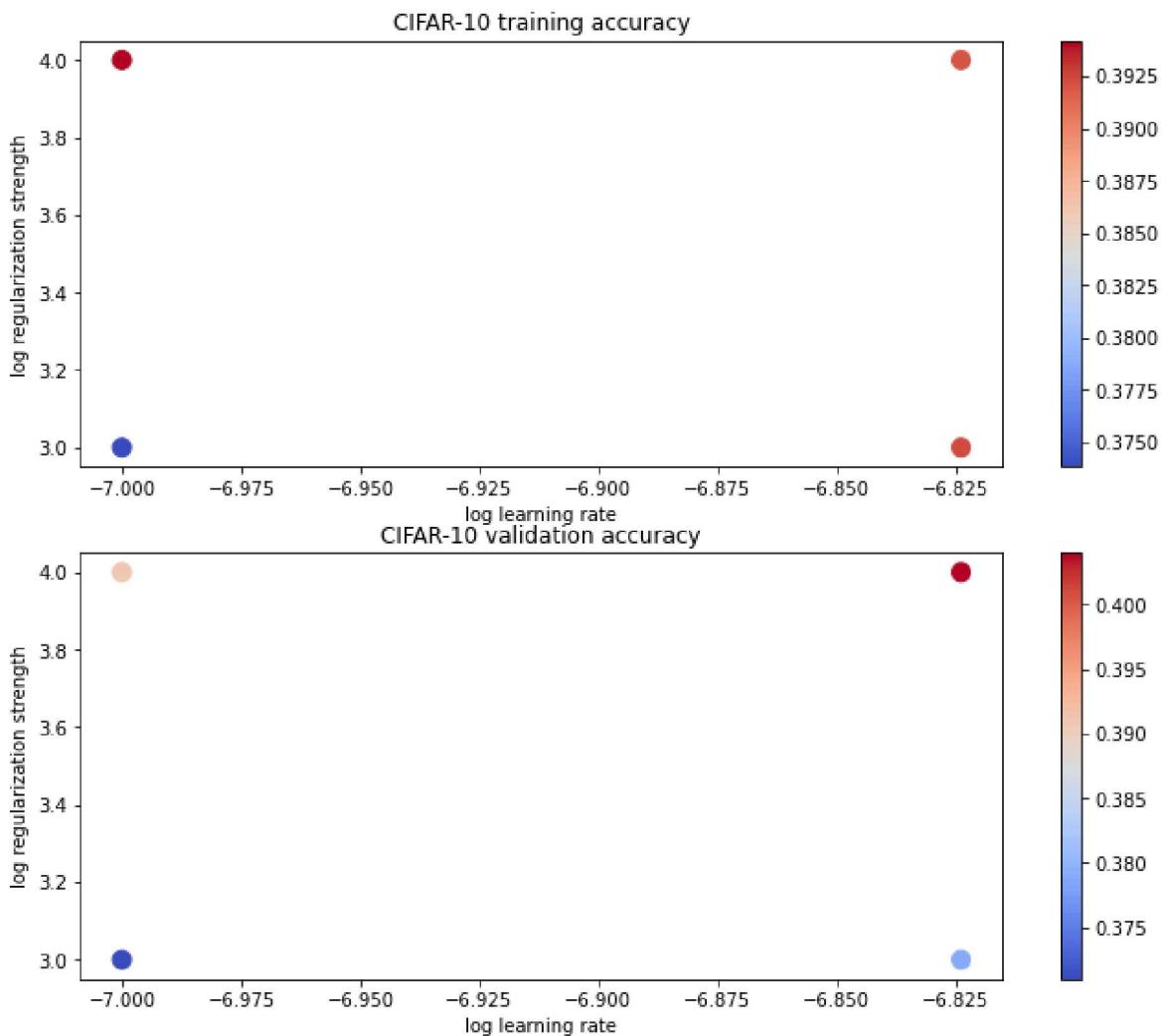
```

```

plt.colorbar()
plt.xlabel('log learning rate')
plt.ylabel('log regularization strength')
plt.title('CIFAR-10 training accuracy')

# plot validation accuracy
colors = [results[x][1] for x in results] # default size of markers is 20
plt.subplot(2, 1, 2)
plt.scatter(x_scatter, y_scatter, marker_size, c=colors, cmap=plt.cm.coolwarm)
plt.colorbar()
plt.xlabel('log learning rate')
plt.ylabel('log regularization strength')
plt.title('CIFAR-10 validation accuracy')
plt.show()

```



In [20]:

```

# Evaluate the best svm on test set
y_test_pred = best_svm.predict(X_test)
test_accuracy = np.mean(y_test == y_test_pred)
print('linear SVM on raw pixels final test set accuracy: %f' % test_accuracy)

```

linear SVM on raw pixels final test set accuracy: 0.386000

In [21]:

```

# Visualize the Learned weights for each class.
# Depending on your choice of Learning rate and regularization strength, these m
# or may not be nice to look at.
w = best_svm.W[:-1, :] # strip out the bias
w = w.reshape(32, 32, 3, 10)
w_min, w_max = np.min(w), np.max(w)
classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship'

```

```
for i in range(10):
    plt.subplot(2, 5, i + 1)

    # Rescale the weights to be between 0 and 255
    wimg = 255.0 * (w[:, :, :, i].squeeze() - w_min) / (w_max - w_min)
    plt.imshow(wimg.astype('uint8'))
    plt.axis('off')
    plt.title(classes[i])
```



Inline question 2

Describe what your visualized SVM weights look like, and offer a brief explanation for why they look they way that they do.

Your Answer : You can treat SVM as a kind of template matching. Because during the process, you can find the hyper-plane to distinguish the different region.

Softmax exercise

Complete and hand in this completed worksheet (including its outputs and any supporting code outside of the worksheet) with your assignment submission. For more details see the [assignments page](#) on the course website.

This exercise is analogous to the SVM exercise. You will:

- implement a fully-vectorized **loss function** for the Softmax classifier
- implement the fully-vectorized expression for its **analytic gradient**
- **check your implementation** with numerical gradient
- use a validation set to **tune the learning rate and regularization strength**
- **optimize** the loss function with **SGD**
- **visualize** the final learned weights

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

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2
```



```
In [2]: def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000, num_workers=3):
    """
    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 = 'cs231n/datasets/cifar-10-batches-py'

    # Cleaning up variables to prevent Loading data multiple times (which may cause memory issues)
    try:
        del X_train, y_train
        del X_test, y_test
        print('Clear previously loaded data.')
    except:
        pass

    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,)

Softmax Classifier

Your code for this section will all be written inside
 cs231n/classifiers/softmax.py .

