

CS7643: Deep Learning

Fall 2019

HW1 Solutions

James Hahn

September 26, 2019

1 Gradient Descent

1.

$$\begin{aligned} & \arg \min_w f(w^{(t)}) + \langle w - w^{(t)}, \nabla f(w^{(t)}) \rangle + \frac{\lambda}{2} \|w - w^{(t)}\|^2 \\ \implies & \frac{\delta}{\delta w} = \nabla_w f(w^{(t)}) + \lambda(w - w^{(t)}) \\ \implies & \nabla_w f(w^{(t)}) + \lambda(w - w^{(t)}) = 0 \\ \implies & \lambda w = \lambda w^{(t)} - \nabla_w f(w^{(t)}) \\ \implies & w^* = w^{(t)} - \frac{1}{\lambda} \nabla_w f(w^{(t)}) \end{aligned}$$

This tells us the gradient descent update rule helps find the solution to the true, underlying weights by approximating a function, rather than solving the optimization problem through direct differentiation. The relationship between λ and η is that $\frac{1}{\lambda} = \eta$. If we see λ to be a regularization parameter, this could indicate the lower the λ (regularization parameter), the higher the overall regularization on the weights. This occurs because a lower λ produces higher $\frac{1}{\lambda}$, placing more weight on per-iteration weight differences, leading to a slower learning/approximation process for our weights.

2.

We want to show $\sum_{t=1}^T \langle w^{(t)} - w^*, v_t \rangle \leq \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2$. As a preliminary step, let's simplify the problem. We will prove a lemma based on the update rule provided to us.

Lemma 1.1. $w^{(t+1)} = -\eta \sum_{i=1}^t v_t$

We are given $w^{(1)} = 0$. For a proof by induction, the base case for $t = 2$ is as follows: $w^{(2)} = w^{(1)} - \eta v_1 = -\eta v_1 = -\eta \sum_{i=1}^1 v_t$, which is the trivial base case. For a non-trivial base case, such as $t = 3$, we see the following: $w^{(3)} = w^{(2)} - \eta v_1 = -\eta v_1 - \eta v_2 = -\eta \sum_{i=1}^2 v_t$. As such, we will show this follows for future t . Assume we have some $w^{(t+1)} = -\eta \sum_{i=1}^t v_t$. Now, $w^{(t+2)} = w^{(t+1)} - \eta v_{t+1} = -\eta \sum_{i=1}^t v_t - \eta v_{t+1} = -\eta[(\sum_{i=1}^t v_t) + v_{t+1}] = -\eta \sum_{i=1}^{t+1} v_t$. As such, the lemma holds. \square

First, we decompose the LHS of the inequality. Namely, $\sum_{t=1}^T \langle w^{(t)} - w^*, v_t \rangle = \sum_{t=1}^T \langle w^{(t)}, v_t \rangle - \sum_{t=1}^T \langle w^*, v_t \rangle = \sum_{t=1}^T \langle w^{(t)}, v_t \rangle - \langle w^*, \sum_{t=1}^T v_t \rangle$.

Second, by utilizing Lemma 1.1, we observe the following transformation of the first term on the RHS of the above equation:

$$\begin{aligned} \text{RHS first term} &= \sum_{t=1}^T \langle w^{(t)} - w^*, v_t \rangle \\ &= \sum_{t=1}^T \langle -\eta \sum_{u=1}^{t-1} v_u, v_t \rangle \\ &= -\eta \sum_{t=1}^T \langle \sum_{u=1}^{t-1} v_u, v_t \rangle \\ &= -\eta \left(\frac{\sum_{t=1}^T v_t \cdot \sum_{t=1}^T v_t}{2} - \frac{\sum_{t=1}^T \|v_t\|^2}{2} \right) \\ &= -\frac{\eta}{2} \|\sum_{t=1}^T v_t\|^2 + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2 \\ &= -\frac{1}{2\eta} \|w^{(t+1)}\|^2 + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2 \end{aligned}$$

