

knn

December 30, 2025

```
[1]: # This mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive')

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'icv83551/assignments/assignment1/'
FOLDERNAME = 'icv83551/assignments/assignment1/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

# Now that we've mounted your Drive, this ensures that
# the Python interpreter of the Colab VM can load
# python files from within it.
import sys
sys.path.append('/content/drive/My Drive/{}'.format(FOLDERNAME))

# This downloads the CIFAR-10 dataset to your Drive
# if it doesn't already exist.
%cd /content/drive/My\ Drive/$FOLDERNAME/icv83551/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

```
Mounted at /content/drive
/content/drive/My Drive/icv83551/assignments/assignment1/icv83551/datasets
/content/drive/My Drive/icv83551/assignments/assignment1
```

1 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. 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.

```
[2]: # Run some setup code for this notebook.

import random
import numpy as np
from icv83551.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'

%%load_ext autoreload
%%autoreload 2
```

```
[3]: # Load the raw CIFAR-10 data.
cifar10_dir = 'icv83551/datasets/cifar-10-batches-py'

# Cleaning up variables to prevent loading data multiple times (which may cause
↳ memory issue)
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)
```

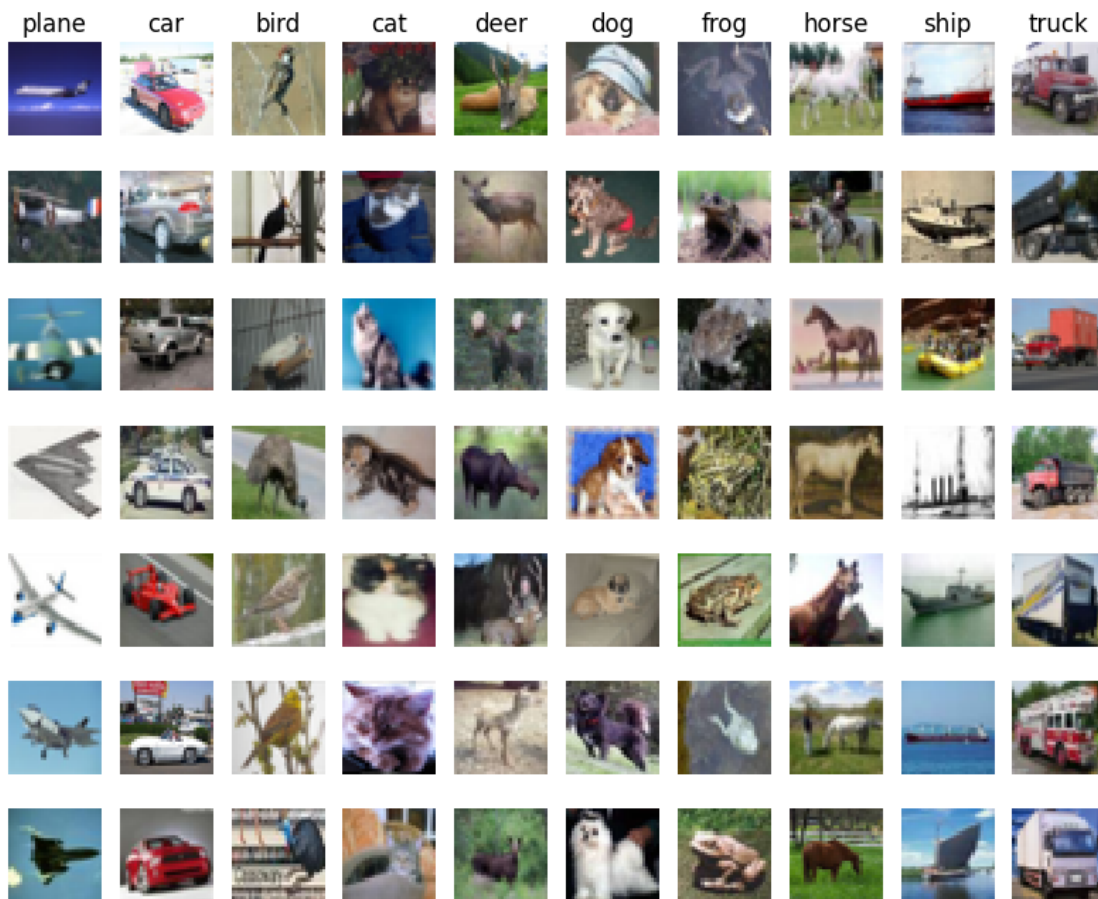
```
Training data shape: (50000, 32, 32, 3)
Training labels shape: (50000,)
Test data shape: (10000, 32, 32, 3)
Test labels shape: (10000,)
```

```
[4]: # Visualize some examples from the dataset.
# We show a few examples of training images from each class.
classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse',
↳ 'ship', 'truck']
```

```

num_classes = len(classes)
samples_per_class = 7
for y, cls in enumerate(classes):
    idxs = np.flatnonzero(y_train == y)
    idxs = np.random.choice(idxs, samples_per_class, replace=False)
    for i, idx in enumerate(idxs):
        plt_idx = i * num_classes + y + 1
        plt.subplot(samples_per_class, num_classes, plt_idx)
        plt.imshow(X_train[idx].astype('uint8'))
        plt.axis('off')
        if i == 0:
            plt.title(cls)
plt.show()

```



```

[5]: # 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)

```

[6]: from icv83551.classifiers import KNearestNeighbor

# Create a kNN classifier instance.
# Remember that training a kNN classifier is a noop:
# 'noop' stands for "no operation": in this context, training a kNN classifier
# does not perform any real learning step, it simply stores the labeled data
# so it can be used later to find nearest neighbors during prediction.

# 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 **N_{tr}** training examples and **N_{te}** test examples, this stage should result in a **N_{te} x N_{tr}** 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 `icv83551/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.

```

[7]: # Open icv83551/classifiers/k_nearest_neighbor.py and implement
# compute_distances_two_loops.

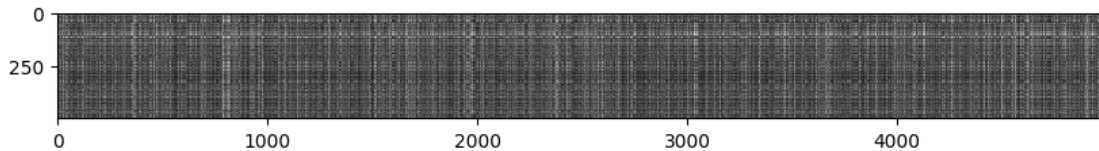
# Test your implementation:
dists = classifier.compute_distances_two_loops(X_test)

```

```
print(dists.shape)
```

(500, 5000)

```
[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 visibly 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 : The Y-axis are the test samples and the X-axis are the train samples, hence each row is the L2 distance of a single test sample and every train sample. Each column is the L2 distance of a single train sample and all the test samples. * Bright rows are caused when a test sample is very different from every train sample, because a row is the comparison of a test sample and train data, if there is a big difference (for example the test is a different class from most train data) the L2 distance of the pixels will be great for almost every column in the row and hence the color scheme and row will be bright.

- Bright columns are caused by a similar reason. When a train sample is very different from most of the test samples, the L2 distance will be great and most of the column will be bright.

```
[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:

```
[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. To clarify, both training and test examples are preprocessed in the same way.

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, which means rotating all the images by the same angle. Empty regions in the image caused by rotation are padded with a same pixel value and no interpolation is performed.

Your Answer : Preprocessing steps that will not change the performance: 1, 2, 3

Your Explanation : L1 is computed by $|p^{(a)} - p^{(b)}|$ where a, b are pixel indices.

1. The mean will cancel out each other because both pixels will be subtracted by the same constant, hence it will not affect the computation.
2. The mean will be canceled out, same as in question 1.
3. As we found subtracting by the mean does not affect the performance, dividing it by σ which is a constant won't affect the performance also because the constant for all the pixels is the same and all divided by the same value, so the comparison between pixels will be the same except the scale which will be scaled by σ .
4. Unlike question 3, here each distance is scaled by a different σ , so the scaling is different for all distances, meaning the comparison is compromised as different scalings can affect greatly the differences.

5. Because the image is not with a size $N \times N$, but $M \times N$, when we rotate we lose data (because previously filled data will be zero padded), and it will affect the distance between pixels.

```
[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

```
[12]: # 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

```
[13]: # 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 execute.
    """
    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
↪ implementation!

# 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 41.498916 seconds
One loop version took 55.363695 seconds
No loop version took 0.756936 seconds

1.0.1 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.

```

[14]: 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.
      #####
      X_train_folds = np.array_split(X_train, num_folds)
      y_train_folds = np.array_split(y_train, num_folds)

      # 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.
#####
for k in k_choices:
    k_to_accuracies[k] = []
    for i in range(num_folds):
        # concatenate all folds except fold i
        X_train_cv = np.concatenate([X_train_folds[j] for j in range(num_folds)
↪if j != i])
        y_train_cv = np.concatenate([y_train_folds[j] for j in range(num_folds)
↪if j != i])

        # Set fold i as the validation set
        X_val_cv = X_train_folds[i]
        y_val_cv = y_train_folds[i]

        # Train the classifier
        classifier = KNearestNeighbor()
        classifier.train(X_train_cv, y_train_cv)

        # Predict labels for the validation set with the no loops
        dists = classifier.compute_distances_no_loops(X_val_cv)
        y_val_pred = classifier.predict_labels(dists, k=k)

        # Calculate Accuracy
        num_correct = np.sum(y_val_pred == y_val_cv)
        accuracy = float(num_correct) / X_val_cv.shape[0]

        # Append to our results dictionary
        k_to_accuracies[k].append(accuracy)

# 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

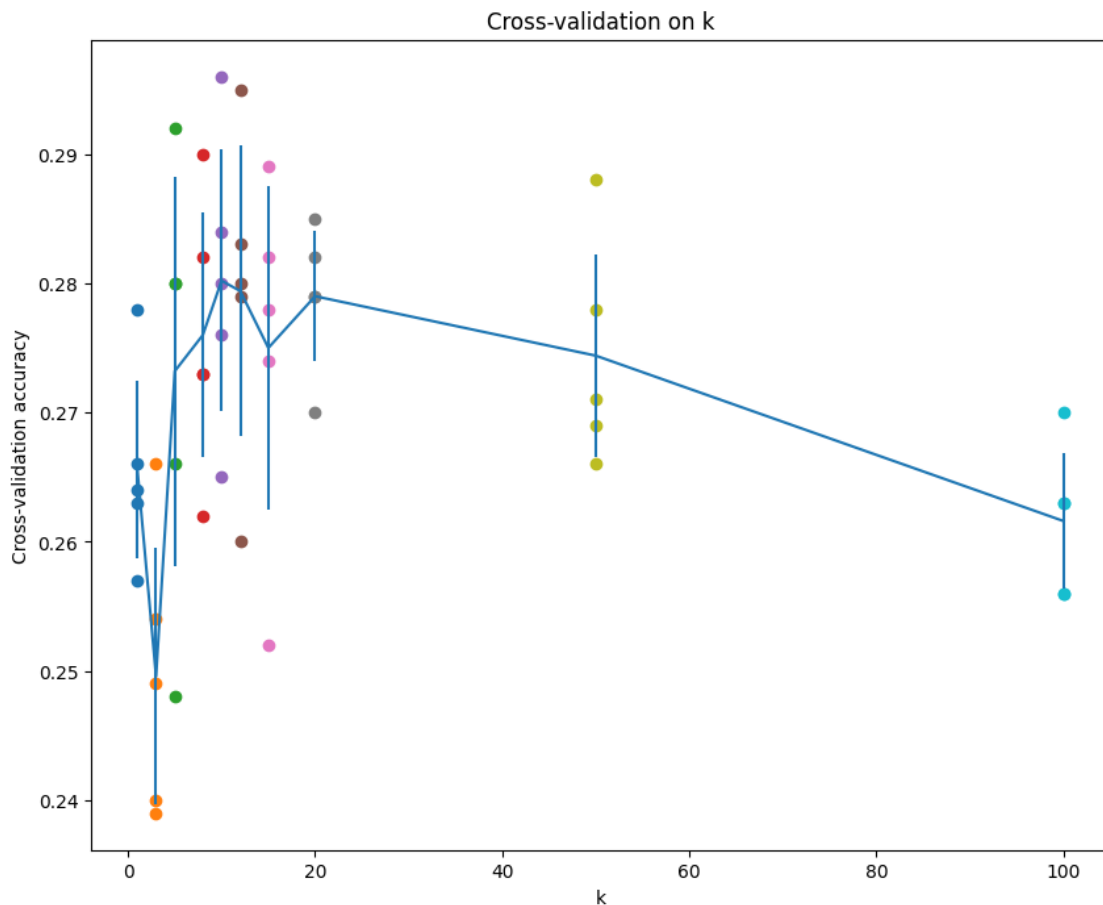
k = 1, accuracy = 0.257000

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

```
[15]: # 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()
```



```
[16]: # 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 = k_choices[accuracies_mean.argmax()]

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 or equal to 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 : True statements: 2, 4

Your Explanation :

1. FALSE: the decision boundary is not linear in n -dimensional space, the boundary is set by the closest neighbours and not by any sort of a hyperplane, hence it can be any complexed form and not linear.
2. TRUE: because, with 1-NN training sample it will predict its label based on the answer of the first closest image label, which is the given image. With 5-NN the training sample will base the answer on the first 5 closest images so there is a possibility to get the wrong label.
3. FALSE: 1-NN can misclassify an image during the test and 5-NN might get the correct class even if the closest image has a wrong label.
4. TRUE: each test sample is compared with every training sample during the test so when we increase the training set the test sample has more samples to be compared with, meaning longer comparison time.
5. False: 1,3 are false.

softmax

December 30, 2025

```
[1]: # This mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive')

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'icv83551/assignments/assignment1/'
FOLDERNAME = 'icv83551/assignments/assignment1/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

# Now that we've mounted your Drive, this ensures that
# the Python interpreter of the Colab VM can load
# python files from within it.
import sys
sys.path.append('/content/drive/My Drive/{}'.format(FOLDERNAME))

# This downloads the CIFAR-10 dataset to your Drive
# if it doesn't already exist.
%cd /content/drive/My\ Drive/$FOLDERNAME/icv83551/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

```
Mounted at /content/drive
/content/drive/My Drive/icv83551/assignments/assignment1/icv83551/datasets
/content/drive/My Drive/icv83551/assignments/assignment1
```

1 Softmax Classifier exercise

*Complete and hand in this completed worksheet (including its outputs and any supporting code outside of the worksheet) with your assignment submission.

In this 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** 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

```
[2]: # Run some setup code for this notebook.
import random
import numpy as np
from icv83551.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'

#%load_ext autoreload
#%autoreload 2
```

1.1 CIFAR-10 Data Loading and Preprocessing

```
[3]: # Load the raw CIFAR-10 data.
cifar10_dir = 'icv83551/datasets/cifar-10-batches-py'

# Cleaning up variables to prevent loading data multiple times (which may cause
↳memory issue)
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)
```

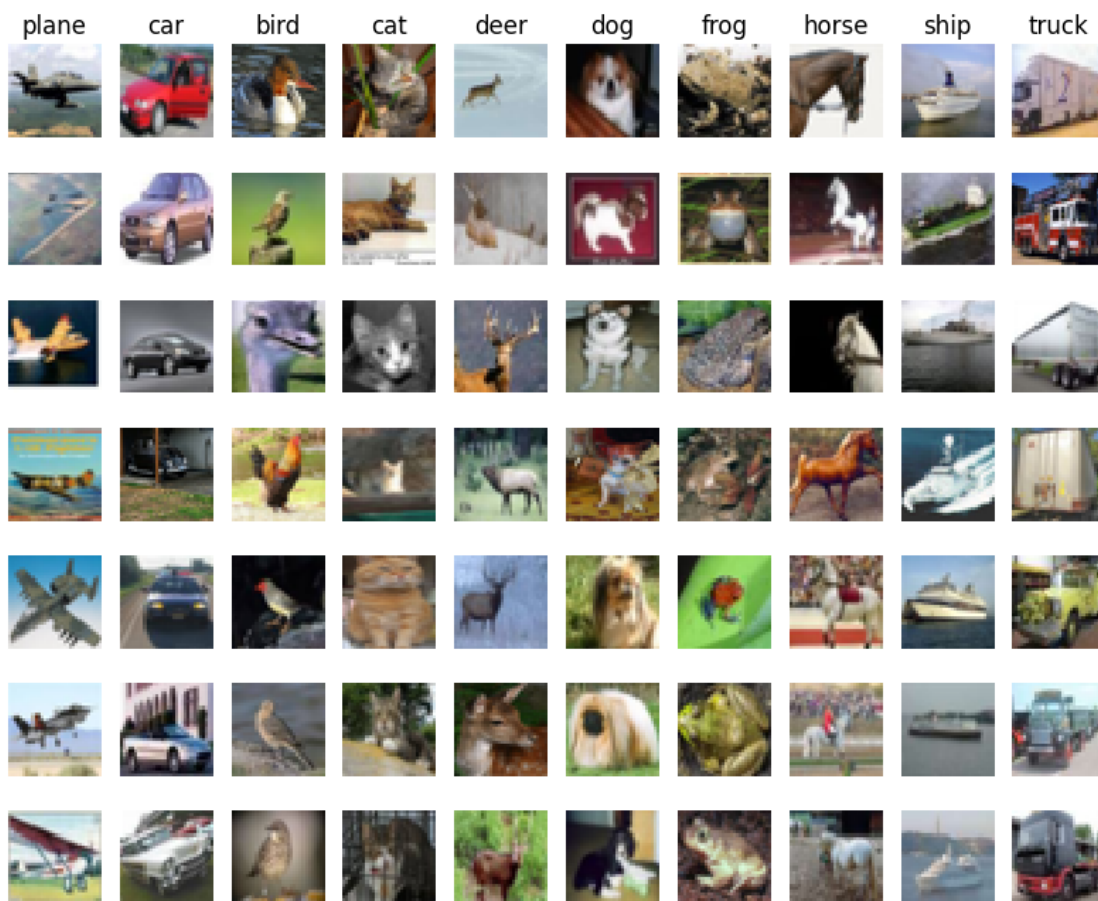
```
Training data shape: (50000, 32, 32, 3)
Training labels shape: (50000,)
Test data shape: (10000, 32, 32, 3)
Test labels shape: (10000,)
```

```
[4]: # Visualize some examples from the dataset.
# We show a few examples of training images from each class.
classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', '
↳ship', 'truck']
```

```

num_classes = len(classes)
samples_per_class = 7
for y, cls in enumerate(classes):
    idxs = np.flatnonzero(y_train == y)
    idxs = np.random.choice(idxs, samples_per_class, replace=False)
    for i, idx in enumerate(idxs):
        plt_idx = i * num_classes + y + 1
        plt.subplot(samples_per_class, num_classes, plt_idx)
        plt.imshow(X_train[idx].astype('uint8'))
        plt.axis('off')
        if i == 0:
            plt.title(cls)
plt.show()