```
In [12]: # First implement the naive softmax Loss function with nested loops.
# Open the file cs231n/classifiers/softmax.py and implement the
# softmax_loss_naive function.

from cs231n.classifiers.softmax import softmax_loss_naive
import time

# Generate a random softmax weight matrix and use it to compute the Loss.
W = np.random.randn(3073, 10) * 0.0001
loss, grad = softmax_loss_naive(W, X_dev, y_dev, 0.0)

# As a rough sanity check, our Loss should be something close to -Log(0.1).
print('loss: %f' % loss)
print('sanity check: %f' % (-np.log(0.1)))
```

loss: 2.363396
 sanity check: 2.302585

Inline Question 1

Why do we expect our loss to be close to $-\log(0.1)$? Explain briefly.**

Your Answer :

Since the weight matrix W is uniform randomly selected, then the predicted prob of each class is identically equal to $\frac{1}{10}$, where 10 is the number of classes.
 Therefore, from the def of the cross entropy, we can get for each class the result is equal to $-\log(0.1)$, which should be equal to the loss.

```
In [13]: # Complete the implementation of softmax_loss_naive and implement a (naive)
# version of the gradient that uses nested loops.
loss, grad = softmax_loss_naive(W, X_dev, y_dev, 0.0)

# As we did for the SVM, use numeric gradient checking as a debugging tool.
# The numeric gradient should be close to the analytic gradient.
from cs231n.gradient_check import grad_check_sparse
f = lambda w: softmax_loss_naive(w, X_dev, y_dev, 0.0)[0]
grad_numerical = grad_check_sparse(f, W, grad, 10)

# similar to SVM case, do another gradient check with regularization
loss, grad = softmax_loss_naive(W, X_dev, y_dev, 5e1)
f = lambda w: softmax_loss_naive(w, X_dev, y_dev, 5e1)[0]
grad_numerical = grad_check_sparse(f, W, grad, 10)
```

```
numerical: 0.682840 analytic: 0.682840, relative error: 5.625545e-08
numerical: -1.813164 analytic: -1.813164, relative error: 1.167975e-08
numerical: -0.085388 analytic: -0.085388, relative error: 2.590980e-07
numerical: 1.755841 analytic: 1.755841, relative error: 4.711730e-09
numerical: -0.629661 analytic: -0.629661, relative error: 8.946373e-09
numerical: 2.097693 analytic: 2.097693, relative error: 1.145659e-08
numerical: 0.898483 analytic: 0.898483, relative error: 6.996410e-08
numerical: -1.142693 analytic: -1.142693, relative error: 3.744289e-08
numerical: 0.023064 analytic: 0.023064, relative error: 8.981415e-07
numerical: 0.848872 analytic: 0.848872, relative error: 4.278324e-09
numerical: -0.119443 analytic: -0.119443, relative error: 6.096139e-07
numerical: -1.627726 analytic: -1.627726, relative error: 1.915281e-08
numerical: -0.363066 analytic: -0.363066, relative error: 1.058739e-07
numerical: 2.311269 analytic: 2.311269, relative error: 1.119167e-09
numerical: 1.738757 analytic: 1.738757, relative error: 3.599364e-08
numerical: 0.077635 analytic: 0.077635, relative error: 2.577279e-07
numerical: 1.041345 analytic: 1.041345, relative error: 2.913598e-08
numerical: -3.048888 analytic: -3.048888, relative error: 1.455938e-08
numerical: -4.474690 analytic: -4.474690, relative error: 6.577992e-09
numerical: 2.479441 analytic: 2.479441, relative error: 2.906516e-08
```

In [14]:

```
# Now that we have a naive implementation of the softmax loss function and its gradients,
# implement a vectorized version in softmax_loss_vectorized.
# The two versions should compute the same results, but the vectorized version should be much faster.
tic = time.time()
loss_naive, grad_naive = softmax_loss_naive(W, X_dev, y_dev, 0.000005)
toc = time.time()
print('naive loss: %e computed in %fs' % (loss_naive, toc - tic))

from cs231n.classifiers.softmax import softmax_loss_vectorized
tic = time.time()
loss_vectorized, grad_vectorized = softmax_loss_vectorized(W, X_dev, y_dev, 0.000005)
toc = time.time()
print('vectorized loss: %e computed in %fs' % (loss_vectorized, toc - tic))

# As we did for the SVM, we use the Frobenius norm to compare the two versions
# of the gradient.
grad_difference = np.linalg.norm(grad_naive - grad_vectorized, ord='fro')
print('Loss difference: %f' % np.abs(loss_naive - loss_vectorized))
print('Gradient difference: %f' % grad_difference)
```

```
naive loss: 2.363396e+00 computed in 0.096997s
vectorized loss: 2.363396e+00 computed in 0.002001s
Loss difference: 0.000000
Gradient difference: 0.000000
```

In [15]:

```
# Use the validation set to tune hyperparameters (regularization strength and learning rate).
# You should experiment with different ranges for the learning rates and regularization strengths; if you are careful you should be able to get a classification accuracy of over 0.35 on the validation set.

from cs231n.classifiers import Softmax
results = {}
best_val = -1
best_softmax = None

#####
# TODO: # Use the validation set to set the learning rate and regularization strength. #
#####
```

```
# This should be identical to the validation that you did for the SVM; save      #
# the best trained softmax classifier in best_softmax.                          #
#####
# Provided as a reference. You may or may not want to change these hyperparameters#
learning_rates = [1e-7, 5e-7]                                                 #
regularization_strengths = [2.5e4, 5e4]                                         #

# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
for rs in regularization_strengths:
    for lr in learning_rates:
        softmax = Softmax()
        loss_hist = softmax.train(X_train, y_train, lr, rs, num_iters=5000)
        y_train_pred = softmax.predict(X_train)
        train_accuracy = np.mean(y_train_pred == y_train)

        y_val_pred = softmax.predict(X_val)
        test_accuracy = np.mean(y_val == y_val_pred)
        if test_accuracy > best_val:
            best_val = test_accuracy
            best_softmax = softmax
        results[(lr,rs)] = train_accuracy, test_accuracy

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

# Print out results.
for lr, reg in sorted(results):
    train_accuracy, val_accuracy = results[(lr, reg)]
    print('lr %e reg %e train accuracy: %f val accuracy: %f' % (
        lr, reg, train_accuracy, val_accuracy))

print('best validation accuracy achieved during cross-validation: %f' % best_val)
```

lr 1.000000e-07 reg 2.500000e+04 train accuracy: 0.351510 val accuracy: 0.359000
 lr 1.000000e-07 reg 5.000000e+04 train accuracy: 0.335245 val accuracy: 0.353000
 lr 5.000000e-07 reg 2.500000e+04 train accuracy: 0.353388 val accuracy: 0.364000
 lr 5.000000e-07 reg 5.000000e+04 train accuracy: 0.333122 val accuracy: 0.334000
 best validation accuracy achieved during cross-validation: 0.364000

In [16]: # evaluate on test set
 # Evaluate the best softmax on test set
 y_test_pred = best_softmax.predict(X_test)
 test_accuracy = np.mean(y_test == y_test_pred)
 print('softmax on raw pixels final test set accuracy: %f' % (test_accuracy,))

softmax on raw pixels final test set accuracy: 0.356000

Inline Question 2 - True or False

Suppose the overall training loss is defined as the sum of the per-datapoint loss over all training examples. It is possible to add a new datapoint to a training set that would leave the SVM loss unchanged, but this is not the case with the Softmax classifier loss.

Your Answer :

True.

Your Explanation :

That is because, for a SVM classifier, if you add a new data point, and the model

classifies correctly, then it will not have an influence on the svm loss. But for softmax model, because of its function, we need to calculate the sum of the exp scores for all data points. Therefore, adding a new datapoint will have an influence on softmax loss.