Third, we simplify the second term on the RHS of the first equation with the use of Lemma 1.1:

$$\begin{aligned} \text{RHS second term} &= \langle w^*, \sum_{t=1}^T v_t \rangle \\ &= \langle w^*, -\frac{1}{\eta} w^{(t+1)} \rangle \quad (\text{from Lemma 1.1: } w^{(t+1)} = -\eta \sum_{i=1}^t v_t \implies \sum_{i=1}^t v_t = -\frac{1}{\eta} w^{(t+1)}) \\ &= -\frac{1}{\eta} \langle w^*, w^{(t+1)} \rangle \end{aligned}$$

So, we have finally expanded the LHS of the inequality:

$$\text{LHS} = -\frac{1}{2\eta} \|w^{(t+1)}\|^2 + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2 + \frac{1}{\eta} \langle w^*, w^{(t+1)} \rangle$$

The key thing to note here is that $\frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2$ already appears on the given RHS of the inequality, so they cancel each other out. To simplify the problem and to further prove the inequality, all we need to show is $-\frac{1}{2\eta} \|w^{(t+1)}\|^2 + \frac{1}{\eta} \langle w^*, w^{(t+1)} \rangle \leq \frac{\|w^*\|^2}{2\eta}$.

First, we see $-\frac{1}{2\eta} \|w^{(t+1)}\|^2$ is negative since it's the product of a negative fraction and a positive distance. Second, we see $\frac{1}{\eta} \langle w^*, w^{(t+1)} \rangle = \frac{1}{2\eta} \langle w^*, w^{(t+1)} \rangle + \frac{1}{2\eta} \langle w^*, w^{(t+1)} \rangle$. Third, $\langle w^*, w^{(t+1)} \rangle \leq \langle w^{(t+1)}, w^{(t+1)} \rangle = \|w^{(t+1)}\|^2$ and $\langle w^*, w^{(t+1)} \rangle \leq \langle w^*, w^* \rangle$. Therefore, with these observed properties, we know $-\frac{1}{2\eta} \|w^{(t+1)}\|^2 \leq -\frac{1}{2\eta} \langle w^*, w^{(t+1)} \rangle \implies -\frac{1}{2\eta} \|w^{(t+1)}\|^2 + \frac{1}{2\eta} \langle w^*, w^{(t+1)} \rangle \leq 0$. Additionally, $\frac{1}{2\eta} \langle w^*, w^{(t+1)} \rangle \leq \frac{1}{2\eta} \|w^*\|^2$.

Finally, since we know $-\frac{1}{2\eta} \|w^{(t+1)}\|^2 + \frac{1}{2\eta} \langle w^*, w^{(t+1)} \rangle \leq 0$ and $\frac{1}{2\eta} \langle w^*, w^{(t+1)} \rangle \leq \frac{1}{2\eta} \|w^*\|^2$, we get the summation of a negative term and term that's less than the RHS of the inequality we're trying to prove, so we observe $-\frac{1}{2\eta} \|w^{(t+1)}\|^2 + \frac{1}{\eta} \langle w^*, w^{(t+1)} \rangle \leq \frac{\|w^*\|^2}{2\eta}$. So, we have proved the simpler inequality we established above.

$$\therefore \sum_{t=1}^T \langle w^{(t)} - w^*, v_t \rangle \leq \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2 \quad \square$$

3.

We must show for $\bar{w} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$, we get $f(\bar{w}) - f(w^*) \leq \frac{1}{T} \sum_{t=1}^T \langle w^{(t)} - w^*, \nabla f(w^{(t)}) \rangle$. In the problem statement for question 1, we were told “When f is convex, this approximation forms a lower bound of f , referring to equation 2. As such, if we assume f is convex, which is reasonable since w^* minimizes f , we get

$$\begin{aligned} f(w^*) &\geq f(w^{(t)}) + \langle w^* - w^{(t)}, \nabla f(w^{(t)}) \rangle \\ \implies f(w^*) - f(w^{(t)}) &\geq \langle w^* - w^{(t)}, \nabla f(w^{(t)}) \rangle \\ \implies f(w^{(t)}) - f(w^*) &\leq \langle w^{(t)} - w^*, \nabla f(w^{(t)}) \rangle \\ \implies \frac{1}{T} \sum_{t=1}^T [f(w^{(t)}) - f(w^*)] &\leq \frac{1}{T} \sum_{t=1}^T \langle w^{(t)} - w^*, \nabla f(w^{(t)}) \rangle \\ \implies \frac{1}{T} \sum_{t=1}^T [f(w^{(t)})] - f(w^*) &\leq \frac{1}{T} \sum_{t=1}^T \langle w^{(t)} - w^*, \nabla f(w^{(t)}) \rangle \end{aligned}$$