```



```

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

```

```

num_validation = 1000
num_test = 1000
num_dev = 500

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

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

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

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

print('Train data shape: ', X_train.shape)
print('Train labels shape: ', y_train.shape)
print('Validation data shape: ', X_val.shape)
print('Validation labels shape: ', y_val.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)

```

```

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,)

```

```

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

```

```

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

```

```

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

```

```

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

# 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_
    ↪classifier
# 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))])

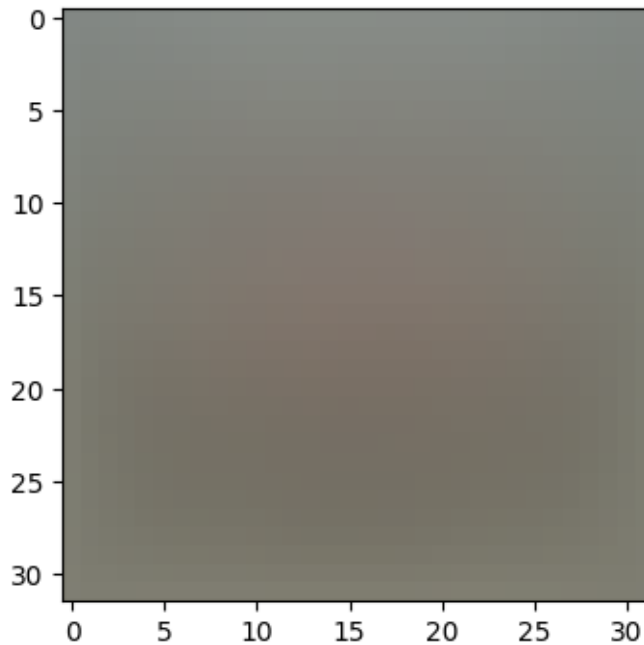
print(X_train.shape, X_val.shape, X_test.shape, X_dev.shape)

```

```

[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)

1.2 Softmax Classifier

Your code for this section will all be written inside `icv83551/classifiers/softmax.py`.

As you can see, we have prefilled the function `softmax_loss_naive` which uses for loops to evaluate the softmax loss function.

```
[8]: # Evaluate the naive implementation of the loss we provided for you:
from icv83551.classifiers.softmax import softmax_loss_naive
import time

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

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

# 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: 2.352047

loss: 2.352047

sanity check: 2.302585

Inline Question 1

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

Your Answer : Because if we have the same number of examples for each class, there is 10% probability (0.1 score) to classify the sample correctly if we guess randomly. The weights initialized randomly and naive softmax turns scores into probabilities, and the probability of a correct guess is 0.1.

The `grad` returned from the function above is right now all zero. Derive and implement the gradient for the softmax loss function and implement it inline inside the function `softmax_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, 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:

```
[9]: # 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 = softmax_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 icv83551.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)  
  
# do the gradient check once again with regularization turned on  
# you didn't forget the regularization gradient did you?  
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)
```

```
numerical: 0.824100 analytic: 0.824100, relative error: 1.613092e-09  
numerical: 0.087802 analytic: 0.087802, relative error: 3.895543e-07  
numerical: 1.469102 analytic: 1.469102, relative error: 4.505513e-08  
numerical: 0.228007 analytic: 0.228007, relative error: 1.984521e-07  
numerical: 0.710678 analytic: 0.710678, relative error: 2.062902e-08  
numerical: 0.014400 analytic: 0.014400, relative error: 2.505322e-06  
numerical: 1.901800 analytic: 1.901800, relative error: 1.361942e-08  
numerical: -0.785273 analytic: -0.785273, relative error: 1.047504e-07  
numerical: -2.373142 analytic: -2.373142, relative error: 1.348296e-08  
numerical: -2.916686 analytic: -2.916686, relative error: 2.524578e-09  
numerical: 0.824146 analytic: 0.824146, relative error: 1.608028e-09  
numerical: 0.461067 analytic: 0.461067, relative error: 4.441184e-08  
numerical: 0.178980 analytic: 0.178980, relative error: 4.831935e-07
```

```

numerical: -2.681250 analytic: -2.681250, relative error: 2.495805e-09
numerical: 1.879327 analytic: 1.879327, relative error: 2.152786e-08
numerical: 0.395466 analytic: 0.395466, relative error: 8.679832e-08
numerical: 3.218606 analytic: 3.218606, relative error: 8.061012e-09
numerical: 1.955164 analytic: 1.955164, relative error: 2.785039e-08
numerical: -1.240219 analytic: -1.240219, relative error: 5.490651e-09
numerical: -0.607289 analytic: -0.607289, relative error: 5.753195e-08

```

Inline Question 2

Although gradcheck is reliable softmax loss, it is possible that for SVM loss, 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 svm loss gradient check could fail? How would change the margin affect of the frequency of this happening?

Note that SVM loss for a sample (x_i, y_i) is defined as:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + \Delta)$$

Where j iterates over all classes except the correct class y_i and s_j denotes the classifier score for j^{th} class. Δ is a scalar margin. *Hint: the SVM loss function is not strictly speaking differentiable.*

Your Answer :

1. This discrepancy could be caused by the $\max(0, z)$ function, where $z = s_j - s_{y_i} + \Delta$. This function (Hinge Loss) is not differentiable at $z = 0$, so the analytical and numerical answers can be slightly different (z is less than the gradient step size).
2. It is not a reason for concern if both values are very close to each other and it isn't happening very often.
3. $f(z) = \max(0, z)$, $h = 0.00001$ where h is the gradient step size. let z be 0.000001 . The analytical gradient is $f'(z) = 1$ because $z > 0$. But the Numerical gradient will be

$$\frac{f(z+h) - f(z-h)}{2h} = \frac{\max(0, 0.000001 + 0.00001) - \max(0, 0.000001 - 0.00001)}{2 * 0.00001} = \frac{\max(0, 0.000011) - \max(0, -0.000009)}{0.00002}$$

As we can see $1 \neq 0.55$.

4. The change of margin Δ won't affect the frequency of that happening, only the location because it happens when $s_j - s_{y_i} = -\Delta$, the frequency does not depend on the value.

```

[10]: # Next implement the function softmax_loss_vectorized; for now only compute the
      ↪ loss;
      # we will implement the gradient in a moment.
      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 icv83551.classifiers.softmax import softmax_loss_vectorized
      tic = time.time()

```

```

loss_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))

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

```

Naive loss: 2.352047e+00 computed in 0.071032s
 Vectorized loss: 2.352047e+00 computed in 0.011881s
 difference: -0.000000

```

[11]: # Complete the implementation of softmax_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 = softmax_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 = softmax_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 0.050906s
 Vectorized loss and gradient: computed in 0.006796s
 difference: 0.000000

1.2.1 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 `icv83551/classifiers/linear_classifier.py`.

```

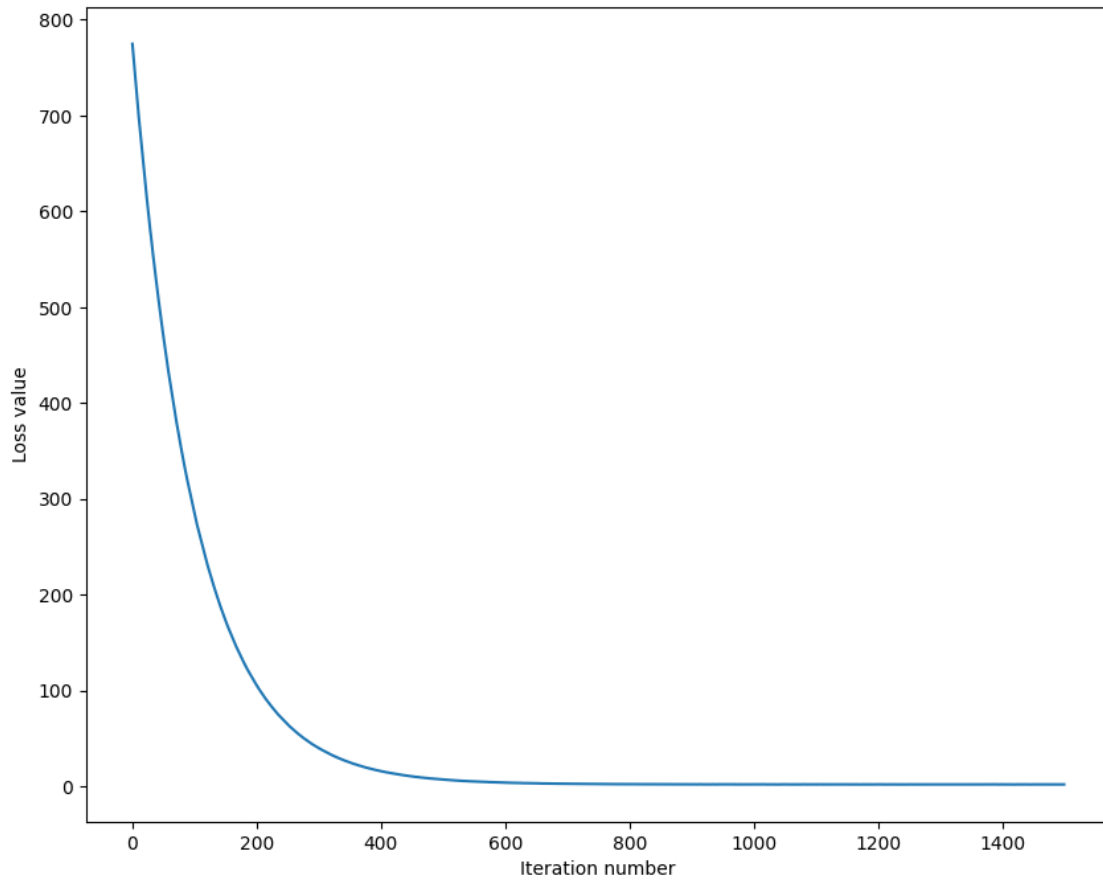
[12]: # In the file linear_classifier.py, implement SGD in the function
      # LinearClassifier.train() and then run it with the code below.
      from icv83551.classifiers import Softmax
      softmax = Softmax()

```

```
tic = time.time()
loss_hist = softmax.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 774.549019
iteration 100 / 1500: loss 284.496040
iteration 200 / 1500: loss 105.595929
iteration 300 / 1500: loss 39.879496
iteration 400 / 1500: loss 15.907653
iteration 500 / 1500: loss 7.217231
iteration 600 / 1500: loss 3.921068
iteration 700 / 1500: loss 2.776863
iteration 800 / 1500: loss 2.311333
iteration 900 / 1500: loss 2.206896
iteration 1000 / 1500: loss 2.136169
iteration 1100 / 1500: loss 2.110088
iteration 1200 / 1500: loss 2.055868
iteration 1300 / 1500: loss 2.093111
iteration 1400 / 1500: loss 2.128530
That took 5.128370s
```

```
[13]: # 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()
```



```
[14]: # Write the LinearClassifier.predict function and evaluate the performance on
# both the training and validation set
# You should get validation accuracy of about 0.34 (> 0.33).
y_train_pred = softmax.predict(X_train)
print(y_train_pred)
print('training accuracy: %f' % (np.mean(y_train == y_train_pred), ))
y_val_pred = softmax.predict(X_val)
print('validation accuracy: %f' % (np.mean(y_val == y_val_pred), ))
```

```
[6 8 9 ... 5 6 8]
training accuracy: 0.329265
validation accuracy: 0.351000
```

```
[15]: # Save the trained model for autograder.
softmax.save("softmax.npy")
```

```
softmax.npy saved.
```

```
[46]: # 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 0.365 (> 0.36) 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_softmax = None # The Softmax 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 Softmax 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 Softmax object that achieves this. #
# accuracy in best_softmax. #
# #
# Hint: You should use a small value for num_iters as you develop your #
# validation code so that the classifiers don't take much time to train; once #
# you are confident that your validation code works, you should rerun the #
# 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-6, 1e-5]
regularization_strengths = [1e4, 2.5e3]

import itertools

for lr, reg in itertools.product(learning_rates, regularization_strengths):
    print(f'lr = {lr} reg = {reg}')
    # Softmax classifier
    softmax = Softmax()
    softmax.train(X_train, y_train, lr, reg, num_iters=1000)

    # Softmax predictions
    y_train_pred = softmax.predict(X_train)
```

```

y_val_pred = softmax.predict(X_val)
results[(lr, reg)] = np.mean(y_train == y_train_pred), np.mean(y_val ==
↪y_val_pred)

# Save if validation accuracy is the best
if results[(lr, reg)][1] > best_val:
    best_val = results[(lr, reg)][1]
    best_softmax = softmax

# 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 = 1e-06 reg = 10000.0
lr = 1e-06 reg = 2500.0
lr = 1e-05 reg = 10000.0
lr = 1e-05 reg = 2500.0
lr 1.000000e-06 reg 2.500000e+03 train accuracy: 0.382102 val accuracy: 0.384000
lr 1.000000e-06 reg 1.000000e+04 train accuracy: 0.352612 val accuracy: 0.349000
lr 1.000000e-05 reg 2.500000e+03 train accuracy: 0.210612 val accuracy: 0.217000
lr 1.000000e-05 reg 1.000000e+04 train accuracy: 0.190673 val accuracy: 0.197000
best validation accuracy achieved during cross-validation: 0.384000

```

```

[47]: # Visualize the cross-validation results
import math
import pdb

# pdb.set_trace()

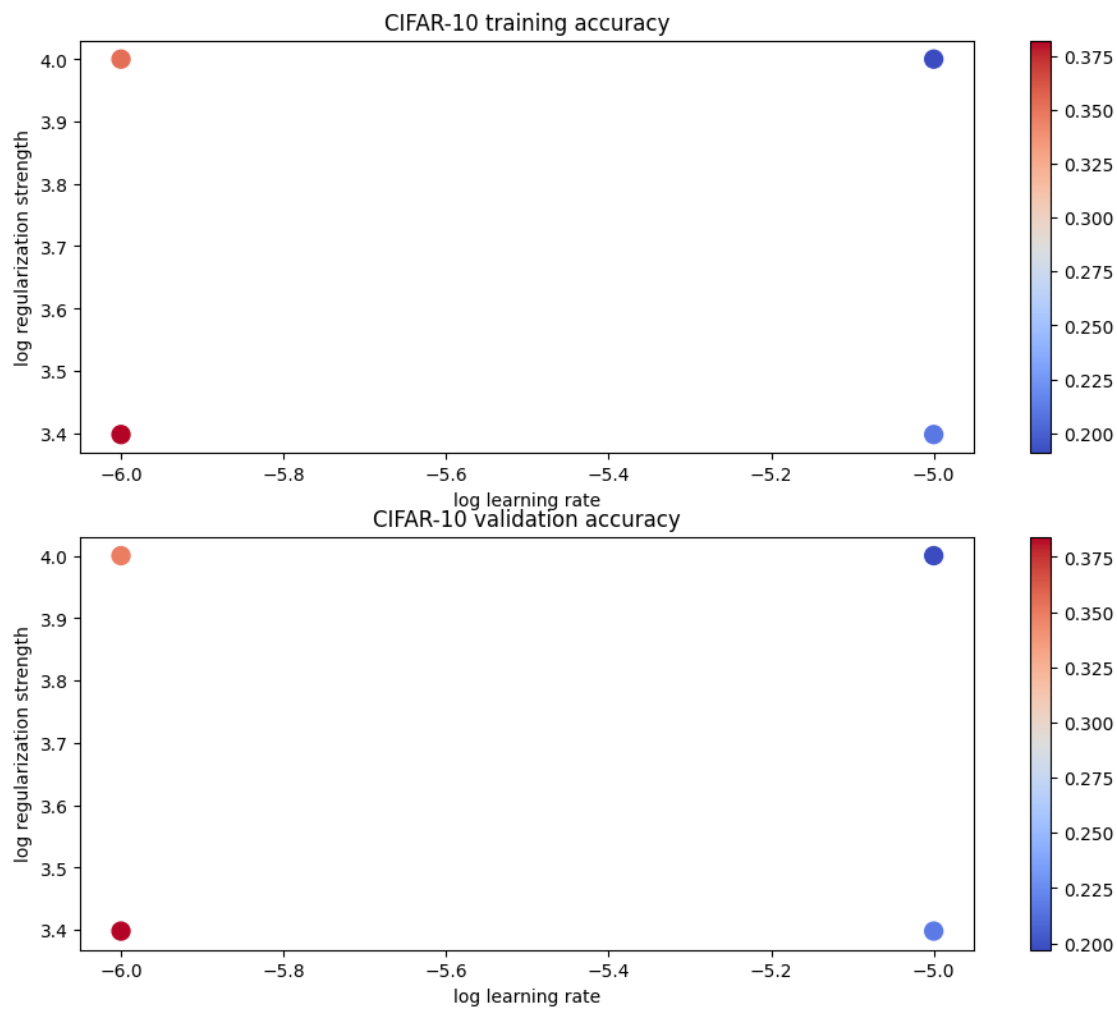
x_scatter = [math.log10(x[0]) for x in results]
y_scatter = [math.log10(x[1]) for x in results]

# 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()
```



```
[48]: # 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 classifier on raw pixels final test set accuracy: %f' %  
      ↪test_accuracy)
```

Softmax classifier on raw pixels final test set accuracy: 0.367000

```
[49]: # Save best softmax model  
best_softmax.save("best_softmax.npy")
```

best_softmax.npy saved.