```
In [17]: # Visualize the Learned weights for each class
w = best_softmax.W[:-1,:,:] # strip out the bias
w = w.reshape(32, 32, 3, 10)

w_min, w_max = np.min(w), np.max(w)

classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship'
for i in range(10):
    plt.subplot(2, 5, i + 1)

    # Rescale the weights to be between 0 and 255
    wimg = 255.0 * (w[:, :, :, i].squeeze() - w_min) / (w_max - w_min)
    plt.imshow(wimg.astype('uint8'))
    plt.axis('off')
    plt.title(classes[i])
```



Implementing a Neural Network

In this exercise we will develop a neural network with fully-connected layers to perform classification, and test it out on the CIFAR-10 dataset.

```
In [2]: # A bit of setup

import numpy as np
import matplotlib.pyplot as plt

from cs231n.classifiers.neural_net import TwoLayerNet

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

def rel_error(x, y):
    """ returns relative error """
    return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))



```

We will use the class `TwoLayerNet` in the file

`cs231n/classifiers/neural_net.py` to represent instances of our network. The network parameters are stored in the instance variable `self.params` where keys are string parameter names and values are numpy arrays. Below, we initialize toy data and a toy model that we will use to develop your implementation.

```
In [3]: # Create a small net and some toy data to check your implementations.
# Note that we set the random seed for repeatable experiments.

input_size = 4
hidden_size = 10
num_classes = 3
num_inputs = 5

def init_toy_model():
    np.random.seed(0)
    return TwoLayerNet(input_size, hidden_size, num_classes, std=1e-1)

def init_toy_data():
    np.random.seed(1)
    X = 10 * np.random.randn(num_inputs, input_size)
    y = np.array([0, 1, 2, 2, 1])
    return X, y

net = init_toy_model()
X, y = init_toy_data()
```

Forward pass: compute scores

Open the file `cs231n/classifiers/neural_net.py` and look at the method `TwoLayerNet.loss`. This function is very similar to the loss functions you have written for the SVM and Softmax exercises: It takes the data and weights and computes the class scores, the loss, and the gradients on the parameters.

Implement the first part of the forward pass which uses the weights and biases to compute the scores for all inputs.

```
In [4]: scores = net.loss(X)
print('Your scores:')
print(scores)
print()
print('correct scores:')
correct_scores = np.asarray([
    [-0.81233741, -1.27654624, -0.70335995],
    [-0.17129677, -1.18803311, -0.47310444],
    [-0.51590475, -1.01354314, -0.8504215 ],
    [-0.15419291, -0.48629638, -0.52901952],
    [-0.00618733, -0.12435261, -0.15226949]])
print(correct_scores)
print()

# The difference should be very small. We get < 1e-7
print('Difference between your scores and correct scores:')
print(np.sum(np.abs(scores - correct_scores)))
```

Your scores:
`[[-0.81233741 -1.27654624 -0.70335995]
[-0.17129677 -1.18803311 -0.47310444]
[-0.51590475 -1.01354314 -0.8504215]
[-0.15419291 -0.48629638 -0.52901952]
[-0.00618733 -0.12435261 -0.15226949]]`

correct scores:
`[[-0.81233741 -1.27654624 -0.70335995]
[-0.17129677 -1.18803311 -0.47310444]
[-0.51590475 -1.01354314 -0.8504215]
[-0.15419291 -0.48629638 -0.52901952]
[-0.00618733 -0.12435261 -0.15226949]]`

Difference between your scores and correct scores:
`3.6802720745909845e-08`

Forward pass: compute loss

In the same function, implement the second part that computes the data and regularization loss.

```
In [5]: loss, _ = net.loss(X, y, reg=0.1)
correct_loss = 1.30378789133
```

```
# should be very small, we get < 1e-12
print('Difference between your loss and correct loss:')
print(np.sum(np.abs(loss - correct_loss)))
```

Difference between your loss and correct loss:
1.7985612998927536e-13

Backward pass

Implement the rest of the function. This will compute the gradient of the loss with respect to the variables `w1`, `b1`, `w2`, and `b2`. Now that you (hopefully!) have a correctly implemented forward pass, you can debug your backward pass using a numeric gradient check:

```
In [6]: from cs231n.gradient_check import eval_numerical_gradient

# Use numeric gradient checking to check your implementation of the backward pass
# If your implementation is correct, the difference between the numeric and
# analytic gradients should be less than 1e-8 for each of w1, w2, b1, and b2.

loss, grads = net.loss(X, y, reg=0.1)