The definition of convexity states if $f(x)$ is convex, then $f(tx) \leq t f(x)$ for some variable or function t . So, we know $f(\bar{w}) = f\left(\frac{1}{T} \sum_{t=1}^T w^{(t)}\right) \leq \frac{1}{T} \sum_{t=1}^T f(w^{(t)})$

$$\implies f(\bar{w}) - f(w^*) \leq \frac{1}{T} \sum_{t=1}^T \langle w^{(t)} - w^*, \nabla f(w^{(t)}) \rangle$$

From question 2, we proved $\sum_{t=1}^T \langle w^{(t)} - w^*, v_t \rangle \leq \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|v_t\|^2$, so we get:

$$\begin{aligned} \implies f(\bar{w}) - f(w^*) &\leq \frac{1}{T} \left[\frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\nabla f(w^{(t)})\|^2 \right] \\ \implies f(\bar{w}) - f(w^*) &\leq \frac{\|w^*\|^2}{2T\eta} + \frac{\eta}{2T} \sum_{t=1}^T \|\nabla f(w^{(t)})\|^2 \end{aligned}$$

Now, if we just treat B and ρ as upper bounds for $\|w\|^2$ and $\|\nabla f(w^{(t)})\|^2$ respectively, as given in the problem, we get

$$\begin{aligned} \implies f(\bar{w}) - f(w^*) &\leq \frac{B^2}{2T\eta} + \frac{\eta}{2T} \sum_{t=1}^T \rho^2 \\ \implies f(\bar{w}) - f(w^*) &\leq \frac{B^2}{2T\eta} + \frac{\eta T \rho^2}{2T} \end{aligned}$$

With simple algebra, the T cancels out on the second term of the RHS, and we can substitute $\eta = \sqrt{\frac{B^2}{\rho^2 T}}$, which is given to us, into the equation:

$$\begin{aligned} \implies f(\bar{w}) - f(w^*) &\leq \frac{B^2}{2T\sqrt{\frac{B^2}{\rho^2 T}}} + \frac{\sqrt{\frac{B^2}{\rho^2 T}} \rho^2}{2} \\ \implies f(\bar{w}) - f(w^*) &\leq \frac{B^2}{2T\frac{B}{\rho}\sqrt{\frac{1}{T}}} + \frac{\frac{B}{\rho} \rho^2 \sqrt{\frac{1}{T}}}{2} \\ \implies f(\bar{w}) - f(w^*) &\leq \frac{B\rho}{2\sqrt{T}} + \frac{B\rho}{2\sqrt{T}} \\ \implies f(\bar{w}) - f(w^*) &\leq \frac{B\rho}{\sqrt{T}} \end{aligned}$$

As shown above, the convergence rate of w^* is $\mathcal{O}(\frac{1}{\sqrt{T}})$ i.e. upper bound of $f(\bar{w}) - f(w^*) \propto \frac{1}{\sqrt{T}}$. \square

4.

No. We will show this with a simple counter-example. We were given $w^{(1)} = 0$ in this problem. Additionally, we are given the objective function $f(w) = \frac{1}{2}(w - 2)^2 + \frac{1}{2}(w + 1)^2$. We want to show the objective function does not always decrease after each iteration, assuming the gradient of the objective function at each iteration is randomly selected from the gradients of the two sub-terms of the objective function. As such, we either have $f'_1(w) = w - 2$ or $f'_2(w) = w + 1$ with equal probability. We use the update rule $w^{(t+1)} = w^{(t)} - \eta f'(w)$.