```
[50]: # Visualize the learned weights for each class.  
# Depending on your choice of learning rate and regularization strength, these  
      ↪may  
# or may not be nice to look at.  
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', 'truck']  
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 3

Describe what your visualized Softmax classifier weights look like, and offer a brief explanation for why they look the way they do.

Your Answer : The visualized weights looks like a very noisy version of each class. This is because the Softmax classifier is a linear transformation and averaging of the training images. It is trying to template match, meaning the classifier checks how well the input image aligns with the weight matrix which was calculated by linear calculations and averaging of different examples (colors, background, etc..) hence a lot of noise in the process.

Inline Question 4 - 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 change the softmax loss, but leave the SVM loss unchanged.

Your Answer : True

Your Explanation : It is possible because of the margin of SVM. For SVM if the new training point satisfy the margin Δ (meaning the correct class wins by a margin) then the loss would be 0. This is not the case for the Softmax loss, as it can never be 0 (Softmax is a probability and it can never be exactly 1 or 0, meaning there will always be loss).

[]:

two_layer_net

December 30, 2025

```
[1]: # This mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive')

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'icv83551/assignments/assignment1/'
FOLDERNAME = 'icv83551/assignments/assignment1/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

# Now that we've mounted your Drive, this ensures that
# the Python interpreter of the Colab VM can load
# python files from within it.
import sys
sys.path.append('/content/drive/My Drive/{}'.format(FOLDERNAME))

# This downloads the CIFAR-10 dataset to your Drive
# if it doesn't already exist.
%cd /content/drive/My\ Drive/$FOLDERNAME/icv83551/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

```
Mounted at /content/drive
/content/drive/My Drive/icv83551/assignments/assignment1/icv83551/datasets
/content/drive/My Drive/icv83551/assignments/assignment1
```

1 Fully-Connected Neural Nets

In this exercise we will implement fully-connected networks using a modular approach. For each layer we will implement a **forward** and a **backward** function. The **forward** function will receive inputs, weights, and other parameters and will return both an output and a **cache** object storing data needed for the backward pass, like this:

```
def layer_forward(x, w):
    """ Receive inputs x and weights w """
    # Do some computations ...
    z = # ... some intermediate value
    # Do some more computations ...
    out = # the output
```

```
cache = (x, w, z, out) # Values we need to compute gradients
```

```
return out, cache
```

The backward pass will receive upstream derivatives and the `cache` object, and will return gradients with respect to the inputs and weights, like this:

```
def layer_backward(dout, cache):
    """
    Receive dout (derivative of loss with respect to outputs) and cache,
    and compute derivative with respect to inputs.
    """
    # Unpack cache values
    x, w, z, out = cache

    # Use values in cache to compute derivatives
    dx = # Derivative of loss with respect to x
    dw = # Derivative of loss with respect to w

    return dx, dw
```

After implementing a bunch of layers this way, we will be able to easily combine them to build classifiers with different architectures.

```
[2]: # As usual, a bit of setup
from __future__ import print_function
import time
import numpy as np
import matplotlib.pyplot as plt
from icv83551.classifiers.fc_net import *
from icv83551.data_utils import get_CIFAR10_data
from icv83551.gradient_check import eval_numerical_gradient, \
    eval_numerical_gradient_array
from icv83551.solver import Solver

%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'

%%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))))
```

```
[3]: # Load the (preprocessed) CIFAR10 data.
```

```
data = get_CIFAR10_data()
for k, v in list(data.items()):
    print('%s: ' % k, v.shape)
```

```
('X_train: ', (49000, 3, 32, 32))
('y_train: ', (49000,))
('X_val: ', (1000, 3, 32, 32))
('y_val: ', (1000,))
('X_test: ', (1000, 3, 32, 32))
('y_test: ', (1000,))
```

2 Affine layer: forward

Open the file `icv83551/layers.py` and implement the `affine_forward` function.

Once you are done you can test your implementation by running the following:

```
[4]: # Test the affine_forward function
```

```
num_inputs = 2
input_shape = (4, 5, 6)
output_dim = 3

input_size = num_inputs * np.prod(input_shape)
weight_size = output_dim * np.prod(input_shape)

x = np.linspace(-0.1, 0.5, num=input_size).reshape(num_inputs, *input_shape)
w = np.linspace(-0.2, 0.3, num=weight_size).reshape(np.prod(input_shape),
↳output_dim)
b = np.linspace(-0.3, 0.1, num=output_dim)

out, _ = affine_forward(x, w, b)
correct_out = np.array([[ 1.49834967,  1.70660132,  1.91485297],
                        [ 3.25553199,  3.5141327,   3.77273342]])

# Compare your output with ours. The error should be around e-9 or less.
print('Testing affine_forward function:')
print('difference: ', rel_error(out, correct_out))
```

```
Testing affine_forward function:
difference: 9.769849468192957e-10
```

3 Affine layer: backward

Now implement the `affine_backward` function and test your implementation using numeric gradient checking.

```
[5]: # Test the affine_backward function
np.random.seed(231)
x = np.random.randn(10, 2, 3)
w = np.random.randn(6, 5)
b = np.random.randn(5)
dout = np.random.randn(10, 5)

dx_num = eval_numerical_gradient_array(lambda x: affine_forward(x, w, b)[0], x,
    ↪dout)
dw_num = eval_numerical_gradient_array(lambda w: affine_forward(x, w, b)[0], w,
    ↪dout)
db_num = eval_numerical_gradient_array(lambda b: affine_forward(x, w, b)[0], b,
    ↪dout)

_, cache = affine_forward(x, w, b)
dx, dw, db = affine_backward(dout, cache)

# The error should be around e-10 or less
print('Testing affine_backward function:')
print('dx error: ', rel_error(dx_num, dx))
print('dw error: ', rel_error(dw_num, dw))
print('db error: ', rel_error(db_num, db))
```

```
Testing affine_backward function:
dx error:  5.399100368651805e-11
dw error:  9.904211865398145e-11
db error:  2.4122867568119087e-11
```

4 ReLU activation: forward

Implement the forward pass for the ReLU activation function in the `relu_forward` function and test your implementation using the following:

```
[6]: # Test the relu_forward function

x = np.linspace(-0.5, 0.5, num=12).reshape(3, 4)

out, _ = relu_forward(x)
correct_out = np.array([[ 0.,          0.,          0.,          0.],
                        [ 0.,          0.,          0.04545455, 0.13636364],
                        [ 0.22727273, 0.31818182, 0.40909091, 0.5]])

# Compare your output with ours. The error should be on the order of e-8
print('Testing relu_forward function:')
print('difference: ', rel_error(out, correct_out))
```

```
Testing relu_forward function:
difference: 4.999999798022158e-08
```

5 ReLU activation: backward

Now implement the backward pass for the ReLU activation function in the `relu_backward` function and test your implementation using numeric gradient checking:

```
[7]: np.random.seed(231)
x = np.random.randn(10, 10)
dout = np.random.randn(*x.shape)

dx_num = eval_numerical_gradient_array(lambda x: relu_forward(x)[0], x, dout)

_, cache = relu_forward(x)
dx = relu_backward(dout, cache)

# The error should be on the order of e-12
print('Testing relu_backward function:')
print('dx error: ', rel_error(dx_num, dx))
```

```
Testing relu_backward function:
dx error: 3.2756349136310288e-12
```

5.1 Inline Question 1:

We've only asked you to implement ReLU, but there are a number of different activation functions that one could use in neural networks, each with its pros and cons. In particular, an issue commonly seen with activation functions is getting zero (or close to zero) gradient flow during backpropagation. Which of the following activation functions have this problem? If you consider these functions in the one dimensional case, what types of input would lead to this behaviour? 1. Sigmoid 2. ReLU 3. Leaky ReLU

Your Answer : Functions 1 and 2 have this problem.

Function 1 - the Sigmoid: when a value is very large or very small the probability is almost 1 or 0, with very little to no change and hence the gradient becomes 0 (or very very close).

Function 2 - the ReLU: if all the inputs to the ReLU function are negative, then $\text{ReLU} = 0$ and there will be no gradient.

6 “Sandwich” layers

There are some common patterns of layers that are frequently used in neural nets. For example, affine layers are frequently followed by a ReLU nonlinearity. To make these common patterns easy, we define several convenience layers in the file `icv83551/layer_utils.py`.

For now take a look at the `affine_relu_forward` and `affine_relu_backward` functions, and run the following to numerically gradient check the backward pass:

```
[8]: from icv83551.layer_utils import affine_relu_forward, affine_relu_backward
np.random.seed(231)
x = np.random.randn(2, 3, 4)
w = np.random.randn(12, 10)
b = np.random.randn(10)
dout = np.random.randn(2, 10)

out, cache = affine_relu_forward(x, w, b)
dx, dw, db = affine_relu_backward(dout, cache)

dx_num = eval_numerical_gradient_array(lambda x: affine_relu_forward(x, w,
    ↪b)[0], x, dout)
dw_num = eval_numerical_gradient_array(lambda w: affine_relu_forward(x, w,
    ↪b)[0], w, dout)
db_num = eval_numerical_gradient_array(lambda b: affine_relu_forward(x, w,
    ↪b)[0], b, dout)

# Relative error should be around e-10 or less
print('Testing affine_relu_forward and affine_relu_backward:')
print('dx error: ', rel_error(dx_num, dx))
print('dw error: ', rel_error(dw_num, dw))
print('db error: ', rel_error(db_num, db))
```

```
Testing affine_relu_forward and affine_relu_backward:
dx error:  2.299579177309368e-11
dw error:  8.162011105764925e-11
db error:  7.826724021458994e-12
```

7 Loss layers: Softmax

Now implement the loss and gradient for softmax in the `softmax_loss` function in `icv83551/layers.py`. These should be similar to what you implemented in `icv83551/classifiers/softmax.py`. Other loss functions (e.g. `svm_loss`) can also be implemented in a modular way, however, it is not required for this assignment.

You can make sure that the implementations are correct by running the following:

```
[9]: np.random.seed(231)
num_classes, num_inputs = 10, 50
x = 0.001 * np.random.randn(num_inputs, num_classes)
y = np.random.randint(num_classes, size=num_inputs)

dx_num = eval_numerical_gradient(lambda x: softmax_loss(x, y)[0], x,
    ↪verbose=False)
loss, dx = softmax_loss(x, y)
```

```
# Test softmax_loss function. Loss should be close to 2.3 and dx error should
    ↪ be around e-8
print('\nTesting softmax_loss:')
print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))
```

```
Testing softmax_loss:
loss: 2.3025458445007376
dx error: 8.234144091578429e-09
```

8 Two-layer network

Open the file `icv83551/classifiers/fc_net.py` and complete the implementation of the `TwoLayerNet` class. Read through it to make sure you understand the API. You can run the cell below to test your implementation.

```
[10]: np.random.seed(231)
N, D, H, C = 3, 5, 50, 7
X = np.random.randn(N, D)
y = np.random.randint(C, size=N)

std = 1e-3
model = TwoLayerNet(input_dim=D, hidden_dim=H, num_classes=C, weight_scale=std)

print('Testing initialization ... ')
W1_std = abs(model.params['W1'].std() - std)
b1 = model.params['b1']
W2_std = abs(model.params['W2'].std() - std)
b2 = model.params['b2']
assert W1_std < std / 10, 'First layer weights do not seem right'
assert np.all(b1 == 0), 'First layer biases do not seem right'
assert W2_std < std / 10, 'Second layer weights do not seem right'
assert np.all(b2 == 0), 'Second layer biases do not seem right'

print('Testing test-time forward pass ... ')
model.params['W1'] = np.linspace(-0.7, 0.3, num=D*H).reshape(D, H)
model.params['b1'] = np.linspace(-0.1, 0.9, num=H)
model.params['W2'] = np.linspace(-0.3, 0.4, num=H*C).reshape(H, C)
model.params['b2'] = np.linspace(-0.9, 0.1, num=C)
X = np.linspace(-5.5, 4.5, num=N*D).reshape(D, N).T
scores = model.loss(X)
correct_scores = np.asarray(
    [[11.53165108, 12.2917344, 13.05181771, 13.81190102, 14.57198434, 15.
    ↪ 33206765, 16.09215096],
    [12.05769098, 12.74614105, 13.43459113, 14.1230412, 14.81149128, 15.
    ↪ 49994135, 16.18839143],
```

```

[12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.05099822, 15.
↪66781506, 16.2846319 ]])
scores_diff = np.abs(scores - correct_scores).sum()
assert scores_diff < 1e-6, 'Problem with test-time forward pass'

print('Testing training loss (no regularization)')
y = np.asarray([0, 5, 1])
loss, grads = model.loss(X, y)
correct_loss = 3.4702243556
assert abs(loss - correct_loss) < 1e-10, 'Problem with training-time loss'

model.reg = 1.0
loss, grads = model.loss(X, y)
correct_loss = 26.5948426952
assert abs(loss - correct_loss) < 1e-10, 'Problem with regularization loss'

# Errors should be around e-7 or less
for reg in [0.0, 0.7]:
    print('Running numeric gradient check with reg = ', reg)
    model.reg = reg
    loss, grads = model.loss(X, y)

    for name in sorted(grads):
        f = lambda _: model.loss(X, y)[0]
        grad_num = eval_numerical_gradient(f, model.params[name], verbose=False)
        print('%s relative error: %.2e' % (name, rel_error(grad_num, grads[name])))

```

```

Testing initialization ...
Testing test-time forward pass ...
Testing training loss (no regularization)
Running numeric gradient check with reg = 0.0
W1 relative error: 1.83e-08
W2 relative error: 3.20e-10
b1 relative error: 9.83e-09
b2 relative error: 4.33e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 2.53e-07
W2 relative error: 2.85e-08
b1 relative error: 1.56e-08
b2 relative error: 9.09e-10

```

9 Solver

Open the file `icv83551/solver.py` and read through it to familiarize yourself with the API. After doing so, use a `Solver` instance to train a `TwoLayerNet` that achieves about 36% accuracy on the validation set.

```
[11]: input_size = 32 * 32 * 3
hidden_size = 50
num_classes = 10
model = TwoLayerNet(input_size, hidden_size, num_classes)
solver = None

#####
# TODO: Use a Solver instance to train a TwoLayerNet that achieves about 36% #
# accuracy on the validation set.                                           #
#####
solver = Solver(model, data, optim_config={'learning_rate': 1e-4})
solver.train()
#####
#                               END OF YOUR CODE                               #
#####
```