# these should all be less than 1e-8 or so
for param_name in grads:
    f = lambda W: net.loss(X, y, reg=0.1)[0]
    param_grad_num = eval_numerical_gradient(f, net.params[param_name], verbose=False)
    print('%s max relative error: %e' % (param_name, rel_error(param_grad_num, g
```

W2 max relative error: 3.440708e-09
b2 max relative error: 4.447625e-11
W1 max relative error: 3.561318e-09
b1 max relative error: 2.738421e-09

Train the network

To train the network we will use stochastic gradient descent (SGD), similar to the SVM and Softmax classifiers. Look at the function `TwoLayerNet.train` and fill in the missing sections to implement the training procedure. This should be very similar to the training procedure you used for the SVM and Softmax classifiers. You will also have to implement `TwoLayerNet.predict`, as the training process periodically performs prediction to keep track of accuracy over time while the network trains.

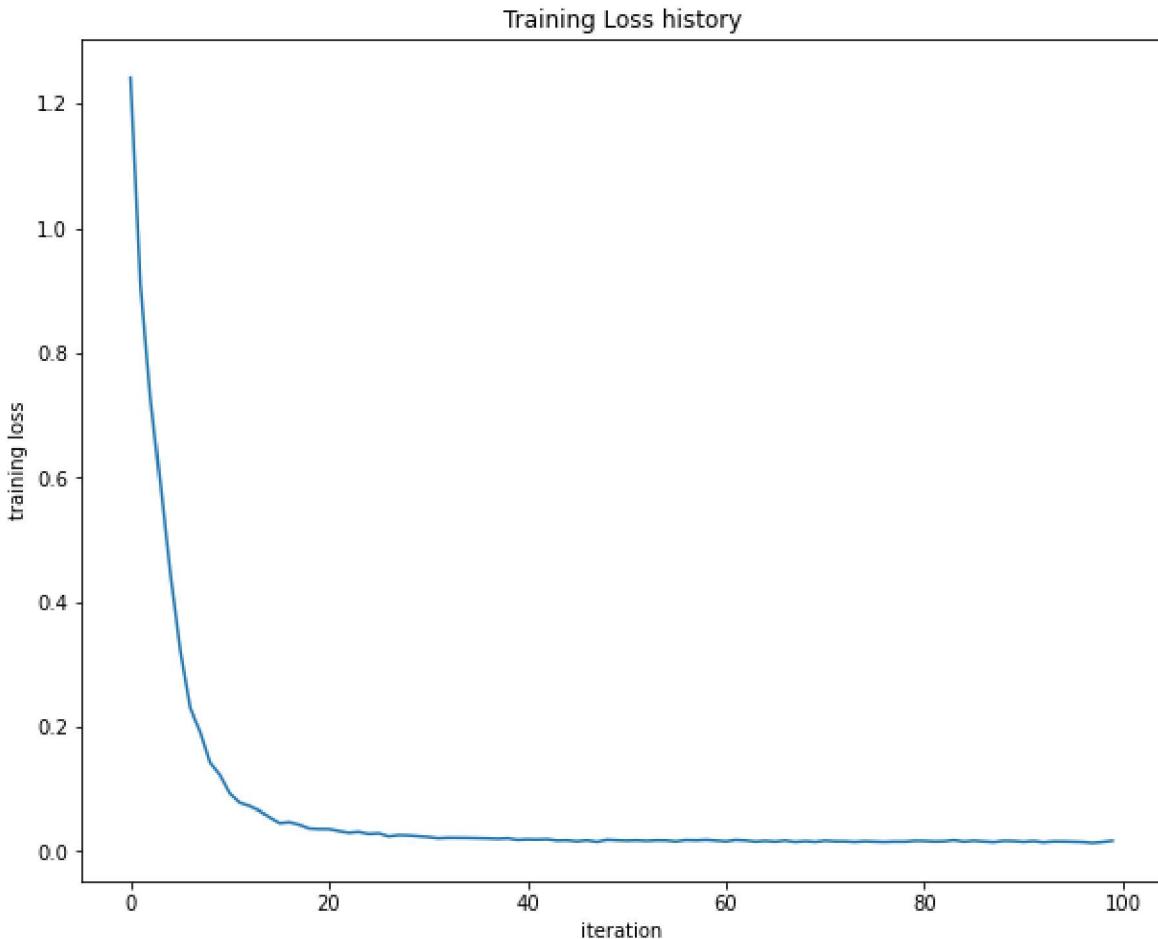
Once you have implemented the method, run the code below to train a two-layer network on toy data. You should achieve a training loss less than 0.02.

```
In [7]: net = init_toy_model()
stats = net.train(X, y, X, y,
                  learning_rate=1e-1, reg=5e-6,
                  num_iters=100, verbose=False)

print('Final training loss: ', stats['loss_history'][-1])
```

```
# plot the loss history
plt.plot(stats['loss_history'])
plt.xlabel('iteration')
plt.ylabel('training loss')
plt.title('Training Loss history')
plt.show()
```

Final training loss: 0.01714364353292381



Load the data

Now that you have implemented a two-layer network that passes gradient checks and works on toy data, it's time to load up our favorite CIFAR-10 data so we can use it to train a classifier on a real dataset.

```
In [8]: from cs231n.data_utils import load_CIFAR10

def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000):
    """
    Load the CIFAR-10 dataset from disk and perform preprocessing to prepare
    it for the two-layer neural net 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 = 'cs231n/datasets/cifar-10-batches-py'

    # Cleaning up variables to prevent loading data multiple times (which may ca
    try:
        del X_train, y_train
```

```

    del X_test, y_test
    print('Clear previously loaded data.')
except:
    pass

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]

# 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

# Reshape data to rows
X_train = X_train.reshape(num_training, -1)
X_val = X_val.reshape(num_validation, -1)
X_test = X_test.reshape(num_test, -1)

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

# Invoke the above function to get our data.
X_train, y_train, X_val, y_val, X_test, y_test = 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)

```

```

Train data shape: (49000, 3072)
Train labels shape: (49000,)
Validation data shape: (1000, 3072)
Validation labels shape: (1000,)
Test data shape: (1000, 3072)
Test labels shape: (1000,)

```

Train a network

To train our network we will use SGD. In addition, we will adjust the learning rate with an exponential learning rate schedule as optimization proceeds; after each epoch, we will reduce the learning rate by multiplying it by a decay rate.

```
In [9]: input_size = 32 * 32 * 3
hidden_size = 50
num_classes = 10
```

```

net = TwoLayerNet(input_size, hidden_size, num_classes)

# Train the network
stats = net.train(X_train, y_train, X_val, y_val,
                   num_iters=1000, batch_size=200,
                   learning_rate=1e-4, learning_rate_decay=0.95,
                   reg=0.25, verbose=True)

# Predict on the validation set
val_acc = (net.predict(X_val) == y_val).mean()
print('Validation accuracy: ', val_acc)

```