First, we see $f(w^{(1)}) = \frac{1}{2}(0 - 2)^2 + \frac{1}{2}(0 + 1)^2 = 2 + 0.5 = 2.5$. So, we see the gradient is

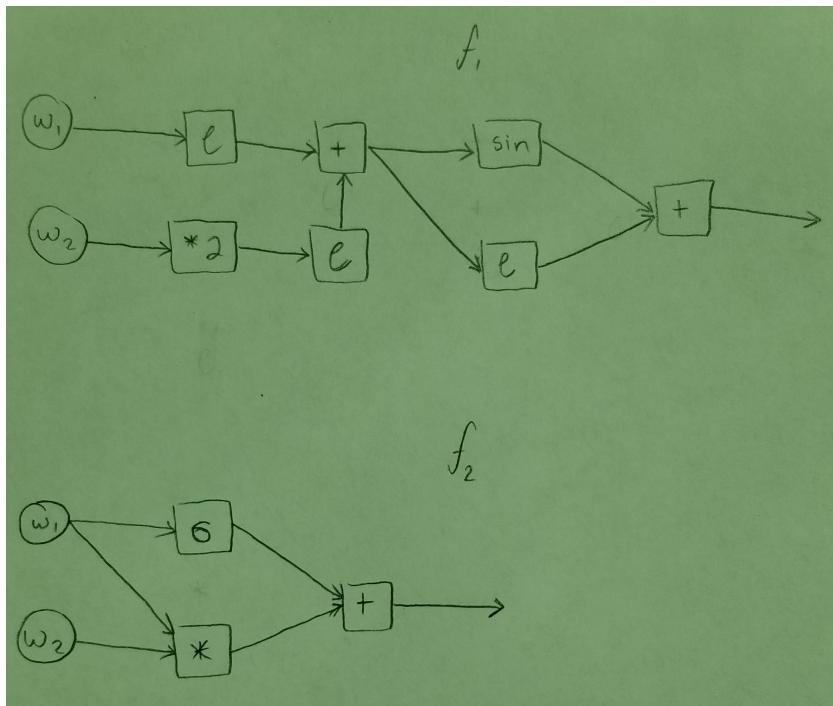
$$f'(w) = \begin{cases} -2 & \text{if } f'_1(w) \\ 1 & \text{if } f'_2(w) \end{cases}$$

Let's use $f'_1(w)$ first. We see $w_1^{(2)} = w^{(1)} - \eta f'_1(w) = 0 - \eta(-2) = 2\eta$. As such, $f(w_1^{(2)}) = \frac{1}{2}(2\eta - 2)^2 + \frac{1}{2}(2\eta + 1)^2 = \frac{1}{2}(4\eta^2 - 8\eta + 4) + \frac{1}{2}(4\eta^2 + 4\eta + 1) = 4\eta^2 - 4\eta + 5$. We know $\eta > 0$, as provided in the problem description, so all we can infer is the new loss function is less than 5, but we cannot assume it increased.

Next, let's observe $f'_2(w)$. We see $w_2^{(2)} = w^{(1)} - \eta f'_2(w) = 0 - \eta(1) = -\eta$. As such, $f(w_2^{(2)}) = \frac{1}{2}(-\eta - 2)^2 + \frac{1}{2}(-\eta + 1)^2 = \frac{1}{2}(\eta^2 + 4\eta + 4) + \frac{1}{2}(\eta^2 - 2\eta + 1) = \eta^2 + \eta + 2.5$. Since $\eta > 0$, we know for sure the new value of the objective function is greater than 2.5. Therefore, the value has increased.
 \therefore We cannot guarantee the loss function will decrease.

2 Automatic Differentiation

5(a).



$$f_1(w) = f_1(1, 2) = e^{e^{w_1} + e^{2w_2}} + \sin(e^{w_1} + e^{2w_2}) = e^{e^1 + e^4} + \sin(e^1 + e^4) = 7.8 \times 10^{24}$$