```
(Iteration 1 / 4900) loss: 2.300089
(Epoch 0 / 10) train acc: 0.138000; val_acc: 0.138000
(Iteration 11 / 4900) loss: 2.297774
(Iteration 21 / 4900) loss: 2.297181
(Iteration 31 / 4900) loss: 2.292076
(Iteration 41 / 4900) loss: 2.295523
(Iteration 51 / 4900) loss: 2.289783
(Iteration 61 / 4900) loss: 2.276536
(Iteration 71 / 4900) loss: 2.265533
(Iteration 81 / 4900) loss: 2.275771
(Iteration 91 / 4900) loss: 2.261915
(Iteration 101 / 4900) loss: 2.240327
(Iteration 111 / 4900) loss: 2.241637
(Iteration 121 / 4900) loss: 2.236025
(Iteration 131 / 4900) loss: 2.222546
(Iteration 141 / 4900) loss: 2.210281
(Iteration 151 / 4900) loss: 2.196274
(Iteration 161 / 4900) loss: 2.153850
(Iteration 171 / 4900) loss: 2.185944
(Iteration 181 / 4900) loss: 2.177319
(Iteration 191 / 4900) loss: 2.173674
(Iteration 201 / 4900) loss: 2.113574
(Iteration 211 / 4900) loss: 2.084054
(Iteration 221 / 4900) loss: 2.153114
(Iteration 231 / 4900) loss: 2.098817
(Iteration 241 / 4900) loss: 2.054590
(Iteration 251 / 4900) loss: 2.023479
(Iteration 261 / 4900) loss: 2.178227
(Iteration 271 / 4900) loss: 2.057826
(Iteration 281 / 4900) loss: 2.035345
(Iteration 291 / 4900) loss: 2.099167
```

(Iteration 301 / 4900) loss: 2.026074
(Iteration 311 / 4900) loss: 2.108364
(Iteration 321 / 4900) loss: 2.154241
(Iteration 331 / 4900) loss: 1.991962
(Iteration 341 / 4900) loss: 2.045785
(Iteration 351 / 4900) loss: 1.987154
(Iteration 361 / 4900) loss: 1.991075
(Iteration 371 / 4900) loss: 2.134392
(Iteration 381 / 4900) loss: 2.123832
(Iteration 391 / 4900) loss: 2.019381
(Iteration 401 / 4900) loss: 2.048385
(Iteration 411 / 4900) loss: 1.986564
(Iteration 421 / 4900) loss: 1.915372
(Iteration 431 / 4900) loss: 1.949682
(Iteration 441 / 4900) loss: 1.925332
(Iteration 451 / 4900) loss: 1.995748
(Iteration 461 / 4900) loss: 1.935236
(Iteration 471 / 4900) loss: 1.963299
(Iteration 481 / 4900) loss: 2.045542
(Epoch 1 / 10) train acc: 0.297000; val_acc: 0.300000
(Iteration 491 / 4900) loss: 1.990057
(Iteration 501 / 4900) loss: 1.978907
(Iteration 511 / 4900) loss: 1.997122
(Iteration 521 / 4900) loss: 1.994085
(Iteration 531 / 4900) loss: 1.970249
(Iteration 541 / 4900) loss: 2.042874
(Iteration 551 / 4900) loss: 1.995102
(Iteration 561 / 4900) loss: 1.889537
(Iteration 571 / 4900) loss: 2.021225
(Iteration 581 / 4900) loss: 1.938430
(Iteration 591 / 4900) loss: 1.870129
(Iteration 601 / 4900) loss: 1.832621
(Iteration 611 / 4900) loss: 1.774099
(Iteration 621 / 4900) loss: 1.861298
(Iteration 631 / 4900) loss: 1.840973
(Iteration 641 / 4900) loss: 1.934974
(Iteration 651 / 4900) loss: 1.794739
(Iteration 661 / 4900) loss: 1.889135
(Iteration 671 / 4900) loss: 1.870638
(Iteration 681 / 4900) loss: 1.877747
(Iteration 691 / 4900) loss: 1.904023
(Iteration 701 / 4900) loss: 1.862299
(Iteration 711 / 4900) loss: 1.975274
(Iteration 721 / 4900) loss: 1.866224
(Iteration 731 / 4900) loss: 1.944338
(Iteration 741 / 4900) loss: 1.920524
(Iteration 751 / 4900) loss: 1.765448
(Iteration 761 / 4900) loss: 1.969160

(Iteration 771 / 4900) loss: 1.970204
(Iteration 781 / 4900) loss: 1.883607
(Iteration 791 / 4900) loss: 1.757966
(Iteration 801 / 4900) loss: 1.788405
(Iteration 811 / 4900) loss: 2.006606
(Iteration 821 / 4900) loss: 1.872685
(Iteration 831 / 4900) loss: 1.824610
(Iteration 841 / 4900) loss: 1.881702
(Iteration 851 / 4900) loss: 1.852286
(Iteration 861 / 4900) loss: 1.775567
(Iteration 871 / 4900) loss: 1.726155
(Iteration 881 / 4900) loss: 1.927214
(Iteration 891 / 4900) loss: 2.045421
(Iteration 901 / 4900) loss: 1.763286
(Iteration 911 / 4900) loss: 1.827155
(Iteration 921 / 4900) loss: 1.768191
(Iteration 931 / 4900) loss: 1.844635
(Iteration 941 / 4900) loss: 1.810668
(Iteration 951 / 4900) loss: 1.833363
(Iteration 961 / 4900) loss: 1.839569
(Iteration 971 / 4900) loss: 1.652194
(Epoch 2 / 10) train acc: 0.358000; val_acc: 0.359000
(Iteration 981 / 4900) loss: 1.924723
(Iteration 991 / 4900) loss: 1.848202
(Iteration 1001 / 4900) loss: 1.712420
(Iteration 1011 / 4900) loss: 1.714966
(Iteration 1021 / 4900) loss: 1.643041
(Iteration 1031 / 4900) loss: 1.817483
(Iteration 1041 / 4900) loss: 1.958801
(Iteration 1051 / 4900) loss: 1.869896
(Iteration 1061 / 4900) loss: 1.800445
(Iteration 1071 / 4900) loss: 1.669287
(Iteration 1081 / 4900) loss: 1.769503
(Iteration 1091 / 4900) loss: 1.722780
(Iteration 1101 / 4900) loss: 1.733109
(Iteration 1111 / 4900) loss: 1.965675
(Iteration 1121 / 4900) loss: 1.748992
(Iteration 1131 / 4900) loss: 1.728268
(Iteration 1141 / 4900) loss: 1.842387
(Iteration 1151 / 4900) loss: 1.630562
(Iteration 1161 / 4900) loss: 1.696379
(Iteration 1171 / 4900) loss: 1.816942
(Iteration 1181 / 4900) loss: 1.750974
(Iteration 1191 / 4900) loss: 1.777530
(Iteration 1201 / 4900) loss: 1.747407
(Iteration 1211 / 4900) loss: 1.813170
(Iteration 1221 / 4900) loss: 1.778816
(Iteration 1231 / 4900) loss: 1.687309

(Iteration 1241 / 4900) loss: 1.722979
(Iteration 1251 / 4900) loss: 1.610165
(Iteration 1261 / 4900) loss: 1.857974
(Iteration 1271 / 4900) loss: 1.815787
(Iteration 1281 / 4900) loss: 1.698776
(Iteration 1291 / 4900) loss: 1.597460
(Iteration 1301 / 4900) loss: 1.788008
(Iteration 1311 / 4900) loss: 1.691238
(Iteration 1321 / 4900) loss: 1.673124
(Iteration 1331 / 4900) loss: 1.839491
(Iteration 1341 / 4900) loss: 1.727903
(Iteration 1351 / 4900) loss: 1.845529
(Iteration 1361 / 4900) loss: 1.784334
(Iteration 1371 / 4900) loss: 1.599927
(Iteration 1381 / 4900) loss: 1.885755
(Iteration 1391 / 4900) loss: 1.832098
(Iteration 1401 / 4900) loss: 1.785752
(Iteration 1411 / 4900) loss: 1.791261
(Iteration 1421 / 4900) loss: 1.806974
(Iteration 1431 / 4900) loss: 1.767754
(Iteration 1441 / 4900) loss: 1.678441
(Iteration 1451 / 4900) loss: 1.600273
(Iteration 1461 / 4900) loss: 1.740208
(Epoch 3 / 10) train acc: 0.404000; val_acc: 0.383000
(Iteration 1471 / 4900) loss: 1.706930
(Iteration 1481 / 4900) loss: 1.678862
(Iteration 1491 / 4900) loss: 1.726915
(Iteration 1501 / 4900) loss: 1.707665
(Iteration 1511 / 4900) loss: 1.671057
(Iteration 1521 / 4900) loss: 1.693702
(Iteration 1531 / 4900) loss: 1.800129
(Iteration 1541 / 4900) loss: 1.709060
(Iteration 1551 / 4900) loss: 1.883944
(Iteration 1561 / 4900) loss: 1.565166
(Iteration 1571 / 4900) loss: 1.716219
(Iteration 1581 / 4900) loss: 1.643147
(Iteration 1591 / 4900) loss: 1.689360
(Iteration 1601 / 4900) loss: 1.750729
(Iteration 1611 / 4900) loss: 1.670908
(Iteration 1621 / 4900) loss: 1.810203
(Iteration 1631 / 4900) loss: 1.707229
(Iteration 1641 / 4900) loss: 1.638586
(Iteration 1651 / 4900) loss: 1.882884
(Iteration 1661 / 4900) loss: 1.559042
(Iteration 1671 / 4900) loss: 1.704065
(Iteration 1681 / 4900) loss: 1.631521
(Iteration 1691 / 4900) loss: 1.668143
(Iteration 1701 / 4900) loss: 1.697952

(Iteration 1711 / 4900) loss: 1.700783
(Iteration 1721 / 4900) loss: 1.655203
(Iteration 1731 / 4900) loss: 1.610352
(Iteration 1741 / 4900) loss: 1.712151
(Iteration 1751 / 4900) loss: 1.751529
(Iteration 1761 / 4900) loss: 1.701241
(Iteration 1771 / 4900) loss: 1.782908
(Iteration 1781 / 4900) loss: 1.748410
(Iteration 1791 / 4900) loss: 1.628487
(Iteration 1801 / 4900) loss: 1.675512
(Iteration 1811 / 4900) loss: 1.645888
(Iteration 1821 / 4900) loss: 1.583925
(Iteration 1831 / 4900) loss: 1.671382
(Iteration 1841 / 4900) loss: 1.614654
(Iteration 1851 / 4900) loss: 1.716625
(Iteration 1861 / 4900) loss: 1.590578
(Iteration 1871 / 4900) loss: 1.664115
(Iteration 1881 / 4900) loss: 1.616829
(Iteration 1891 / 4900) loss: 1.720497
(Iteration 1901 / 4900) loss: 1.776410
(Iteration 1911 / 4900) loss: 1.588803
(Iteration 1921 / 4900) loss: 1.812399
(Iteration 1931 / 4900) loss: 1.593358
(Iteration 1941 / 4900) loss: 1.624999
(Iteration 1951 / 4900) loss: 1.706136
(Epoch 4 / 10) train acc: 0.403000; val_acc: 0.426000
(Iteration 1961 / 4900) loss: 1.458841
(Iteration 1971 / 4900) loss: 1.644322
(Iteration 1981 / 4900) loss: 1.718680
(Iteration 1991 / 4900) loss: 1.685929
(Iteration 2001 / 4900) loss: 1.655020
(Iteration 2011 / 4900) loss: 1.688121
(Iteration 2021 / 4900) loss: 1.677948
(Iteration 2031 / 4900) loss: 1.870887
(Iteration 2041 / 4900) loss: 1.668987
(Iteration 2051 / 4900) loss: 1.611042
(Iteration 2061 / 4900) loss: 1.683642
(Iteration 2071 / 4900) loss: 1.589631
(Iteration 2081 / 4900) loss: 1.720237
(Iteration 2091 / 4900) loss: 1.529303
(Iteration 2101 / 4900) loss: 1.624574
(Iteration 2111 / 4900) loss: 1.615859
(Iteration 2121 / 4900) loss: 1.717363
(Iteration 2131 / 4900) loss: 1.665371
(Iteration 2141 / 4900) loss: 1.789255
(Iteration 2151 / 4900) loss: 1.622238
(Iteration 2161 / 4900) loss: 1.739577
(Iteration 2171 / 4900) loss: 1.530157

(Iteration 2181 / 4900) loss: 1.639242
(Iteration 2191 / 4900) loss: 1.767974
(Iteration 2201 / 4900) loss: 1.844418
(Iteration 2211 / 4900) loss: 1.763672
(Iteration 2221 / 4900) loss: 1.556612
(Iteration 2231 / 4900) loss: 1.703411
(Iteration 2241 / 4900) loss: 1.614494
(Iteration 2251 / 4900) loss: 1.621294
(Iteration 2261 / 4900) loss: 1.606939
(Iteration 2271 / 4900) loss: 1.587882
(Iteration 2281 / 4900) loss: 1.573547
(Iteration 2291 / 4900) loss: 1.683787
(Iteration 2301 / 4900) loss: 1.447170
(Iteration 2311 / 4900) loss: 1.633658
(Iteration 2321 / 4900) loss: 1.603134
(Iteration 2331 / 4900) loss: 1.649711
(Iteration 2341 / 4900) loss: 1.519442
(Iteration 2351 / 4900) loss: 1.706299
(Iteration 2361 / 4900) loss: 1.555587
(Iteration 2371 / 4900) loss: 1.597922
(Iteration 2381 / 4900) loss: 1.577607
(Iteration 2391 / 4900) loss: 1.522055
(Iteration 2401 / 4900) loss: 1.617801
(Iteration 2411 / 4900) loss: 1.501640
(Iteration 2421 / 4900) loss: 1.572696
(Iteration 2431 / 4900) loss: 1.706325
(Iteration 2441 / 4900) loss: 1.681593
(Epoch 5 / 10) train acc: 0.408000; val_acc: 0.431000
(Iteration 2451 / 4900) loss: 1.751486
(Iteration 2461 / 4900) loss: 1.622609
(Iteration 2471 / 4900) loss: 1.781154
(Iteration 2481 / 4900) loss: 1.600690
(Iteration 2491 / 4900) loss: 1.663584
(Iteration 2501 / 4900) loss: 1.567472
(Iteration 2511 / 4900) loss: 1.659099
(Iteration 2521 / 4900) loss: 1.562896
(Iteration 2531 / 4900) loss: 1.589824
(Iteration 2541 / 4900) loss: 1.619890
(Iteration 2551 / 4900) loss: 1.564083
(Iteration 2561 / 4900) loss: 1.428909
(Iteration 2571 / 4900) loss: 1.760323
(Iteration 2581 / 4900) loss: 1.561412
(Iteration 2591 / 4900) loss: 1.517823
(Iteration 2601 / 4900) loss: 1.674657
(Iteration 2611 / 4900) loss: 1.633549
(Iteration 2621 / 4900) loss: 1.666899
(Iteration 2631 / 4900) loss: 1.558462
(Iteration 2641 / 4900) loss: 1.539278

(Iteration 2651 / 4900) loss: 1.644847
(Iteration 2661 / 4900) loss: 1.498245
(Iteration 2671 / 4900) loss: 1.648804
(Iteration 2681 / 4900) loss: 1.721081
(Iteration 2691 / 4900) loss: 1.625495
(Iteration 2701 / 4900) loss: 1.560354
(Iteration 2711 / 4900) loss: 1.574463
(Iteration 2721 / 4900) loss: 1.651460
(Iteration 2731 / 4900) loss: 1.435066
(Iteration 2741 / 4900) loss: 1.525976
(Iteration 2751 / 4900) loss: 1.633446
(Iteration 2761 / 4900) loss: 1.487799
(Iteration 2771 / 4900) loss: 1.475543
(Iteration 2781 / 4900) loss: 1.540993
(Iteration 2791 / 4900) loss: 1.721756
(Iteration 2801 / 4900) loss: 1.703828
(Iteration 2811 / 4900) loss: 1.882478
(Iteration 2821 / 4900) loss: 1.584537
(Iteration 2831 / 4900) loss: 1.628555
(Iteration 2841 / 4900) loss: 1.605570
(Iteration 2851 / 4900) loss: 1.606625
(Iteration 2861 / 4900) loss: 1.631563
(Iteration 2871 / 4900) loss: 1.699675
(Iteration 2881 / 4900) loss: 1.382116
(Iteration 2891 / 4900) loss: 1.467813
(Iteration 2901 / 4900) loss: 1.466214
(Iteration 2911 / 4900) loss: 1.700998
(Iteration 2921 / 4900) loss: 1.645747
(Iteration 2931 / 4900) loss: 1.891789
(Epoch 6 / 10) train acc: 0.451000; val_acc: 0.447000
(Iteration 2941 / 4900) loss: 1.575182
(Iteration 2951 / 4900) loss: 1.589419
(Iteration 2961 / 4900) loss: 1.590541
(Iteration 2971 / 4900) loss: 1.535560
(Iteration 2981 / 4900) loss: 1.738069
(Iteration 2991 / 4900) loss: 1.440787
(Iteration 3001 / 4900) loss: 1.574692
(Iteration 3011 / 4900) loss: 1.656443
(Iteration 3021 / 4900) loss: 1.514266
(Iteration 3031 / 4900) loss: 1.629376
(Iteration 3041 / 4900) loss: 1.596641
(Iteration 3051 / 4900) loss: 1.583588
(Iteration 3061 / 4900) loss: 1.547650
(Iteration 3071 / 4900) loss: 1.478644
(Iteration 3081 / 4900) loss: 1.508567
(Iteration 3091 / 4900) loss: 1.515584
(Iteration 3101 / 4900) loss: 1.521753
(Iteration 3111 / 4900) loss: 1.436926

(Iteration 3121 / 4900) loss: 1.650917
(Iteration 3131 / 4900) loss: 1.459509
(Iteration 3141 / 4900) loss: 1.720790
(Iteration 3151 / 4900) loss: 1.473725
(Iteration 3161 / 4900) loss: 1.667788
(Iteration 3171 / 4900) loss: 1.418745
(Iteration 3181 / 4900) loss: 1.588805
(Iteration 3191 / 4900) loss: 1.521273
(Iteration 3201 / 4900) loss: 1.599726
(Iteration 3211 / 4900) loss: 1.559156
(Iteration 3221 / 4900) loss: 1.538093
(Iteration 3231 / 4900) loss: 1.700120
(Iteration 3241 / 4900) loss: 1.569651
(Iteration 3251 / 4900) loss: 1.506105
(Iteration 3261 / 4900) loss: 1.496333
(Iteration 3271 / 4900) loss: 1.759554
(Iteration 3281 / 4900) loss: 1.570052
(Iteration 3291 / 4900) loss: 1.460758
(Iteration 3301 / 4900) loss: 1.504628
(Iteration 3311 / 4900) loss: 1.534143
(Iteration 3321 / 4900) loss: 1.606255
(Iteration 3331 / 4900) loss: 1.436728
(Iteration 3341 / 4900) loss: 1.542199
(Iteration 3351 / 4900) loss: 1.460853
(Iteration 3361 / 4900) loss: 1.567534
(Iteration 3371 / 4900) loss: 1.691367
(Iteration 3381 / 4900) loss: 1.550098
(Iteration 3391 / 4900) loss: 1.680122
(Iteration 3401 / 4900) loss: 1.390429
(Iteration 3411 / 4900) loss: 1.521238
(Iteration 3421 / 4900) loss: 1.499048
(Epoch 7 / 10) train acc: 0.469000; val_acc: 0.451000
(Iteration 3431 / 4900) loss: 1.535694
(Iteration 3441 / 4900) loss: 1.456257
(Iteration 3451 / 4900) loss: 1.631885
(Iteration 3461 / 4900) loss: 1.587840
(Iteration 3471 / 4900) loss: 1.555226
(Iteration 3481 / 4900) loss: 1.457427
(Iteration 3491 / 4900) loss: 1.711124
(Iteration 3501 / 4900) loss: 1.714735
(Iteration 3511 / 4900) loss: 1.702186
(Iteration 3521 / 4900) loss: 1.681523
(Iteration 3531 / 4900) loss: 1.555984
(Iteration 3541 / 4900) loss: 1.628258
(Iteration 3551 / 4900) loss: 1.480815
(Iteration 3561 / 4900) loss: 1.516137
(Iteration 3571 / 4900) loss: 1.607945
(Iteration 3581 / 4900) loss: 1.464008

(Iteration 3591 / 4900) loss: 1.522531
(Iteration 3601 / 4900) loss: 1.386743
(Iteration 3611 / 4900) loss: 1.541850
(Iteration 3621 / 4900) loss: 1.596365
(Iteration 3631 / 4900) loss: 1.470899
(Iteration 3641 / 4900) loss: 1.504619
(Iteration 3651 / 4900) loss: 1.549541
(Iteration 3661 / 4900) loss: 1.479146
(Iteration 3671 / 4900) loss: 1.569918
(Iteration 3681 / 4900) loss: 1.525105
(Iteration 3691 / 4900) loss: 1.626181
(Iteration 3701 / 4900) loss: 1.828679
(Iteration 3711 / 4900) loss: 1.677590
(Iteration 3721 / 4900) loss: 1.440147
(Iteration 3731 / 4900) loss: 1.557177
(Iteration 3741 / 4900) loss: 1.533165
(Iteration 3751 / 4900) loss: 1.641170
(Iteration 3761 / 4900) loss: 1.509822
(Iteration 3771 / 4900) loss: 1.471518
(Iteration 3781 / 4900) loss: 1.516810
(Iteration 3791 / 4900) loss: 1.506589
(Iteration 3801 / 4900) loss: 1.548596
(Iteration 3811 / 4900) loss: 1.499590
(Iteration 3821 / 4900) loss: 1.435734
(Iteration 3831 / 4900) loss: 1.696941
(Iteration 3841 / 4900) loss: 1.709301
(Iteration 3851 / 4900) loss: 1.507234
(Iteration 3861 / 4900) loss: 1.483066
(Iteration 3871 / 4900) loss: 1.548691
(Iteration 3881 / 4900) loss: 1.479723
(Iteration 3891 / 4900) loss: 1.535084
(Iteration 3901 / 4900) loss: 1.448461
(Iteration 3911 / 4900) loss: 1.527180
(Epoch 8 / 10) train acc: 0.462000; val_acc: 0.467000
(Iteration 3921 / 4900) loss: 1.465790
(Iteration 3931 / 4900) loss: 1.513631
(Iteration 3941 / 4900) loss: 1.378480
(Iteration 3951 / 4900) loss: 1.568631
(Iteration 3961 / 4900) loss: 1.612067
(Iteration 3971 / 4900) loss: 1.304457
(Iteration 3981 / 4900) loss: 1.387809
(Iteration 3991 / 4900) loss: 1.434775
(Iteration 4001 / 4900) loss: 1.433746
(Iteration 4011 / 4900) loss: 1.322788
(Iteration 4021 / 4900) loss: 1.690131
(Iteration 4031 / 4900) loss: 1.769119
(Iteration 4041 / 4900) loss: 1.367324
(Iteration 4051 / 4900) loss: 1.663580

(Iteration 4061 / 4900) loss: 1.550479
(Iteration 4071 / 4900) loss: 1.559759
(Iteration 4081 / 4900) loss: 1.563596
(Iteration 4091 / 4900) loss: 1.503273
(Iteration 4101 / 4900) loss: 1.635513
(Iteration 4111 / 4900) loss: 1.606411
(Iteration 4121 / 4900) loss: 1.660035
(Iteration 4131 / 4900) loss: 1.754815
(Iteration 4141 / 4900) loss: 1.280390
(Iteration 4151 / 4900) loss: 1.311317
(Iteration 4161 / 4900) loss: 1.503973
(Iteration 4171 / 4900) loss: 1.588227
(Iteration 4181 / 4900) loss: 1.522910
(Iteration 4191 / 4900) loss: 1.429910
(Iteration 4201 / 4900) loss: 1.449076
(Iteration 4211 / 4900) loss: 1.404975
(Iteration 4221 / 4900) loss: 1.543679
(Iteration 4231 / 4900) loss: 1.466628
(Iteration 4241 / 4900) loss: 1.594160
(Iteration 4251 / 4900) loss: 1.692660
(Iteration 4261 / 4900) loss: 1.610414
(Iteration 4271 / 4900) loss: 1.467253
(Iteration 4281 / 4900) loss: 1.410744
(Iteration 4291 / 4900) loss: 1.308937
(Iteration 4301 / 4900) loss: 1.535322
(Iteration 4311 / 4900) loss: 1.319935
(Iteration 4321 / 4900) loss: 1.555955
(Iteration 4331 / 4900) loss: 1.470758
(Iteration 4341 / 4900) loss: 1.597489
(Iteration 4351 / 4900) loss: 1.503501
(Iteration 4361 / 4900) loss: 1.408904
(Iteration 4371 / 4900) loss: 1.472090
(Iteration 4381 / 4900) loss: 1.356957
(Iteration 4391 / 4900) loss: 1.621898
(Iteration 4401 / 4900) loss: 1.376167
(Epoch 9 / 10) train acc: 0.489000; val_acc: 0.457000
(Iteration 4411 / 4900) loss: 1.738783
(Iteration 4421 / 4900) loss: 1.290290
(Iteration 4431 / 4900) loss: 1.597577
(Iteration 4441 / 4900) loss: 1.545853
(Iteration 4451 / 4900) loss: 1.543258
(Iteration 4461 / 4900) loss: 1.344586
(Iteration 4471 / 4900) loss: 1.601376
(Iteration 4481 / 4900) loss: 1.550930
(Iteration 4491 / 4900) loss: 1.379485
(Iteration 4501 / 4900) loss: 1.580567
(Iteration 4511 / 4900) loss: 1.544913
(Iteration 4521 / 4900) loss: 1.334137