```

iteration 0 / 1000: loss 2.302762
iteration 100 / 1000: loss 2.302358
iteration 200 / 1000: loss 2.297404
iteration 300 / 1000: loss 2.258897
iteration 400 / 1000: loss 2.202975
iteration 500 / 1000: loss 2.116816
iteration 600 / 1000: loss 2.049789
iteration 700 / 1000: loss 1.985711
iteration 800 / 1000: loss 2.003726
iteration 900 / 1000: loss 1.948076
Validation accuracy:  0.287

```

Debug the training

With the default parameters we provided above, you should get a validation accuracy of about 0.29 on the validation set. This isn't very good.

One strategy for getting insight into what's wrong is to plot the loss function and the accuracies on the training and validation sets during optimization.

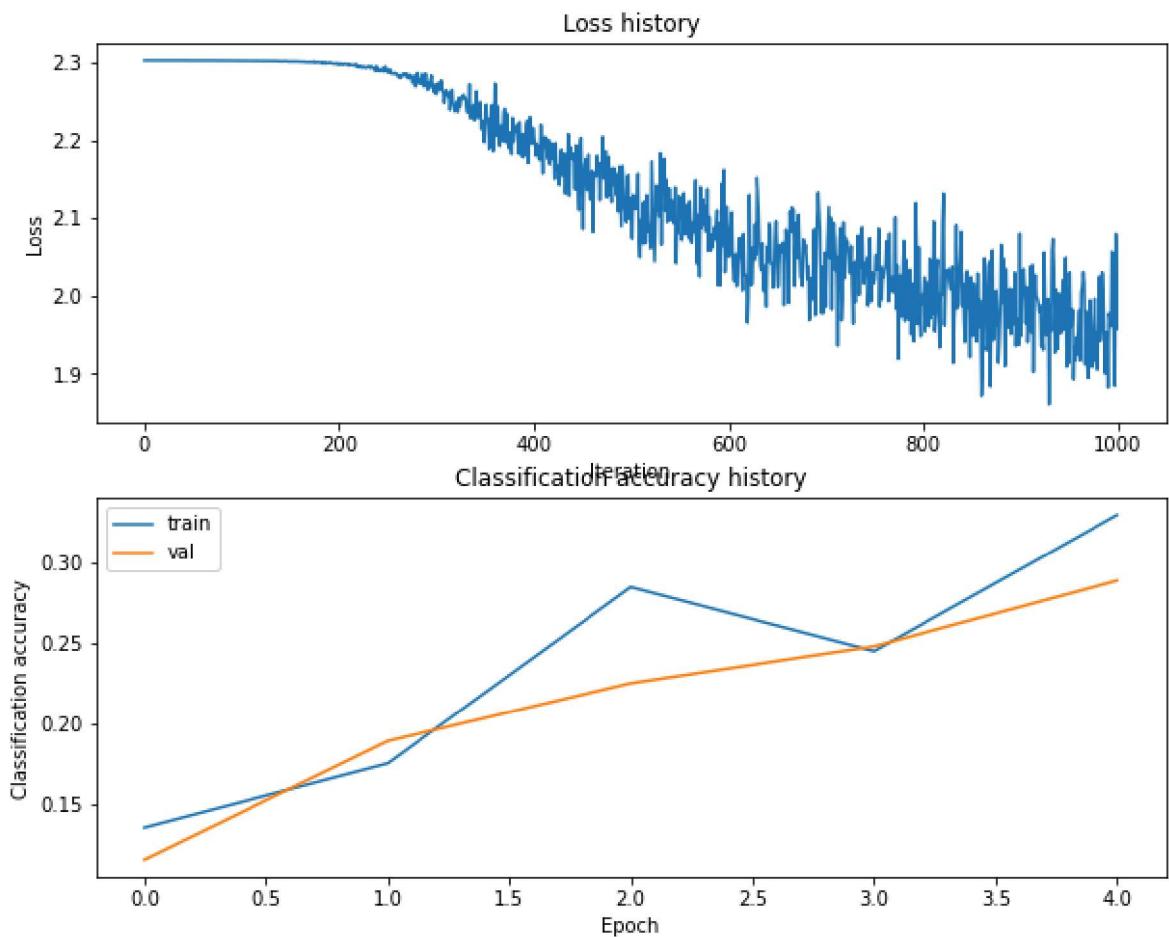
Another strategy is to visualize the weights that were learned in the first layer of the network. In most neural networks trained on visual data, the first layer weights typically show some visible structure when visualized.

```

In [10]: # Plot the loss function and train / validation accuracies
plt.subplot(2, 1, 1)
plt.plot(stats['loss_history'])
plt.title('Loss history')
plt.xlabel('Iteration')
plt.ylabel('Loss')

plt.subplot(2, 1, 2)
plt.plot(stats['train_acc_history'], label='train')
plt.plot(stats['val_acc_history'], label='val')
plt.title('Classification accuracy history')
plt.xlabel('Epoch')
plt.ylabel('Classification accuracy')
plt.legend()
plt.show()

```

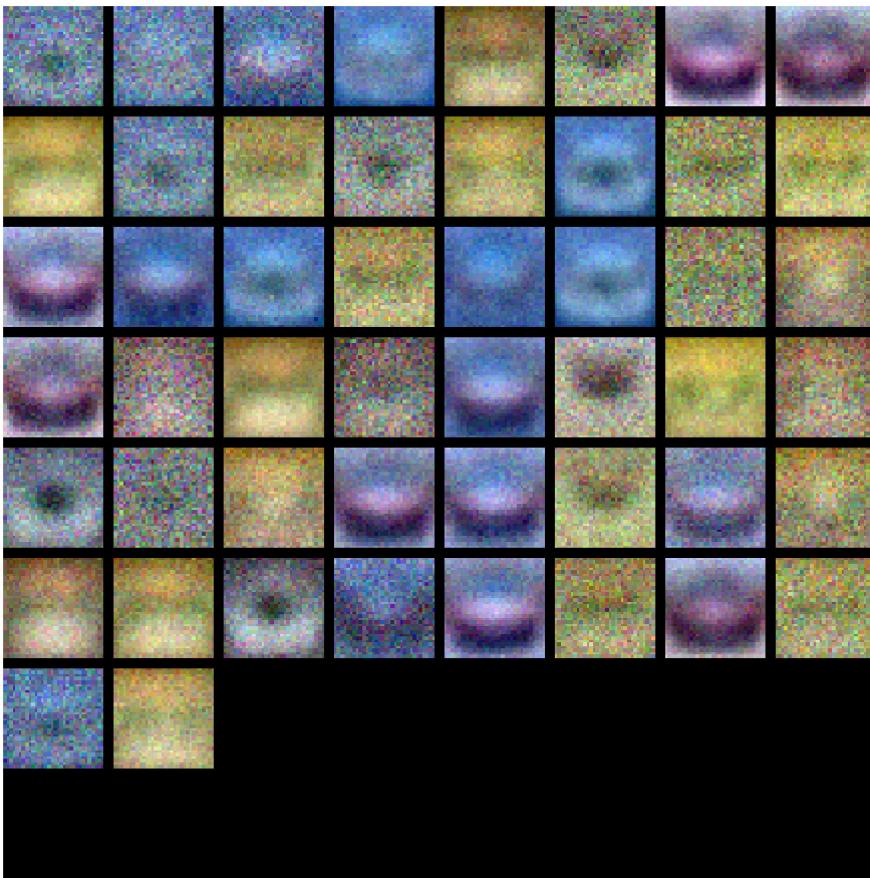