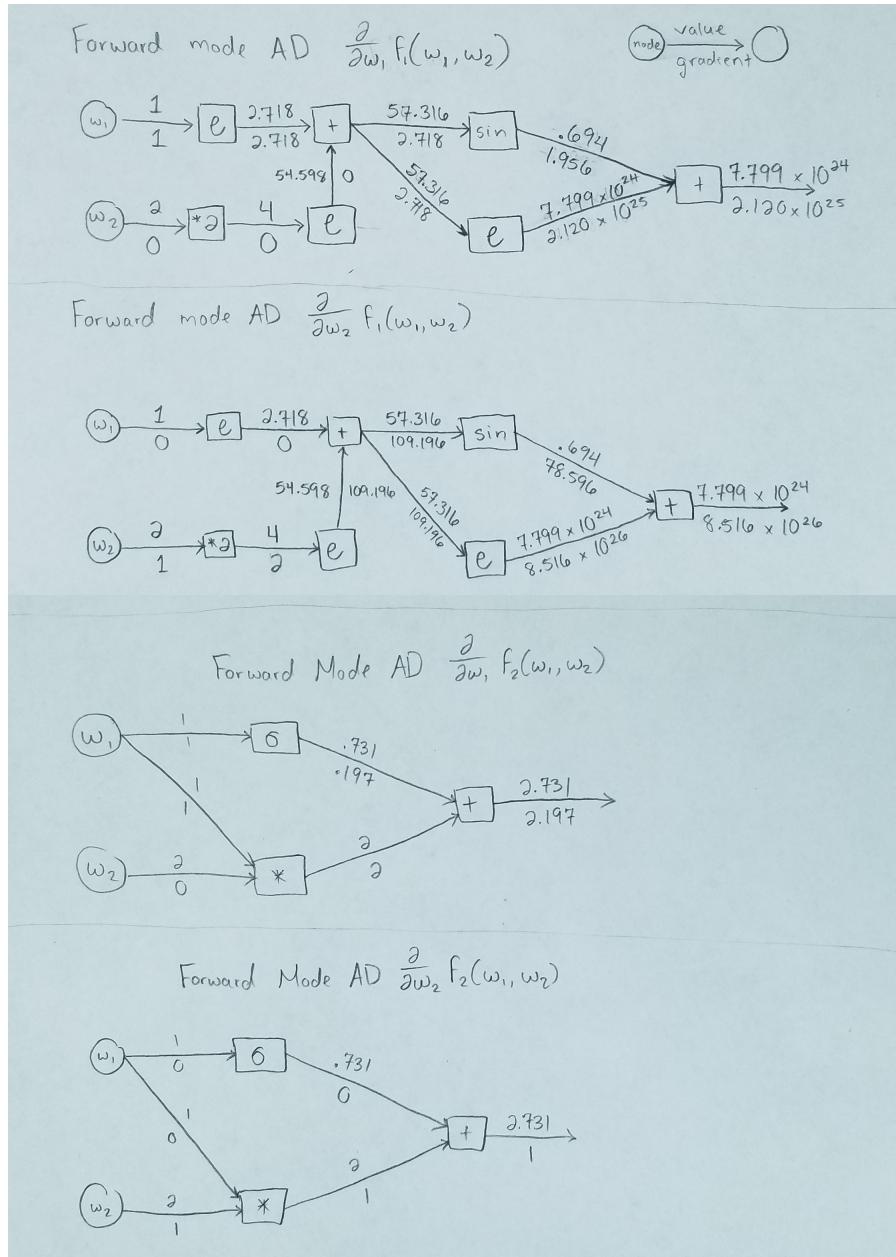
$$f_2(w) = f_2(1, 2) = w_1 w_2 + \sigma(w_1) = (1)(2) + \sigma(1) = 2 + \frac{1}{1+e^{-1}} = 2 + 0.731 = 2.731$$

5(b).

$$\begin{aligned}
 \frac{f_1(w_1+0.01, w_2) - f_1(w_1, w_2)}{0.01} &= \frac{f_1(1.01, 2) - f_1(1, 2)}{0.01} = \frac{(8.018 \times 10^{24}) - (7.802 \times 10^{24})}{0.01} = \frac{2.162 \times 10^{23}}{0.01} = 2.162 \times 10^{25} \\
 \frac{f_1(w_1, w_2+0.01) - f_1(w_1, w_2)}{0.01} &= \frac{f_1(1, 2.01) - f_1(1, 2)}{0.01} = \frac{(2.351 \times 10^{25}) - (7.802 \times 10^{24})}{0.01} = \frac{1.571 \times 10^{25}}{0.01} = 1.571 \times 10^{27} \\
 \frac{f_2(w_1+0.01, w_2) - f_2(w_1, w_2)}{0.01} &= \frac{f_2(1.01, 2) - f_2(1, 2)}{0.01} = \frac{2.753 - 2.731}{0.01} = \frac{0.022}{0.01} = 2.2 \\
 \frac{f_2(w_1, w_2+0.01) - f_2(w_1, w_2)}{0.01} &= \frac{f_2(1, 2.01) - f_2(1, 2)}{0.01} = \frac{2.741 - 2.731}{0.01} = \frac{0.01}{0.01} = 1
 \end{aligned}$$