```
(Iteration 4531 / 4900) loss: 1.862890
(Iteration 4541 / 4900) loss: 1.426022
(Iteration 4551 / 4900) loss: 1.504377
(Iteration 4561 / 4900) loss: 1.617898
(Iteration 4571 / 4900) loss: 1.529963
(Iteration 4581 / 4900) loss: 1.521095
(Iteration 4591 / 4900) loss: 1.417431
(Iteration 4601 / 4900) loss: 1.423031
(Iteration 4611 / 4900) loss: 1.695511
(Iteration 4621 / 4900) loss: 1.500624
(Iteration 4631 / 4900) loss: 1.625314
(Iteration 4641 / 4900) loss: 1.471987
(Iteration 4651 / 4900) loss: 1.625075
(Iteration 4661 / 4900) loss: 1.750594
(Iteration 4671 / 4900) loss: 1.680011
(Iteration 4681 / 4900) loss: 1.438443
(Iteration 4691 / 4900) loss: 1.361646
(Iteration 4701 / 4900) loss: 1.545790
(Iteration 4711 / 4900) loss: 1.436357
(Iteration 4721 / 4900) loss: 1.461629
(Iteration 4731 / 4900) loss: 1.420825
(Iteration 4741 / 4900) loss: 1.477342
(Iteration 4751 / 4900) loss: 1.440879
(Iteration 4761 / 4900) loss: 1.425149
(Iteration 4771 / 4900) loss: 1.408250
(Iteration 4781 / 4900) loss: 1.497814
(Iteration 4791 / 4900) loss: 1.455349
(Iteration 4801 / 4900) loss: 1.458274
(Iteration 4811 / 4900) loss: 1.589643
(Iteration 4821 / 4900) loss: 1.665653
(Iteration 4831 / 4900) loss: 1.645997
(Iteration 4841 / 4900) loss: 1.304858
(Iteration 4851 / 4900) loss: 1.318106
(Iteration 4861 / 4900) loss: 1.272946
(Iteration 4871 / 4900) loss: 1.528606
(Iteration 4881 / 4900) loss: 1.385325
(Iteration 4891 / 4900) loss: 1.378594
(Epoch 10 / 10) train acc: 0.449000; val_acc: 0.476000
```

10 Debug the training

With the default parameters we provided above, you should get a validation accuracy of about 0.36 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.

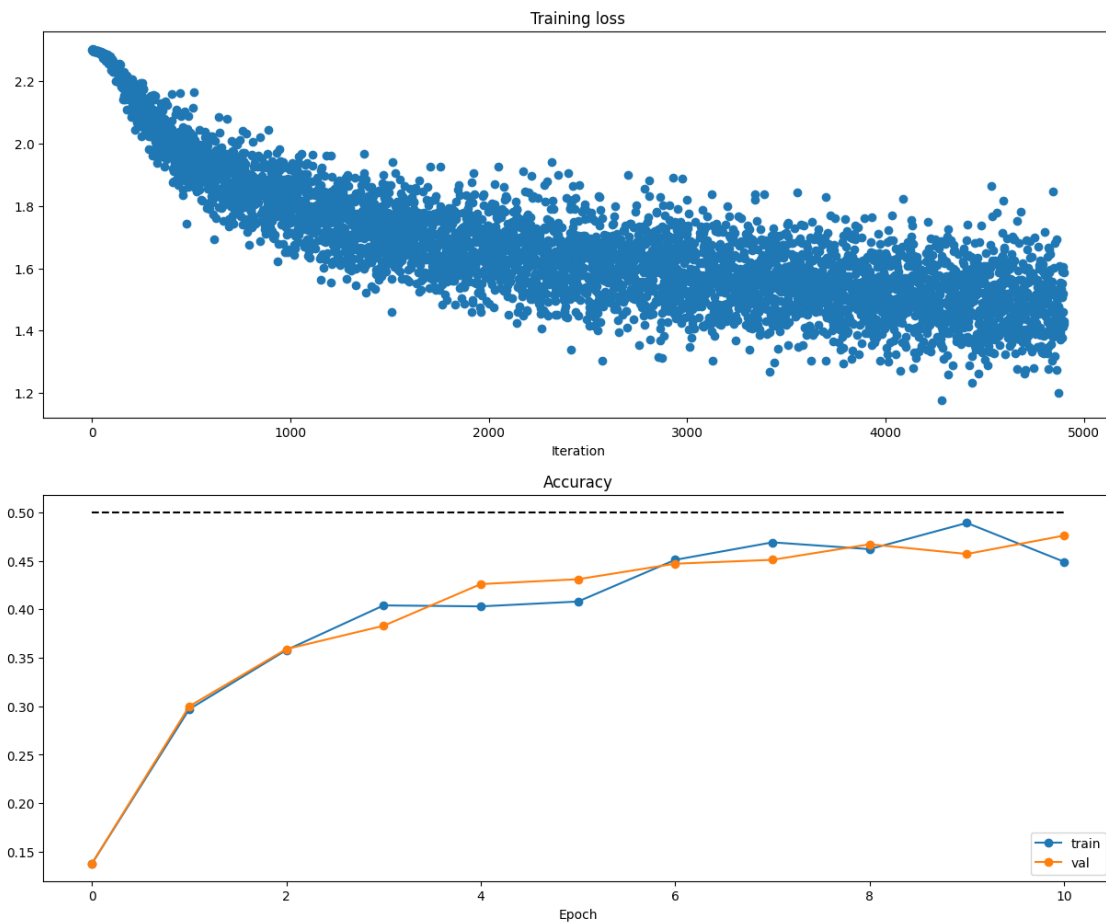
Another strategy is to visualize the weights that were learned in the first layer of the network. In

most neural networks trained on visual data, the first layer weights typically show some visible structure when visualized.

```
[12]: # Run this cell to visualize training loss and train / val accuracy
```

```
plt.subplot(2, 1, 1)
plt.title('Training loss')
plt.plot(solver.loss_history, 'o')
plt.xlabel('Iteration')

plt.subplot(2, 1, 2)
plt.title('Accuracy')
plt.plot(solver.train_acc_history, '-o', label='train')
plt.plot(solver.val_acc_history, '-o', label='val')
plt.plot([0.5] * len(solver.val_acc_history), 'k--')
plt.xlabel('Epoch')
plt.legend(loc='lower right')
plt.gcf().set_size_inches(15, 12)
plt.show()
```

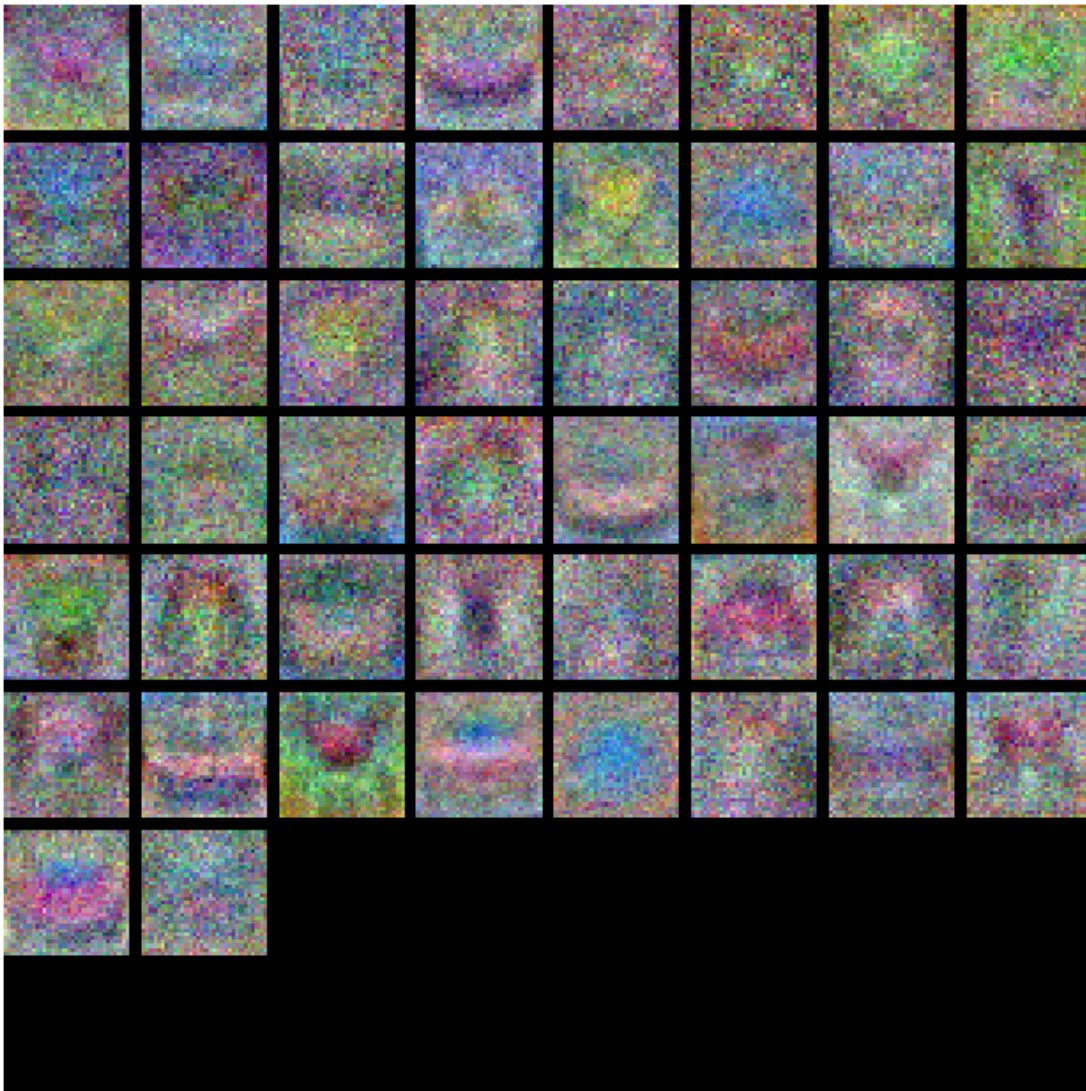


```
[13]: from icv83551.vis_utils import visualize_grid

# Visualize the weights of the network

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

show_net_weights(model)
```



11 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 be 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.).

```
[15]: best_model = None

#####
# TODO: Tune hyperparameters using the validation set. Store your best trained
#
# model in best_model.
#
#
# 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 thes previous exercises.
↪ #
#####

best_val = -1
results = {}

learning_rates = [1e-4, 1e-3, 1e-2]
regularization_strengths = [1e-5, 1e-4, 1e-3]

import itertools
for lr, reg in itertools.product(learning_rates, regularization_strengths):
    print(f'lr = {lr} reg = {reg}')
    # Create a new TwoLayerNet instance
    model = TwoLayerNet(hidden_dim=30, reg=reg)

    # Solver
    solver = Solver(model, data, optim_config={'learning_rate': lr},
↪ num_epochs=10, verbose=False)
    solver.train()

    # Performance on validation set
    results[(lr, reg)] = solver.best_val_acc

    if results[(lr, reg)] > best_val: # Save the best model
        best_val = results[(lr, reg)]
        best_model = model

# Print results
for lr, reg in sorted(results):
    val_acc = results[(lr, reg)]
    print(f'lr {lr} reg {reg} val accuracy: {val_acc}')

print(f'best validation accuracy achieved: {best_val}')
#####
#                               END OF YOUR CODE                               #
#####

```

```

lr = 0.0001 reg = 1e-05
lr = 0.0001 reg = 0.0001
lr = 0.0001 reg = 0.001
lr = 0.001 reg = 1e-05
lr = 0.001 reg = 0.0001
lr = 0.001 reg = 0.001
lr = 0.01 reg = 1e-05

```

/content/drive/MyDrive/icv83551/assignments/assignment1/icv83551/layers.py:731:
RuntimeWarning: divide by zero encountered in log

```

/content/drive/MyDrive/icv83551/assignments/assignment1/icv83551/layers.py:33:
RuntimeWarning: overflow encountered in matmul
    x_flat = x.reshape(N, -1)
/content/drive/MyDrive/icv83551/assignments/assignment1/icv83551/layers.py:729:
RuntimeWarning: overflow encountered in subtract
    """

/content/drive/MyDrive/icv83551/assignments/assignment1/icv83551/layers.py:729:
RuntimeWarning: invalid value encountered in subtract
    """

lr = 0.01 reg = 0.0001
lr = 0.01 reg = 0.001
lr 0.0001 reg 1e-05 val accuracy: 0.475
lr 0.0001 reg 0.0001 val accuracy: 0.469
lr 0.0001 reg 0.001 val accuracy: 0.46
lr 0.001 reg 1e-05 val accuracy: 0.483
lr 0.001 reg 0.0001 val accuracy: 0.48
lr 0.001 reg 0.001 val accuracy: 0.489
lr 0.01 reg 1e-05 val accuracy: 0.135
lr 0.01 reg 0.0001 val accuracy: 0.115
lr 0.01 reg 0.001 val accuracy: 0.131
best validation accuracy achieved: 0.489

```

12 Test your model!

Run your best model on the validation and test sets. You should achieve above 48% accuracy on the validation set and the test set.

```

[16]: y_val_pred = np.argmax(best_model.loss(data['X_val']), axis=1)
      print('Validation set accuracy: ', (y_val_pred == data['y_val']).mean())

```

Validation set accuracy: 0.489

```

[17]: y_test_pred = np.argmax(best_model.loss(data['X_test']), axis=1)
      print('Test set accuracy: ', (y_test_pred == data['y_test']).mean())

```

Test set accuracy: 0.469

```

[18]: # Save best model
      best_model.save("best_two_layer_net.npy")

```

best_two_layer_net.npy saved.

12.1 Inline Question 2:

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 and 3 can decrease the gap.

Your Explanation : This gap is an overfitting to the training data.

If we apply 1 - train a larger dataset the model will be exposed to more data and will have a harder time to learn features which are noises (as it won't be on all the class samples).

If we apply 2 - add more hidden units: will probably worsen the model, as it will learn more complex features and perform even better on the training set and overfit more.

If we apply 3 - increase the regularization strength: weights will be penalized and it will prevent the model from rely only on the larger weights.

[]:

features

December 30, 2025

```
[1]: # This mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive')

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'icv83551/assignments/assignment1/'
FOLDERNAME = 'icv83551/assignments/assignment1/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

# Now that we've mounted your Drive, this ensures that
# the Python interpreter of the Colab VM can load
# python files from within it.
import sys
sys.path.append('/content/drive/My Drive/{}'.format(FOLDERNAME))

# This downloads the CIFAR-10 dataset to your Drive
# if it doesn't already exist.
%cd /content/drive/My\ Drive/$FOLDERNAME/icv83551/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

```
Mounted at /content/drive
/content/drive/My Drive/icv83551/assignments/assignment1/icv83551/datasets
/content/drive/My Drive/icv83551/assignments/assignment1
```

1 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.