```
In [11]: from cs231n.vis_utils import visualize_grid

# Visualize the weights of the network

def show_net_weights(net):
    W1 = net.params['W1']
    W1 = W1.reshape(32, 32, 3, -1).transpose(3, 0, 1, 2)
    plt.imshow(visualize_grid(W1, padding=3).astype('uint8'))
    plt.gca().axis('off')
    plt.show()

show_net_weights(net)
```



Tune your hyperparameters

What's wrong? Looking at the visualizations above, we see that the loss is decreasing more or less linearly, which seems to suggest that the learning rate may be too low. Moreover, there is no gap between the training and validation accuracy, suggesting that the model we used has low capacity, and that we should increase its size. On the other hand, with a very large model we would expect to see more overfitting, which would manifest itself as a very large gap between the training and validation accuracy.

Tuning. Tuning the hyperparameters and developing intuition for how they affect the final performance is a large part of using Neural Networks, so we want you to get a lot of practice. Below, you should experiment with different values of the various hyperparameters, including hidden layer size, learning rate, number of training epochs, and regularization strength. You might also consider tuning the learning rate decay, but you should be able to get good performance using the default value.

Approximate results. You should aim to achieve a classification accuracy of greater than 48% on the validation set. Our best network gets over 52% on the validation set.

Experiment: Your goal in this exercise is to get as good of a result on CIFAR-10 as you can (52% could serve as a reference), with a fully-connected Neural Network.

Feel free implement your own techniques (e.g. PCA to reduce dimensionality, or adding dropout, or adding features to the solver, etc.).

Explain your hyperparameter tuning process below.

Your Answer :

Perform gridded tuning to filter hyperparameters by partitioning the validation set from the dataset

```
In [12]: best_net = None # store the best model into this

#####
# TODO: Tune hyperparameters using the validation set. Store your best trained
# model in best_net.
#
# To help debug your network, it may help to use visualizations similar to the
# ones we used above; these visualizations will have significant qualitative
# differences from the ones we saw above for the poorly tuned network.
#
# Tweaking hyperparameters by hand can be fun, but you might find it useful to
# write code to sweep through possible combinations of hyperparameters
# automatically like we did on the previous exercises.
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

results = {}
best_val = -1 # The highest validation accuracy that we have seen so far.
learning_rates = [1e-3, 1e-4, 5e-4]
regularization_strengths = [0.75, 1]
hidden_size = [50, 75]
for hs in hidden_size:
    for rs in regularization_strengths:
        for lr in learning_rates:
            net = TwoLayerNet(input_size, hs, num_classes)
            loss_hist = net.train(X_train, y_train, X_val, y_val, learning_rate=lr,
            y_train_pred = net.predict(X_train)
            train_accuracy = np.mean(y_train_pred == y_train)

            y_val_pred = net.predict(X_val)
            test_accuracy = np.mean(y_val == y_val_pred)
            if test_accuracy > best_val:
                best_val = test_accuracy
                best_net = net
            results[(lr,rs)] = test_accuracy

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
```

```
In [13]: # Print your validation accuracy: this should be above 48%
val_acc = (best_net.predict(X_val) == y_val).mean()
print('Validation accuracy: ', val_acc)
```

Validation accuracy: 0.522

```
In [14]: # Visualize the weights of the best network
show_net_weights(best_net)
```



Run on the test set

When you are done experimenting, you should evaluate your final trained network on the test set; you should get above 48%.

```
In [18]: # Print your test accuracy: this should be above 48%
test_acc = (best_net.predict(X_test) == y_test).mean()
print('Test accuracy: ', test_acc)
```

Test accuracy: 0.539

Inline Question

Now that you have trained a Neural Network classifier, you may find that your testing accuracy is much lower than the training accuracy. In what ways can we decrease this gap? Select all that apply.

1. Train on a larger dataset.
2. Add more hidden units.
3. Increase the regularization strength.
4. None of the above.

Your Answer :

1, 3

Your Explanation :

1. With a larger dataset, neural network can have a better training and learning, improving its performance while working on testing dataset.
2. Adding more hiddern units has double sides. Maybe it can improve your performance, because of more complex layers. But at the same time, it may lead to overfitting problem with too many units.
3. Increasing the regularization strength can also reduce the occurrence of the overfitting problems by penalizing overly complex models.

Image features exercise

Complete and hand in this completed worksheet (including its outputs and any supporting code outside of the worksheet) with your assignment submission. For more details see the [assignments page](#) on the course website.

We have seen that we can achieve reasonable performance on an image classification task by training a linear classifier on the pixels of the input image. In this exercise we will show that we can improve our classification performance by training linear classifiers not on raw pixels but on features that are computed from the raw pixels.

All of your work for this exercise will be done in this notebook.

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

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2
```

Load data

Similar to previous exercises, we will load CIFAR-10 data from disk.

```
In [2]: from cs231n.features import color_histogram_hsv, hog_feature

def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000):
    # Load the raw CIFAR-10 data
    cifar10_dir = 'cs231n/datasets/cifar-10-batches-py'

    # Cleaning up variables to prevent loading data multiple times (which may cause errors)
    del X_train, y_train
    del X_test, y_test
    print('Clear previously loaded data.')
except:
    pass

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]

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

X_train, y_train, X_val, y_val, X_test, y_test = get_CIFAR10_data()

```

Extract Features

For each image we will compute a Histogram of Oriented Gradients (HOG) as well as a color histogram using the hue channel in HSV color space. We form our final feature vector for each image by concatenating the HOG and color histogram feature vectors.

Roughly speaking, HOG should capture the texture of the image while ignoring color information, and the color histogram represents the color of the input image while ignoring texture. As a result, we expect that using both together ought to work better than using either alone. Verifying this assumption would be a good thing to try for your own interest.

The `hog_feature` and `color_histogram_hsv` functions both operate on a single image and return a feature vector for that image. The `extract_features` function takes a set of images and a list of feature functions and evaluates each feature function on each image, storing the results in a matrix where each column is the concatenation of all feature vectors for a single image.

```

In [3]: from cs231n.features import *

num_color_bins = 10 # Number of bins in the color histogram
feature_fns = [hog_feature, lambda img: color_histogram_hsv(img, nbin=num_color_bins)]
X_train_feats = extract_features(X_train, feature_fns, verbose=True)
X_val_feats = extract_features(X_val, feature_fns)
X_test_feats = extract_features(X_test, feature_fns)

# Preprocessing: Subtract the mean feature
mean_feat = np.mean(X_train_feats, axis=0, keepdims=True)
X_train_feats -= mean_feat
X_val_feats -= mean_feat
X_test_feats -= mean_feat

# Preprocessing: Divide by standard deviation. This ensures that each feature
# has roughly the same scale.
std_feat = np.std(X_train_feats, axis=0, keepdims=True)
X_train_feats /= std_feat
X_val_feats /= std_feat
X_test_feats /= std_feat

# Preprocessing: Add a bias dimension