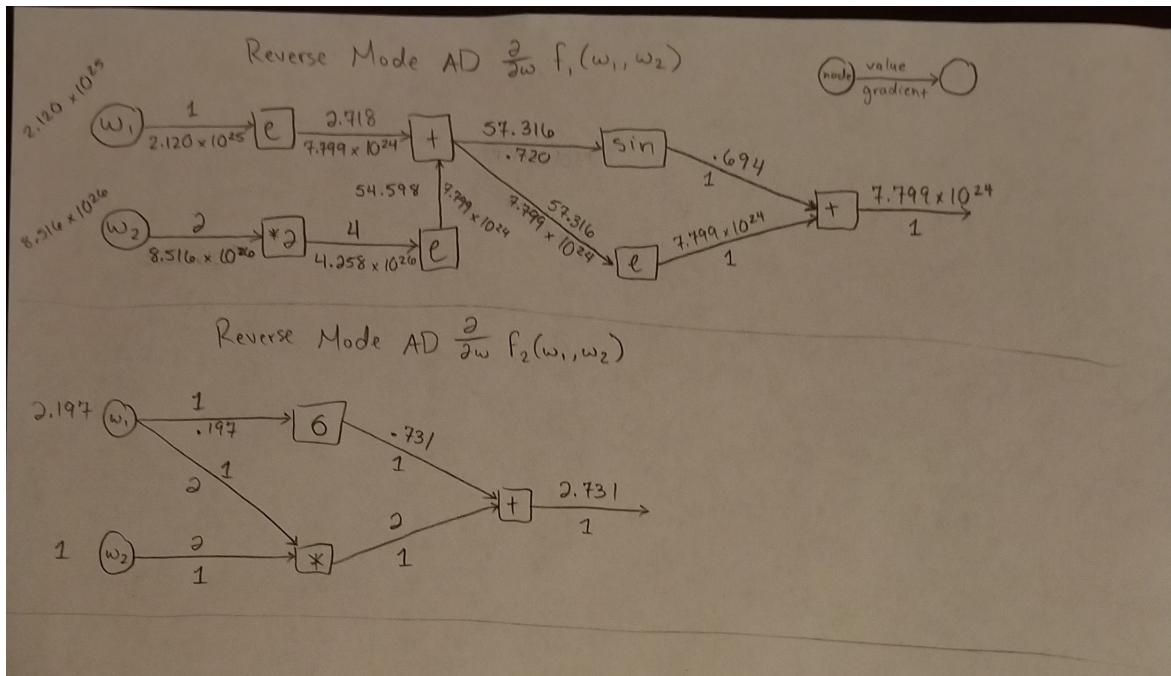
$$J = \begin{bmatrix} 2.162 \times 10^{25} & 1.571 \times 10^{27} \\ 2.2 & 1 \end{bmatrix}$$

5(c).



$$J = \begin{bmatrix} 2.120 \times 10^{25} & 8.516 \times 10^{26} \\ 2.197 & 1 \end{bmatrix}$$

5(d).



$$J = \begin{bmatrix} 2.120 \times 10^{25} & 8.516 \times 10^{26} \\ 2.197 & 1 \end{bmatrix}$$

$5(e)$.

This is torture. Yes.

3 Directed Acyclic Graphs (DAG)

6.

We want to show if a graph G is a DAG, it has a topological ordering.

Since G is a DAG, $\exists v_i \in V_G$ s.t. v_i has no incoming cycle.

Now, let $V' = \emptyset$ be an empty queue of vertices.

Remove v_i from V_G and add it to V' such that $V' = \{v_i\}$ and $G = G - v_i$.

G is still a DAG since removing a vertex does not add any edges to a graph.

Since G is still a DAG, $\exists v_j \in V_G$ s.t. v_j has no incoming cycles.

Remove v_j from V_G and add it to V' such that $V' = \{v_i, v_j\}$ and $G = G - v_j$.

G is still a DAG from the aforementioned logic.

We can repeat this process until $V_G = \emptyset$ and $|V'| = n$, indicating all vertices have been placed in this queue.