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.

```
[2]: import random
import numpy as np
from icv83551.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'

#%%load_ext autoreload
#%%autoreload 2
```

1.1 Load data

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

```
[3]: from icv83551.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 = 'icv83551/datasets/cifar-10-batches-py'

    # Cleaning up variables to prevent loading data multiple times (which may
    # cause memory issue)
    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]

    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()
```

1.2 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.

```
[4]: from icv83551.features import *

# num_color_bins = 10 # Number of bins in the color histogram
num_color_bins = 25 # 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

[illegible]

1.3 Train Softmax classifier on features

Using the Softmax code developed earlier in the assignment, train Softmax classifiers on top of the features extracted above; this should achieve better results than training them directly on top of raw pixels.

```
[5]: # Use the validation set to tune the learning rate and regularization strength

from icv83551.classifiers.linear_classifier import Softmax

learning_rates = [1e-5, 1e-4, 1e-3, 1e-2, 1e-1]
regularization_strengths = [1e-3, 1e-2, 1e-1]

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 Softmax; save#
# the best trained classifier in best_softmax. If you carefully tune the model, #
# you should be able to get accuracy of above 0.42 on the validation set.      #
#####
import itertools
for lr, reg in itertools.product(learning_rates, regularization_strengths):
    # Create a Softmax classifier
    model = Softmax()
    model.train(X_train_feats, y_train, lr, reg, num_iters=1000)

    y_train_pred, y_val_pred = model.predict(X_train_feats), model.
    ↪predict(X_val_feats)

    # Check performance on validation and train set
    results[(lr, reg)] = np.mean(y_train == y_train_pred), np.mean(y_val ==
    ↪y_val_pred)

    if results[(lr, reg)][1] > best_val: # Save the best model
        best_val = results[(lr, reg)][1]
        best_softmax = model

# 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: %f' % best_val)
```

```

lr 1.000000e-05 reg 1.000000e-03 train accuracy: 0.231163 val accuracy: 0.246000
lr 1.000000e-05 reg 1.000000e-02 train accuracy: 0.212388 val accuracy: 0.214000
lr 1.000000e-05 reg 1.000000e-01 train accuracy: 0.236265 val accuracy: 0.254000
lr 1.000000e-04 reg 1.000000e-03 train accuracy: 0.414714 val accuracy: 0.403000
lr 1.000000e-04 reg 1.000000e-02 train accuracy: 0.417469 val accuracy: 0.413000
lr 1.000000e-04 reg 1.000000e-01 train accuracy: 0.407612 val accuracy: 0.407000
lr 1.000000e-03 reg 1.000000e-03 train accuracy: 0.440980 val accuracy: 0.438000
lr 1.000000e-03 reg 1.000000e-02 train accuracy: 0.441796 val accuracy: 0.444000
lr 1.000000e-03 reg 1.000000e-01 train accuracy: 0.441776 val accuracy: 0.444000
lr 1.000000e-02 reg 1.000000e-03 train accuracy: 0.505102 val accuracy: 0.493000
lr 1.000000e-02 reg 1.000000e-02 train accuracy: 0.503612 val accuracy: 0.499000
lr 1.000000e-02 reg 1.000000e-01 train accuracy: 0.485388 val accuracy: 0.479000
lr 1.000000e-01 reg 1.000000e-03 train accuracy: 0.521122 val accuracy: 0.526000
lr 1.000000e-01 reg 1.000000e-02 train accuracy: 0.516204 val accuracy: 0.511000
lr 1.000000e-01 reg 1.000000e-01 train accuracy: 0.486898 val accuracy: 0.484000
best validation accuracy achieved: 0.526000

```

```

[6]: # Evaluate your trained Softmax on the test set: you should be able to get at_
      ↪ least 0.42
y_test_pred = best_softmax.predict(X_test_feats)
test_accuracy = np.mean(y_test == y_test_pred)
print(test_accuracy)

```

0.492

```

[7]: # Save best softmax model
best_softmax.save("best_softmax_features.npy")

```

best_softmax_features.npy saved.

```

[8]: # 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', 'truck']
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)

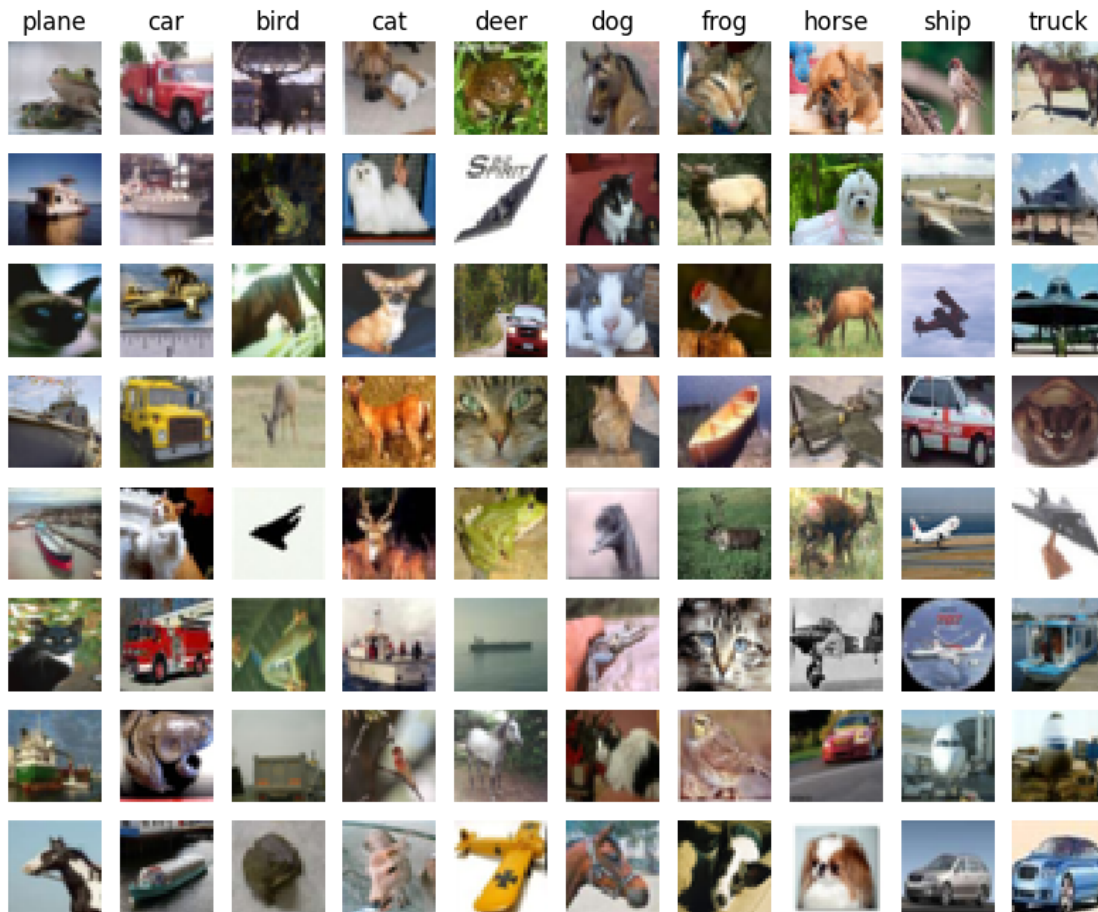
        plt.imshow(X_test[idx].astype('uint8'))
        plt.axis('off')

```

```

if i == 0:
    plt.title(cls_name)
plt.show()

```



1.3.1 Inline question 1:

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

Your Answer : Some of the misclassifications make sense because of some visual or geometric overlaps. For example a bird misclassified as a frog likely due to its beak looking like frog face, or a bird misclassified as a plane likely due to the blue sky background and wings.

1.4 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.

```
[9]: # 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, 170)
```

```
(49000, 169)
```

```
[11]: from icv83551.classifiers.fc_net import TwoLayerNet
from icv83551.solver import Solver

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

data = {
    'X_train': X_train_feats,
    'y_train': y_train,
    'X_val': X_val_feats,
    'y_val': y_val,
    'X_test': X_test_feats,
    'y_test': y_test,
}

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. #
#####

results = {}
best_val = -1

learning_rates = [1e-5, 1e-4, 1e-3, 1e-2, 1e-1]
regularization_strengths = [1e-3, 1e-2, 1e-1]

import itertools
```

```

for lr, reg in itertools.product(learning_rates, regularization_strengths):
    print(f'currently running lr {lr} and reg {reg}')

    # Create Two Layer Net and train it with Solver
    current_net = TwoLayerNet(input_dim, hidden_dim, num_classes, reg=reg)
    solver = Solver(current_net, data, optim_config={'learning_rate': lr},
    num_epochs=10, verbose=False)
    solver.train()

    # Compute validation set accuracy and append to the dictionary
    results[(lr, reg)] = solver.best_val_acc

    if results[(lr, reg)] > best_val: # Save if validation accuracy is the best
        best_val = results[(lr, reg)]
        best_net = current_net

# Print out results.
for lr, reg in sorted(results):
    val_accuracy = results[(lr, reg)]
    print(f'lr {lr} reg {reg} val accuracy: {val_accuracy}')

print(f'best validation: {best_val}')

```

```

currently running lr 1e-05 and reg 0.001
currently running lr 1e-05 and reg 0.01
currently running lr 1e-05 and reg 0.1
currently running lr 0.0001 and reg 0.001
currently running lr 0.0001 and reg 0.01
currently running lr 0.0001 and reg 0.1
currently running lr 0.001 and reg 0.001
currently running lr 0.001 and reg 0.01
currently running lr 0.001 and reg 0.1
currently running lr 0.01 and reg 0.001
currently running lr 0.01 and reg 0.01
currently running lr 0.01 and reg 0.1
currently running lr 0.1 and reg 0.001
currently running lr 0.1 and reg 0.01
currently running lr 0.1 and reg 0.1
lr 1e-05 reg 0.001 val accuracy: 0.077
lr 1e-05 reg 0.01 val accuracy: 0.121
lr 1e-05 reg 0.1 val accuracy: 0.089
lr 0.0001 reg 0.001 val accuracy: 0.159
lr 0.0001 reg 0.01 val accuracy: 0.15
lr 0.0001 reg 0.1 val accuracy: 0.174
lr 0.001 reg 0.001 val accuracy: 0.282
lr 0.001 reg 0.01 val accuracy: 0.278
lr 0.001 reg 0.1 val accuracy: 0.202

```

```
lr 0.01 reg 0.001 val accuracy: 0.513
lr 0.01 reg 0.01 val accuracy: 0.503
lr 0.01 reg 0.1 val accuracy: 0.417
lr 0.1 reg 0.001 val accuracy: 0.59
lr 0.1 reg 0.01 val accuracy: 0.559
lr 0.1 reg 0.1 val accuracy: 0.444
best validation: 0.59
```

```
[12]: # Run your best neural net classifier on the test set. You should be able
# to get more than 58% accuracy. It is also possible to get >60% accuracy
# with careful tuning.
```

```
y_test_pred = np.argmax(best_net.loss(data['X_test']), axis=1)
test_acc = (y_test_pred == data['y_test']).mean()
print(test_acc)
```

0.574

```
[13]: # Save best model
best_net.save("best_two_layer_net_features.npy")
```

best_two_layer_net_features.npy saved.

```
[ ]:
```

FullyConnectedNets

December 30, 2025

```
[1]: # This mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive')

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'icv83551/assignments/assignment1/'
FOLDERNAME = 'icv83551/assignments/assignment1/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

# Now that we've mounted your Drive, this ensures that
# the Python interpreter of the Colab VM can load
# python files from within it.
import sys
sys.path.append('/content/drive/My Drive/{}'.format(FOLDERNAME))

# This downloads the CIFAR-10 dataset to your Drive
# if it doesn't already exist.
%cd /content/drive/My\ Drive/$FOLDERNAME/icv83551/datasets/
# %cd /content/drive/My\ Drive/$FOLDERNAME
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call `drive.mount("/content/drive", force_remount=True)`.

/content/drive/My Drive/icv83551/assignments/assignment1/icv83551/datasets

/content/drive/My Drive/icv83551/assignments/assignment1

1 Multi-Layer Fully Connected Network

In this exercise, you will implement a fully connected network with an arbitrary number of hidden layers.

```
[2]: # from google.colab import drive
# drive.mount('/content/drive')
```

Read through the `FullyConnectedNet` class in the file `icv83551/classifiers/fc_net.py`.

Implement the network initialization, forward pass, and backward pass. Throughout this assign-

ment, you will be implementing layers in `icv83551/layers.py`. You can re-use your implementations for `affine_forward`, `affine_backward`, `relu_forward`, `relu_backward`, and `softmax_loss` from before. For right now, don't worry about implementing dropout or batch/layer normalization yet, as you will add those features later.

```
[3]: # Setup cell.
import time
import numpy as np
import matplotlib.pyplot as plt
from icv83551.classifiers.fc_net import *
from icv83551.data_utils import get_CIFAR10_data
from icv83551.gradient_check import eval_numerical_gradient, \
    eval_numerical_gradient_array
from icv83551.solver import Solver

%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"

#%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))))
```

```
[4]: # Load the (preprocessed) CIFAR-10 data.
data = get_CIFAR10_data()
for k, v in list(data.items()):
    print(f"{k}: {v.shape}")
```

```
X_train: (49000, 3, 32, 32)
y_train: (49000,)
X_val: (1000, 3, 32, 32)
y_val: (1000,)
X_test: (1000, 3, 32, 32)
y_test: (1000,)
```

1.1 Initial Loss and Gradient Check

As a sanity check, run the following to check the initial loss and to gradient check the network both with and without regularization. This is a good way to see if the initial losses seem reasonable.

For gradient checking, you should expect to see errors around $1e-7$ or less.

```
[5]: np.random.seed(231)
N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
```

```

y = np.random.randint(C, size=(N,))

for reg in [0, 3.14]:
    print("Running check with reg = ", reg)
    model = FullyConnectedNet(
        [H1, H2],
        input_dim=D,
        num_classes=C,
        reg=reg,
        weight_scale=5e-2,
        dtype=np.float64
    )

    loss, grads = model.loss(X, y)
    print("Initial loss: ", loss)

    # Most of the errors should be on the order of e-7 or smaller.
    # NOTE: It is fine however to see an error for W2 on the order of e-5
    # for the check when reg = 0.0
    for name in sorted(grads):
        f = lambda _: model.loss(X, y)[0]
        grad_num = eval_numerical_gradient(f, model.params[name],
        verbose=False, h=1e-5)
        print(f"{name} relative error: {rel_error(grad_num, grads[name])}")

```

```

Running check with reg = 0
Initial loss: 2.300479089768492
W1 relative error: 1.0252674471656573e-07
W2 relative error: 2.2120479295080622e-05
W3 relative error: 4.5623278736665505e-07
b1 relative error: 4.6600944653202505e-09
b2 relative error: 2.085654276112763e-09
b3 relative error: 1.689724888469736e-10
Running check with reg = 3.14
Initial loss: 7.052114776533016
W1 relative error: 1.409028728052923e-08
W2 relative error: 6.86942277940646e-08
W3 relative error: 2.131129859578198e-08
b1 relative error: 1.4752427965311745e-08
b2 relative error: 1.7223751746766738e-09
b3 relative error: 2.378772438198909e-10