```

```
X_train_feats = np.hstack([X_train_feats, np.ones((X_train_feats.shape[0], 1))])
X_val_feats = np.hstack([X_val_feats, np.ones((X_val_feats.shape[0], 1))])
X_test_feats = np.hstack([X_test_feats, np.ones((X_test_feats.shape[0], 1))])
```

```
Done extracting features for 1000 / 49000 images
Done extracting features for 2000 / 49000 images
Done extracting features for 3000 / 49000 images
Done extracting features for 4000 / 49000 images
Done extracting features for 5000 / 49000 images
Done extracting features for 6000 / 49000 images
Done extracting features for 7000 / 49000 images
Done extracting features for 8000 / 49000 images
Done extracting features for 9000 / 49000 images
Done extracting features for 10000 / 49000 images
Done extracting features for 11000 / 49000 images
Done extracting features for 12000 / 49000 images
Done extracting features for 13000 / 49000 images
Done extracting features for 14000 / 49000 images
Done extracting features for 15000 / 49000 images
Done extracting features for 16000 / 49000 images
Done extracting features for 17000 / 49000 images
Done extracting features for 18000 / 49000 images
Done extracting features for 19000 / 49000 images
Done extracting features for 20000 / 49000 images
Done extracting features for 21000 / 49000 images
Done extracting features for 22000 / 49000 images
Done extracting features for 23000 / 49000 images
Done extracting features for 24000 / 49000 images
Done extracting features for 25000 / 49000 images
Done extracting features for 26000 / 49000 images
Done extracting features for 27000 / 49000 images
Done extracting features for 28000 / 49000 images
Done extracting features for 29000 / 49000 images
Done extracting features for 30000 / 49000 images
Done extracting features for 31000 / 49000 images
Done extracting features for 32000 / 49000 images
Done extracting features for 33000 / 49000 images
Done extracting features for 34000 / 49000 images
Done extracting features for 35000 / 49000 images
Done extracting features for 36000 / 49000 images
Done extracting features for 37000 / 49000 images
Done extracting features for 38000 / 49000 images
Done extracting features for 39000 / 49000 images
Done extracting features for 40000 / 49000 images
Done extracting features for 41000 / 49000 images
Done extracting features for 42000 / 49000 images
Done extracting features for 43000 / 49000 images
Done extracting features for 44000 / 49000 images
Done extracting features for 45000 / 49000 images
Done extracting features for 46000 / 49000 images
Done extracting features for 47000 / 49000 images
Done extracting features for 48000 / 49000 images
Done extracting features for 49000 / 49000 images
```

Train SVM on features

Using the multiclass SVM code developed earlier in the assignment, train SVMs on top of the features extracted above; this should achieve better results than

training SVMs directly on top of raw pixels.

```
In [7]: # Use the validation set to tune the Learning rate and regularization strength

from cs231n.classifiers.linear_classifier import LinearSVM

learning_rates = [1e-9, 1e-8, 1e-7]
regularization_strengths = [5e4, 5e5, 5e6]

results = {}
best_val = -1
best_svm = None

#####
# TODO:
# Use the validation set to set the Learning rate and regularization strength.
# This should be identical to the validation that you did for the SVM; save
# the best trained classifier in best_svm. You might also want to play
# with different numbers of bins in the color histogram. If you are careful
# you should be able to get accuracy of near 0.44 on the validation set.
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

for rs in regularization_strengths:
    for lr in learning_rates:
        svm = LinearSVM()
        loss_hist = svm.train(X_train_feats, y_train, lr, rs, num_iters=5000)
        y_train_pred = svm.predict(X_train_feats)
        train_accuracy = np.mean(y_train_pred == y_train)

        y_val_pred = svm.predict(X_val_feats)
        test_accuracy = np.mean(y_val == y_val_pred)
        if test_accuracy > best_val:
            best_val = test_accuracy
            best_svm = svm
        results[(lr,rs)] = train_accuracy, test_accuracy

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****


# Print out results.
for lr, reg in sorted(results):
    train_accuracy, val_accuracy = results[(lr, reg)]
    print('lr %e reg %e train accuracy: %f val accuracy: %f' %
          (lr, reg, train_accuracy, val_accuracy))

print('best validation accuracy achieved during cross-validation: %f' % best_val)

lr 1.000000e-09 reg 5.000000e+04 train accuracy: 0.092082 val accuracy: 0.092000
lr 1.000000e-09 reg 5.000000e+05 train accuracy: 0.096102 val accuracy: 0.089000
lr 1.000000e-09 reg 5.000000e+06 train accuracy: 0.415204 val accuracy: 0.420000
lr 1.000000e-08 reg 5.000000e+04 train accuracy: 0.136388 val accuracy: 0.137000
lr 1.000000e-08 reg 5.000000e+05 train accuracy: 0.415163 val accuracy: 0.416000
lr 1.000000e-08 reg 5.000000e+06 train accuracy: 0.416653 val accuracy: 0.419000
lr 1.000000e-07 reg 5.000000e+04 train accuracy: 0.413796 val accuracy: 0.424000
lr 1.000000e-07 reg 5.000000e+05 train accuracy: 0.406612 val accuracy: 0.399000
lr 1.000000e-07 reg 5.000000e+06 train accuracy: 0.362755 val accuracy: 0.357000
best validation accuracy achieved during cross-validation: 0.424000
```

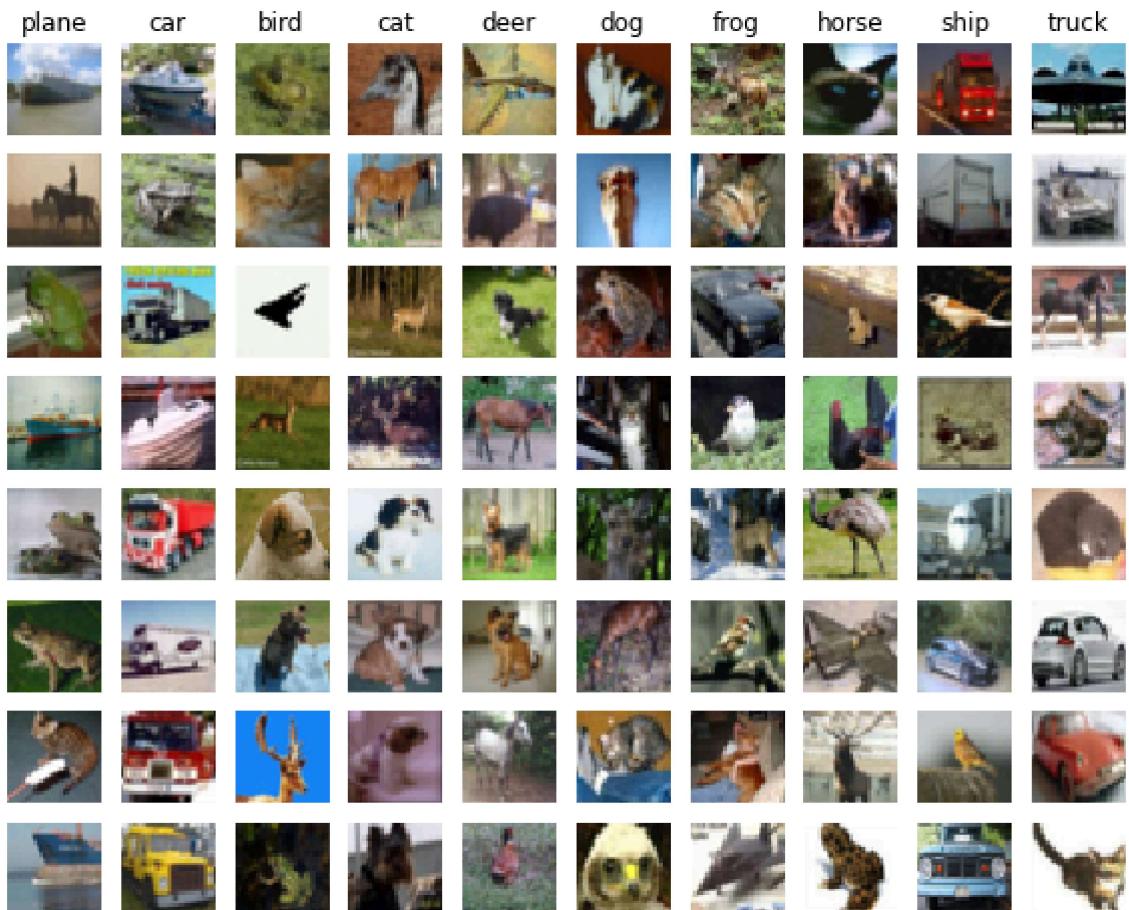
```
In [8]: # Evaluate your trained SVM on the test set: you should be able to get at least
y_test_pred = best_svm.predict(X_test_feats)
```

```
test_accuracy = np.mean(y_test == y_test_pred)
print(test_accuracy)
```

0.426

In [9]: # An important way to gain intuition about how an algorithm works is to
visualize the mistakes that it makes. In this visualization, we show examples
of images that are misclassified by our current system. The first column
shows images that our system labeled as "plane" but whose true label is
something other than "plane".

```
examples_per_class = 8
classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship'
for cls, cls_name in enumerate(classes):
    idxs = np.where((y_test != cls) & (y_test_pred == cls))[0]
    idxs = np.random.choice(idxs, examples_per_class, replace=False)
    for i, idx in enumerate(idxs):
        plt.subplot(examples_per_class, len(classes), i * len(classes) + cls + 1)
        plt.imshow(X_test[idx].astype('uint8'))
        plt.axis('off')
        if i == 0:
            plt.title(cls_name)
plt.show()
```



Inline question 1:

Describe the misclassification results that you see. Do they make sense?

Your Answer :

Do make sense, because the main kinds of features we use are HOG feature and color histogram feature. For example, there are lots of frog images misclassified

as plane, and they have same shape and color. Also, for birds, many misclassified images have color background. Therefore, they do make sense.

Neural Network on image features

Earlier in this assignment we saw that training a two-layer neural network on raw pixels achieved better classification performance than linear classifiers on raw pixels. In this notebook we have seen that linear classifiers on image features outperform linear classifiers on raw pixels.

For completeness, we should also try training a neural network on image features. This approach should outperform all previous approaches: you should easily be able to achieve over 55% classification accuracy on the test set; our best model achieves about 60% classification accuracy.

```
In [10]: # Preprocessing: Remove the bias dimension
# Make sure to run this cell only ONCE
print(X_train_feats.shape)
X_train_feats = X_train_feats[:, :-1]
X_val_feats = X_val_feats[:, :-1]
X_test_feats = X_test_feats[:, :-1]

print(X_train_feats.shape)
```

(49000, 155)
(49000, 154)

```
In [14]: from cs231n.classifiers.neural_net import TwoLayerNet

input_dim = X_train_feats.shape[1]
hidden_dim = 500
num_classes = 10

net = TwoLayerNet(input_dim, hidden_dim, num_classes)
best_net = None

#####
# TODO: Train a two-Layer neural network on image features. You may want to #
# cross-validate various parameters as in previous sections. Store your best #
# model in the best_net variable. #
#####

# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
results = {}
best_val = -1
best_net = None

learning_rates = [1e-2, 1e-1, 5e-1, 1, 5]
regularization_strengths = [1e-3, 5e-3, 1e-2, 1e-1, 0.5, 1]

for lr in learning_rates:
    for reg in regularization_strengths:
        net = TwoLayerNet(input_dim, hidden_dim, num_classes)
        # Train the network
        stats = net.train(X_train_feats, y_train, X_val_feats, y_val,
                          num_iters=1500, learning_rate=lr, reg=reg)
        val_acc = (net.predict(X_val_feats) == y_val).mean()

        if val_acc > best_val:
            best_val = val_acc
            best_net = net
```

```

    if val_acc > best_val:
        best_val = val_acc
        best_net = net
    results[(lr,reg)] = val_acc

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

```

```

e:\CS280\CS280-Fall24-Assignment1\Homework1_partB\cs231n\classifiers\neural_net.py:104: RuntimeWarning: divide by zero encountered in log
    loss = -np.sum(np.log(exp_scores[range(N), list(y)])))
e:\CS280\CS280-Fall24-Assignment1\Homework1_partB\cs231n\classifiers\neural_net.py:102: RuntimeWarning: overflow encountered in subtract
    shift_scores = scores - np.max(scores, axis = 1, keepdims= True)
e:\CS280\CS280-Fall24-Assignment1\Homework1_partB\cs231n\classifiers\neural_net.py:102: RuntimeWarning: invalid value encountered in subtract
    shift_scores = scores - np.max(scores, axis = 1, keepdims= True)
e:\CS280\CS280-Fall24-Assignment1\Homework1_partB\cs231n\classifiers\neural_net.py:106: RuntimeWarning: overflow encountered in scalar add
    loss += 0.5 * reg * (np.sum(W1* W1) + np.sum(W2* W2))
c:\Users\86136\anaconda3\lib\site-packages\numpy\core\fromnumeric.py:88: RuntimeWarning: overflow encountered in reduce
    return ufunc.reduce(obj, axis, dtype, out, **passkwargs)
e:\CS280\CS280-Fall24-Assignment1\Homework1_partB\cs231n\classifiers\neural_net.py:106: RuntimeWarning: overflow encountered in multiply
    loss += 0.5 * reg * (np.sum(W1* W1) + np.sum(W2* W2))

```

In [15]: # Run your best neural net classifier on the test set. You should be able
to get more than 55% accuracy.

```

test_acc = (best_net.predict(X_test_feats) == y_test).mean()
print(test_acc)

```

0.577