```

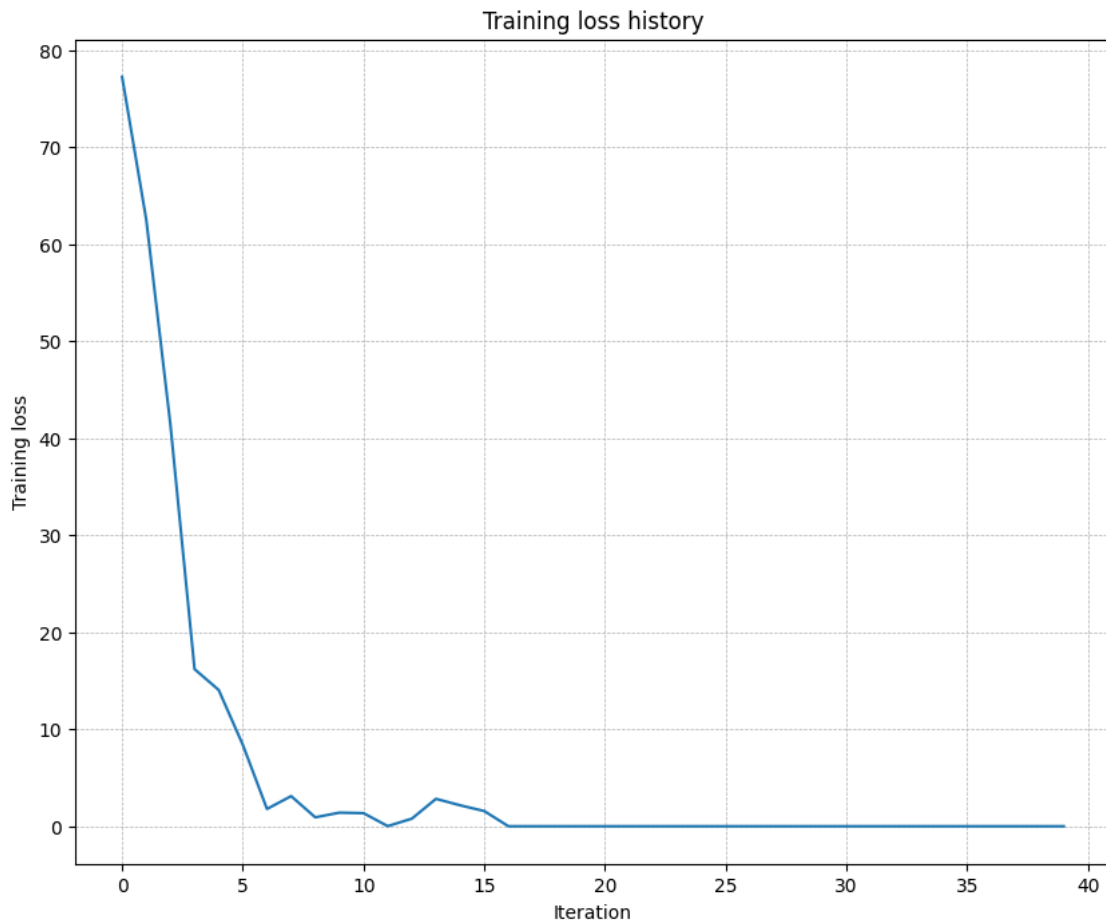
As another sanity check, make sure your network can overfit on a small dataset of 50 images. First, we will try a three-layer network with 100 units in each hidden layer. In the following cell, tweak the **learning rate** and **weight initialization scale** to overfit and achieve 100% training accuracy within 20 epochs.

```
[6]: # TODO: Use a three-layer Net to overfit 50 training examples by  
# tweaking just the learning rate and initialization scale.
```

```
num_train = 50  
small_data = {  
    "X_train": data["X_train"][:num_train],  
    "y_train": data["y_train"][:num_train],  
    "X_val": data["X_val"],  
    "y_val": data["y_val"],  
}  
  
weight_scale = 6e-2    # Experiment with this!  
learning_rate = 1e-3    # Experiment with this!  
  
model = FullyConnectedNet(  
    [100, 100],  
    weight_scale=weight_scale,  
    dtype=np.float64  
)  
solver = Solver(  
    model,  
    small_data,  
    print_every=10,  
    num_epochs=20,  
    batch_size=25,  
    update_rule="sgd",  
    optim_config={"learning_rate": learning_rate},  
)  
solver.train()  
  
plt.plot(solver.loss_history)  
plt.title("Training loss history")  
plt.xlabel("Iteration")  
plt.ylabel("Training loss")  
plt.grid(linestyle='--', linewidth=0.5)  
plt.show()
```

```
(Iteration 1 / 40) loss: 77.236010  
(Epoch 0 / 20) train acc: 0.180000; val_acc: 0.120000  
(Epoch 1 / 20) train acc: 0.340000; val_acc: 0.129000  
(Epoch 2 / 20) train acc: 0.500000; val_acc: 0.136000  
(Epoch 3 / 20) train acc: 0.720000; val_acc: 0.147000  
(Epoch 4 / 20) train acc: 0.820000; val_acc: 0.150000  
(Epoch 5 / 20) train acc: 0.860000; val_acc: 0.183000  
(Iteration 11 / 40) loss: 1.348874  
(Epoch 6 / 20) train acc: 0.920000; val_acc: 0.173000
```

```
(Epoch 7 / 20) train acc: 0.920000; val_acc: 0.165000
(Epoch 8 / 20) train acc: 1.000000; val_acc: 0.171000
(Epoch 9 / 20) train acc: 1.000000; val_acc: 0.171000
(Epoch 10 / 20) train acc: 1.000000; val_acc: 0.171000
(Iteration 21 / 40) loss: 0.000137
(Epoch 11 / 20) train acc: 1.000000; val_acc: 0.171000
(Epoch 12 / 20) train acc: 1.000000; val_acc: 0.171000
(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.171000
(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.171000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.171000
(Iteration 31 / 40) loss: 0.000444
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.171000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.171000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.171000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.171000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.172000
```



Now, try to use a five-layer network with 100 units on each layer to overfit on 50 training examples. Again, you will have to adjust the learning rate and weight initialization scale, but you should be

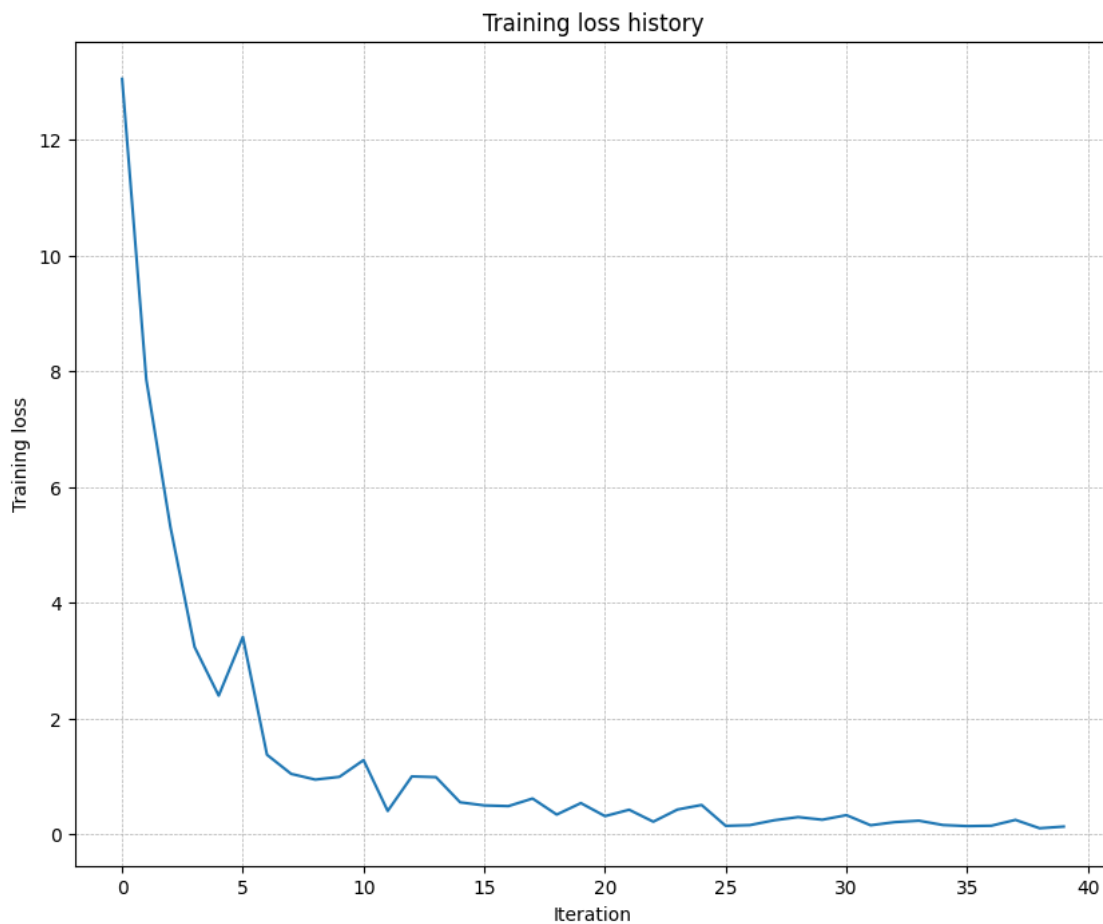
able to achieve 100% training accuracy within 20 epochs.

```
[7]: # TODO: Use a five-layer Net to overfit 50 training examples by  
# tweaking just the learning rate and initialization scale.
```

```
num_train = 50  
small_data = {  
    'X_train': data['X_train'][:num_train],  
    'y_train': data['y_train'][:num_train],  
    'X_val': data['X_val'],  
    'y_val': data['y_val'],  
}  
  
weight_scale = 6e-2    # Experiment with this!  
learning_rate = 1e-3    # Experiment with this!  
  
model = FullyConnectedNet(  
    [100, 100, 100, 100],  
    weight_scale=weight_scale,  
    dtype=np.float64  
)  
solver = Solver(  
    model,  
    small_data,  
    print_every=10,  
    num_epochs=20,  
    batch_size=25,  
    update_rule='sgd',  
    optim_config={'learning_rate': learning_rate},  
)  
solver.train()  
  
plt.plot(solver.loss_history)  
plt.title('Training loss history')  
plt.xlabel('Iteration')  
plt.ylabel('Training loss')  
plt.grid(linestyle='--', linewidth=0.5)  
plt.show()
```

```
(Iteration 1 / 40) loss: 13.054907  
(Epoch 0 / 20) train acc: 0.180000; val_acc: 0.120000  
(Epoch 1 / 20) train acc: 0.260000; val_acc: 0.112000  
(Epoch 2 / 20) train acc: 0.380000; val_acc: 0.112000  
(Epoch 3 / 20) train acc: 0.480000; val_acc: 0.121000  
(Epoch 4 / 20) train acc: 0.660000; val_acc: 0.127000  
(Epoch 5 / 20) train acc: 0.660000; val_acc: 0.112000  
(Iteration 11 / 40) loss: 1.276506
```

(Epoch 6 / 20) train acc: 0.820000; val_acc: 0.122000
(Epoch 7 / 20) train acc: 0.820000; val_acc: 0.127000
(Epoch 8 / 20) train acc: 0.880000; val_acc: 0.121000
(Epoch 9 / 20) train acc: 0.860000; val_acc: 0.119000
(Epoch 10 / 20) train acc: 0.960000; val_acc: 0.129000
(Iteration 21 / 40) loss: 0.307482
(Epoch 11 / 20) train acc: 0.960000; val_acc: 0.122000
(Epoch 12 / 20) train acc: 0.960000; val_acc: 0.132000
(Epoch 13 / 20) train acc: 0.960000; val_acc: 0.126000
(Epoch 14 / 20) train acc: 0.960000; val_acc: 0.126000
(Epoch 15 / 20) train acc: 0.960000; val_acc: 0.129000
(Iteration 31 / 40) loss: 0.324763
(Epoch 16 / 20) train acc: 0.980000; val_acc: 0.130000
(Epoch 17 / 20) train acc: 0.980000; val_acc: 0.127000
(Epoch 18 / 20) train acc: 0.980000; val_acc: 0.131000
(Epoch 19 / 20) train acc: 0.980000; val_acc: 0.121000
(Epoch 20 / 20) train acc: 0.980000; val_acc: 0.126000



1.2 Inline Question 1:

Did you notice anything about the comparative difficulty of training the three-layer network vs. training the five-layer network? In particular, based on your experience, which network seemed more sensitive to the initialization scale? Why do you think that is the case?

1.3 Answer:

It is more difficult for the five-layer network to train than for the three-layer network because for more layers there is more sensitivity to weight initialization and it is easier to find weights that give good accuracy for 3 layers as the weights require less effort. The larger the network \rightarrow the larger the risk of vanishing gradient.

2 Update rules

So far we have used vanilla stochastic gradient descent (SGD) as our update rule. More sophisticated update rules can make it easier to train deep networks. We will implement a few of the most commonly used update rules and compare them to vanilla SGD.

2.1 SGD+Momentum

Stochastic gradient descent with momentum is a widely used update rule that tends to make deep networks converge faster than vanilla stochastic gradient descent.

Open the file `icv83551/optim.py` and read the documentation at the top of the file to make sure you understand the API. Implement the SGD+momentum update rule in the function `sgd_momentum` and run the following to check your implementation. You should see errors less than $e-8$.

```
[8]: from icv83551.optim import sgd_momentum

N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
v = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)

config = {"learning_rate": 1e-3, "velocity": v}
next_w, _ = sgd_momentum(w, dw, config=config)

expected_next_w = np.asarray([
    [ 0.1406,      0.20738947,  0.27417895,  0.34096842,  0.40775789],
    [ 0.47454737,  0.54133684,  0.60812632,  0.67491579,  0.74170526],
    [ 0.80849474,  0.87528421,  0.94207368,  1.00886316,  1.07565263],
    [ 1.14244211,  1.20923158,  1.27602105,  1.34281053,  1.4096    ]])
expected_velocity = np.asarray([
    [ 0.5406,      0.55475789,  0.56891579,  0.58307368,  0.59723158],
    [ 0.61138947,  0.62554737,  0.63970526,  0.65386316,  0.66802105],
    [ 0.68217895,  0.69633684,  0.71049474,  0.72465263,  0.73881053],
    [ 0.75296842,  0.76712632,  0.78128421,  0.79544211,  0.8096    ]])
```

```
# Should see relative errors around e-8 or less
print("next_w error: ", rel_error(next_w, expected_next_w))
print("velocity error: ", rel_error(expected_velocity, config["velocity"]))
```

```
next_w error:  8.882347033505819e-09
velocity error: 4.269287743278663e-09
```

Once you have done so, run the following to train a six-layer network with both SGD and SGD+momentum. You should see the SGD+momentum update rule converge faster.

```
[9]: num_train = 4000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

solvers = {}

for update_rule in ['sgd', 'sgd_momentum']:
    print('Running with ', update_rule)
    model = FullyConnectedNet(
        [100, 100, 100, 100, 100],
        weight_scale=5e-2
    )

    solver = Solver(
        model,
        small_data,
        num_epochs=5,
        batch_size=100,
        update_rule=update_rule,
        optim_config={'learning_rate': 5e-3},
        verbose=True,
    )
    solvers[update_rule] = solver
    solver.train()

fig, axes = plt.subplots(3, 1, figsize=(15, 15))

axes[0].set_title('Training loss')
axes[0].set_xlabel('Iteration')
axes[1].set_title('Training accuracy')
axes[1].set_xlabel('Epoch')
axes[2].set_title('Validation accuracy')
axes[2].set_xlabel('Epoch')
```

```

for update_rule, solver in solvers.items():
    axes[0].plot(solver.loss_history, label=f"loss_{update_rule}")
    axes[1].plot(solver.train_acc_history, label=f"train_acc_{update_rule}")
    axes[2].plot(solver.val_acc_history, label=f"val_acc_{update_rule}")

for ax in axes:
    ax.legend(loc="best", ncol=4)
    ax.grid(linestyle='--', linewidth=0.5)

plt.show()

```

Running with `sgd`

```

(Iteration 1 / 200) loss: 2.559978
(Epoch 0 / 5) train acc: 0.104000; val_acc: 0.107000
(Iteration 11 / 200) loss: 2.356070
(Iteration 21 / 200) loss: 2.214091
(Iteration 31 / 200) loss: 2.205928
(Epoch 1 / 5) train acc: 0.225000; val_acc: 0.193000
(Iteration 41 / 200) loss: 2.132095
(Iteration 51 / 200) loss: 2.118950
(Iteration 61 / 200) loss: 2.116443
(Iteration 71 / 200) loss: 2.132549
(Epoch 2 / 5) train acc: 0.298000; val_acc: 0.260000
(Iteration 81 / 200) loss: 1.977227
(Iteration 91 / 200) loss: 2.007528
(Iteration 101 / 200) loss: 2.004762
(Iteration 111 / 200) loss: 1.885342
(Epoch 3 / 5) train acc: 0.343000; val_acc: 0.287000
(Iteration 121 / 200) loss: 1.891516
(Iteration 131 / 200) loss: 1.923677
(Iteration 141 / 200) loss: 1.957743
(Iteration 151 / 200) loss: 1.966736
(Epoch 4 / 5) train acc: 0.322000; val_acc: 0.305000
(Iteration 161 / 200) loss: 1.801483
(Iteration 171 / 200) loss: 1.973780
(Iteration 181 / 200) loss: 1.666573
(Iteration 191 / 200) loss: 1.909494
(Epoch 5 / 5) train acc: 0.372000; val_acc: 0.319000

```

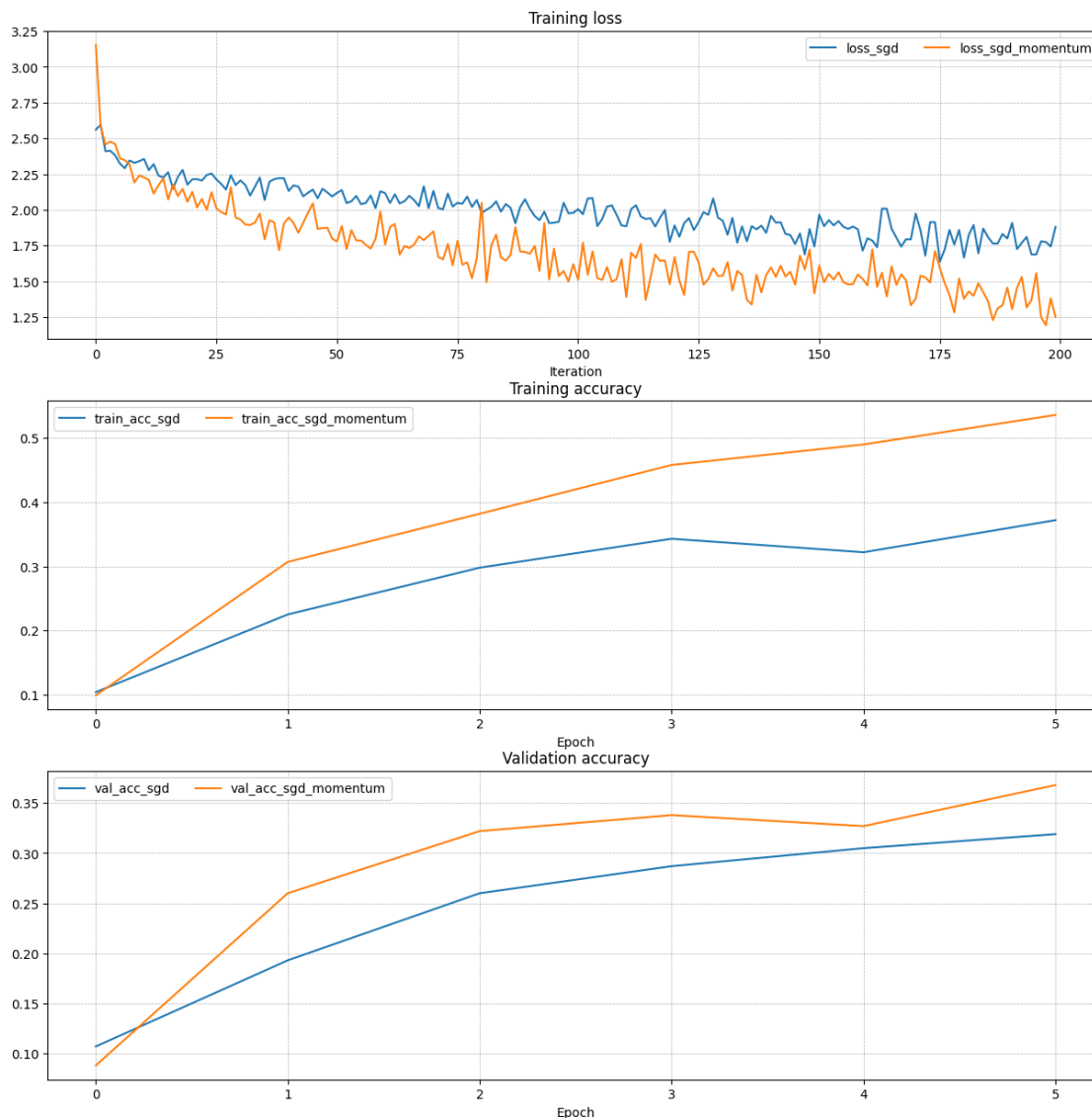
Running with `sgd_momentum`

```

(Iteration 1 / 200) loss: 3.153778
(Epoch 0 / 5) train acc: 0.099000; val_acc: 0.088000
(Iteration 11 / 200) loss: 2.227203
(Iteration 21 / 200) loss: 2.125706
(Iteration 31 / 200) loss: 1.932695
(Epoch 1 / 5) train acc: 0.307000; val_acc: 0.260000
(Iteration 41 / 200) loss: 1.946488
(Iteration 51 / 200) loss: 1.778583

```

(Iteration 61 / 200) loss: 1.758119
(Iteration 71 / 200) loss: 1.849137
(Epoch 2 / 5) train acc: 0.382000; val_acc: 0.322000
(Iteration 81 / 200) loss: 2.048671
(Iteration 91 / 200) loss: 1.693223
(Iteration 101 / 200) loss: 1.511693
(Iteration 111 / 200) loss: 1.390754
(Epoch 3 / 5) train acc: 0.458000; val_acc: 0.338000
(Iteration 121 / 200) loss: 1.670614
(Iteration 131 / 200) loss: 1.540271
(Iteration 141 / 200) loss: 1.597365
(Iteration 151 / 200) loss: 1.609851
(Epoch 4 / 5) train acc: 0.490000; val_acc: 0.327000
(Iteration 161 / 200) loss: 1.472687
(Iteration 171 / 200) loss: 1.378620
(Iteration 181 / 200) loss: 1.378175
(Iteration 191 / 200) loss: 1.305934
(Epoch 5 / 5) train acc: 0.536000; val_acc: 0.368000



2.2 RMSProp and Adam

RMSProp [1] and Adam [2] are update rules that set per-parameter learning rates by using a running average of the second moments of gradients.

In the file `icv83551/optim.py`, implement the RMSProp update rule in the `rmsprop` function and implement the Adam update rule in the `adam` function, and check your implementations using the tests below.

NOTE: Please implement the *complete* Adam update rule (with the bias correction mechanism), not the first simplified version mentioned in the course notes.

[1] Tijmen Tieleman and Geoffrey Hinton. “Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude.” COURSE: Neural Networks for Machine Learning 4 (2012).

[2] Diederik Kingma and Jimmy Ba, “Adam: A Method for Stochastic Optimization”, ICLR 2015.

```
[10]: # Test RMSProp implementation
from icv83551.optim import rmsprop

N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
cache = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-2, 'cache': cache}
next_w, _ = rmsprop(w, dw, config=config)

expected_next_w = np.asarray([
    [-0.39223849, -0.34037513, -0.28849239, -0.23659121, -0.18467247],
    [-0.132737,   -0.08078555, -0.02881884,  0.02316247,  0.07515774],
    [ 0.12716641,  0.17918792,  0.23122175,  0.28326742,  0.33532447],
    [ 0.38739248,  0.43947102,  0.49155973,  0.54365823,  0.59576619]])
expected_cache = np.asarray([
    [ 0.5976,      0.6126277,   0.6277108,   0.64284931,   0.65804321],
    [ 0.67329252,  0.68859723,   0.70395734,   0.71937285,   0.73484377],
    [ 0.75037008,  0.7659518,    0.78158892,   0.79728144,   0.81302936],
    [ 0.82883269,  0.84469141,   0.86060554,   0.87657507,   0.8926    ]])

# You should see relative errors around e-7 or less
print('next_w error: ', rel_error(expected_next_w, next_w))
print('cache error: ', rel_error(expected_cache, config['cache']))
```

next_w error: 9.524687511038133e-08

cache error: 2.6477955807156126e-09

```
[11]: # Test Adam implementation
from icv83551.optim import adam

N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
m = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
v = np.linspace(0.7, 0.5, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-2, 'm': m, 'v': v, 't': 5}
next_w, _ = adam(w, dw, config=config)

expected_next_w = np.asarray([
    [-0.40094747, -0.34836187, -0.29577703, -0.24319299, -0.19060977],
    [-0.1380274,  -0.08544591, -0.03286534,  0.01971428,  0.0722929],
    [ 0.1248705,   0.17744702,  0.23002243,  0.28259667,  0.33516969],
    [ 0.38774145,  0.44031188,  0.49288093,  0.54544852,  0.59801459]])
```

```

expected_v = np.asarray([
    [ 0.69966,      0.68908382,  0.67851319,  0.66794809,  0.65738853,],
    [ 0.64683452,  0.63628604,  0.6257431,   0.61520571,  0.60467385,],
    [ 0.59414753,  0.58362676,  0.57311152,  0.56260183,  0.55209767,],
    [ 0.54159906,  0.53110598,  0.52061845,  0.51013645,  0.49966,   ]])
expected_m = np.asarray([
    [ 0.48,          0.49947368,  0.51894737,  0.53842105,  0.55789474],
    [ 0.57736842,   0.59684211,  0.61631579,  0.63578947,  0.65526316],
    [ 0.67473684,   0.69421053,  0.71368421,  0.73315789,  0.75263158],
    [ 0.77210526,   0.79157895,  0.81105263,  0.83052632,  0.85         ]])

# You should see relative errors around e-7 or less
print('next_w error: ', rel_error(expected_next_w, next_w))
print('v error: ', rel_error(expected_v, config['v']))
print('m error: ', rel_error(expected_m, config['m']))

```

```

next_w error:  1.1395691798535431e-07
v error:  4.208314038113071e-09
m error:  4.214963193114416e-09

```

Once you have debugged your RMSProp and Adam implementations, run the following to train a pair of deep networks using these new update rules:

```

[12]: learning_rates = {'rmsprop': 1e-4, 'adam': 1e-3}
for update_rule in ['adam', 'rmsprop']:
    print('Running with ', update_rule)
    model = FullyConnectedNet(
        [100, 100, 100, 100, 100],
        weight_scale=5e-2
    )
    solver = Solver(
        model,
        small_data,
        num_epochs=5,
        batch_size=100,
        update_rule=update_rule,
        optim_config={'learning_rate': learning_rates[update_rule]},
        verbose=True
    )
    solvers[update_rule] = solver
    solver.train()
    print()

fig, axes = plt.subplots(3, 1, figsize=(15, 15))

axes[0].set_title('Training loss')
axes[0].set_xlabel('Iteration')
axes[1].set_title('Training accuracy')

```

```

axes[1].set_xlabel('Epoch')
axes[2].set_title('Validation accuracy')
axes[2].set_xlabel('Epoch')

for update_rule, solver in solvers.items():
    axes[0].plot(solver.loss_history, label=f"{update_rule}")
    axes[1].plot(solver.train_acc_history, label=f"{update_rule}")
    axes[2].plot(solver.val_acc_history, label=f"{update_rule}")

for ax in axes:
    ax.legend(loc='best', ncol=4)
    ax.grid(linestyle='--', linewidth=0.5)

plt.show()

```

Running with adam

```

(Iteration 1 / 200) loss: 3.476928
(Epoch 0 / 5) train acc: 0.126000; val_acc: 0.110000
(Iteration 11 / 200) loss: 2.027712
(Iteration 21 / 200) loss: 2.183357
(Iteration 31 / 200) loss: 1.744257
(Epoch 1 / 5) train acc: 0.363000; val_acc: 0.330000
(Iteration 41 / 200) loss: 1.707951
(Iteration 51 / 200) loss: 1.703835
(Iteration 61 / 200) loss: 2.094758
(Iteration 71 / 200) loss: 1.505558
(Epoch 2 / 5) train acc: 0.419000; val_acc: 0.362000
(Iteration 81 / 200) loss: 1.594431
(Iteration 91 / 200) loss: 1.511452
(Iteration 101 / 200) loss: 1.389237
(Iteration 111 / 200) loss: 1.463575
(Epoch 3 / 5) train acc: 0.497000; val_acc: 0.368000
(Iteration 121 / 200) loss: 1.231313
(Iteration 131 / 200) loss: 1.520198
(Iteration 141 / 200) loss: 1.363221
(Iteration 151 / 200) loss: 1.355143
(Epoch 4 / 5) train acc: 0.543000; val_acc: 0.347000
(Iteration 161 / 200) loss: 1.436402
(Iteration 171 / 200) loss: 1.231426
(Iteration 181 / 200) loss: 1.153575
(Iteration 191 / 200) loss: 1.209479
(Epoch 5 / 5) train acc: 0.619000; val_acc: 0.374000

```

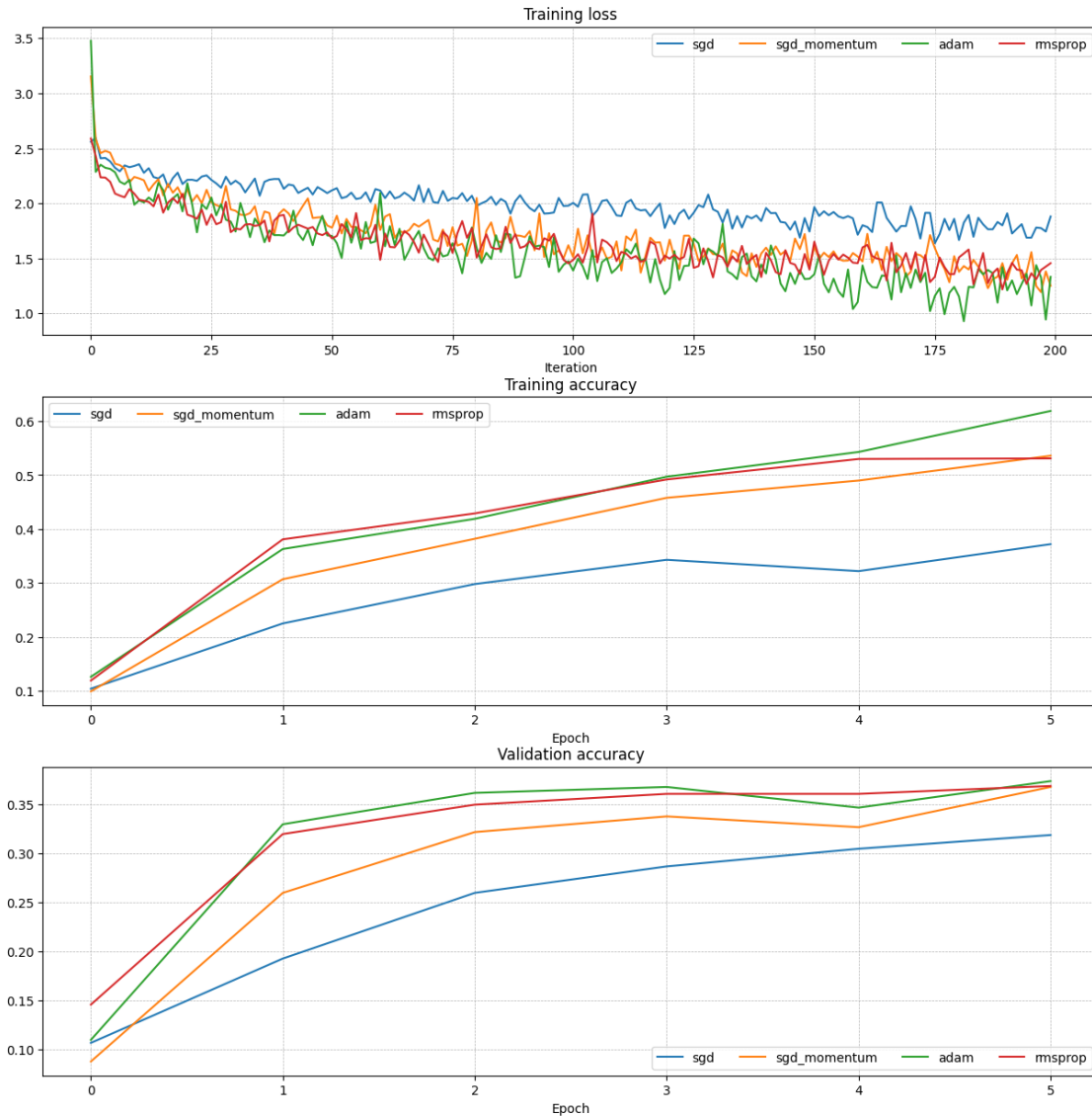
Running with rmsprop

```

(Iteration 1 / 200) loss: 2.589166
(Epoch 0 / 5) train acc: 0.119000; val_acc: 0.146000
(Iteration 11 / 200) loss: 2.032921

```

(Iteration 21 / 200) loss: 1.897277
(Iteration 31 / 200) loss: 1.770793
(Epoch 1 / 5) train acc: 0.381000; val_acc: 0.320000
(Iteration 41 / 200) loss: 1.895731
(Iteration 51 / 200) loss: 1.681091
(Iteration 61 / 200) loss: 1.487204
(Iteration 71 / 200) loss: 1.629973
(Epoch 2 / 5) train acc: 0.429000; val_acc: 0.350000
(Iteration 81 / 200) loss: 1.506686
(Iteration 91 / 200) loss: 1.610742
(Iteration 101 / 200) loss: 1.486124
(Iteration 111 / 200) loss: 1.559454
(Epoch 3 / 5) train acc: 0.492000; val_acc: 0.361000
(Iteration 121 / 200) loss: 1.497406
(Iteration 131 / 200) loss: 1.530736
(Iteration 141 / 200) loss: 1.550957
(Iteration 151 / 200) loss: 1.652046
(Epoch 4 / 5) train acc: 0.530000; val_acc: 0.361000
(Iteration 161 / 200) loss: 1.599574
(Iteration 171 / 200) loss: 1.401073
(Iteration 181 / 200) loss: 1.509365
(Iteration 191 / 200) loss: 1.365772
(Epoch 5 / 5) train acc: 0.531000; val_acc: 0.369000



2.3 Inline Question 2:

AdaGrad, like Adam, is a per-parameter optimization method that uses the following update rule:

```
cache += dw**2
w += - learning_rate * dw / (np.sqrt(cache) + eps)
```

John notices that when he was training a network with AdaGrad that the updates became very small, and that his network was learning slowly. Using your knowledge of the AdaGrad update rule, why do you think the updates would become very small? Would Adam have the same issue?

2.4 Answer:

AdaGrad updates becomes small because monotonically increasing cache. When we update the parameters we divide the gradient by a cache that keeps growing and produces low variables so the steps forward get smaller. Adam does not have this problem as it is using a decaying moving average and has a momentum update that affects the velocity.

3 Train a Good Model!

Train the best fully connected model that you can on CIFAR-10, storing your best model in the `best_model` variable. We require you to get at least 50% accuracy on the validation set using a fully connected network.

If you are careful it should be possible to get accuracies above 55%, but we don't require it for this part and won't assign extra credit for doing so. Later in the next assignment, we will ask you to train the best convolutional network that you can on CIFAR-10, and we would prefer that you spend your effort working on convolutional networks rather than fully connected networks.

Note: In the next assignment, you will learn techniques like BatchNormalization and Dropout which can help you train powerful models.

```
[13]: best_model = None

#####
# TODO: Train the best FullyConnectedNet that you can on CIFAR-10. You might #
# find batch/layer normalization and dropout useful. Store your best model in #
# the best_model variable.                                                    #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

best_val_acc = -1

learning_rates = [1e-3, 5e-4]
weight_scales = [1e-2, 5e-2]
regularization = [1e-3, 1e-4]

num_train = 500 # amount of training data for tests

# Generate small dataset
small_data = {
    "X_train": data["X_train"][:num_train],
    "y_train": data["y_train"][:num_train],
    "X_val": data["X_val"],
    "y_val": data["y_val"],
}

for lr in learning_rates:
    for ws in weight_scales:
        for reg in regularization:
```

```

print(f'lr {lr} ws {ws} reg {reg}')

# Initialize a deep network with Batchnorm and Dropout
model = FullyConnectedNet(
    hidden_dims=[256, 128, 64],
    reg=reg,
    weight_scale=ws,
    dtype=np.float64)

# Solver
solver = Solver(
    model, small_data,
    num_epochs=10,
    print_every=100,
    batch_size=128,
    update_rule='adam',
    optim_config={'learning_rate': lr},
    verbose=False)

solver.train()

# Track the best model based on validation accuracy
if solver.best_val_acc > best_val_acc:
    best_val_acc = solver.best_val_acc
    best_params = {'lr':lr, 'ws':ws, 'reg':reg}
    best_model = model

print(f'Best validation of small data accuracy: {best_val_acc}')

# Create solver to train the model on full dataset
solver = Solver(best_model, data,
    num_epochs=10, batch_size=128,
    update_rule='adam',
    optim_config={'learning_rate': best_params['lr']},
    verbose=False)

solver.train() # train the best model

# Print the final validation accuracy acquired with the best model
print(f'Best validation accuracy using full dataset: {solver.best_val_acc}')

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

```

lr 0.001 ws 0.01 reg 0.001

```
lr 0.001 ws 0.01 reg 0.0001
lr 0.001 ws 0.05 reg 0.001
lr 0.001 ws 0.05 reg 0.0001
lr 0.0005 ws 0.01 reg 0.001
lr 0.0005 ws 0.01 reg 0.0001
lr 0.0005 ws 0.05 reg 0.001
lr 0.0005 ws 0.05 reg 0.0001
Best validation of small data accuracy: 0.327
Best validation accuracy using full dataset: 0.539
```

4 Test Your Model!

Run your best model on the validation and test sets. You should achieve at least 50% accuracy on the validation set and the test set.

```
[15]: y_test_pred = np.argmax(best_model.loss(data['X_test']), axis=1)
      y_val_pred = np.argmax(best_model.loss(data['X_val']), axis=1)
      print('Validation set accuracy: ', (y_val_pred == data['y_val']).mean())
      print('Test set accuracy: ', (y_test_pred == data['y_test']).mean())
```

```
Validation set accuracy:  0.539
Test set accuracy:  0.528
```

```
[14]:
```