At this point, G is still a DAG because an empty set is a DAG by definition.

By removing the vertices of V' in first-in-first-out order, we can easily retrieve a topological ordering of G because each vertex v_m removed from the queue removes some edge (v_m, v_k) , allowing v_k to be removed. This enforces an ordering of the vertices such that in order for a vertex v_k to be removed, all vertices (e.g. v_m) directed into it must be removed first.

As such, if we enforce $m < k$, we have that for all edges (v_m, v_k) we have $m < k$, which is the definition of a topological ordering.

\therefore If G is a DAG, it has a topological ordering.

□

7.

We want to show that if G has a topological ordering, it is a DAG.

Since G has a topological ordering, we know that for all edges $(v_i, v_j) \in E$, we have $i < j$.

Assume some cycle C exists in G . Let v_i be the vertex with the lowest number in the cycle.

The vertex directly preceding v_i , let us name it v_j , forms an edge (v_j, v_i) .

★ Since v_i is the lowest numbered vertex in the cycle, we know $i < j$.

However, in order to enforce the topological ordering assumption, since (v_j, v_i) exists, then j must be less than i ($j < i$).

This directly contradicts ★, so the cycle produces a contradiction of the topological ordering of G .

As such, G must have no cycles (i.e. it is acyclic).

∴ G is a DAG. □

softmax

September 26, 2019

1 Softmax Classifier

This exercise guides you through the process of classifying images using a Softmax classifier. As part of this you will:

- Implement a fully vectorized loss function for the Softmax classifier
- Calculate the analytical gradient using vectorized code
- Tune hyperparameters on a validation set
- Optimize the loss function with Stochastic Gradient Descent (SGD)
- Visualize the learned weights

```
In [1]: # start-up code!
import random

import matplotlib.pyplot as plt
import numpy as np

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

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

In [2]: from load_cifar10_tvt import load_cifar10_train_val

X_train, y_train, X_val, y_val, X_test, y_test = load_cifar10_train_val()
print("Train data shape: ", X_train.shape)
print("Train labels shape: ", y_train.shape)
print("Val data shape: ", X_val.shape)
print("Val labels shape: ", y_val.shape)
print("Test data shape: ", X_test.shape)
print("Test labels shape: ", y_test.shape)
```

```

Train, validation and testing sets have been created as
X_i and y_i where i=train,val,test
Train data shape: (3073, 49000)
Train labels shape: (49000,)
Val data shape: (3073, 1000)
Val labels shape: (1000,)
Test data shape: (3073, 1000)
Test labels shape: (1000,)

```

Code for this section is to be written in cs231n/classifiers/softmax.py

In [9]: # Now, implement the vectorized version in softmax_loss_vectorized.

```

import time

from cs231n.classifiers.softmax import softmax_loss_vectorized

# gradient check.
from cs231n.gradient_check import grad_check_sparse

W = np.random.randn(10, 3073) * 0.0001

tic = time.time()
loss, grad = softmax_loss_vectorized(W, X_train, y_train, 0.00001)
toc = time.time()
print("vectorized loss: %e computed in %fs" % (loss, toc - tic))

# 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)))

f = lambda w: softmax_loss_vectorized(w, X_train, y_train, 0.0)[0]
grad_numerical = grad_check_sparse(f, W, grad, 10)

vectorized loss: 2.358837e+00 computed in 0.625360s
loss: 2.358837
sanity check: 2.302585
numerical: 1.897609 analytic: 1.897361, relative error: 6.550134e-05
numerical: -1.387101 analytic: -1.387062, relative error: 1.401846e-05
numerical: 3.547049 analytic: 3.546799, relative error: 3.526429e-05
numerical: 1.427716 analytic: 1.427785, relative error: 2.413449e-05
numerical: 1.067347 analytic: 1.066835, relative error: 2.399849e-04
numerical: 0.339480 analytic: 0.339402, relative error: 1.148107e-04
numerical: 1.035761 analytic: 1.035384, relative error: 1.821086e-04
numerical: -1.030169 analytic: -1.030336, relative error: 8.081725e-05
numerical: -2.182639 analytic: -2.182752, relative error: 2.577810e-05
numerical: 1.464362 analytic: 1.463994, relative error: 1.256691e-04

```

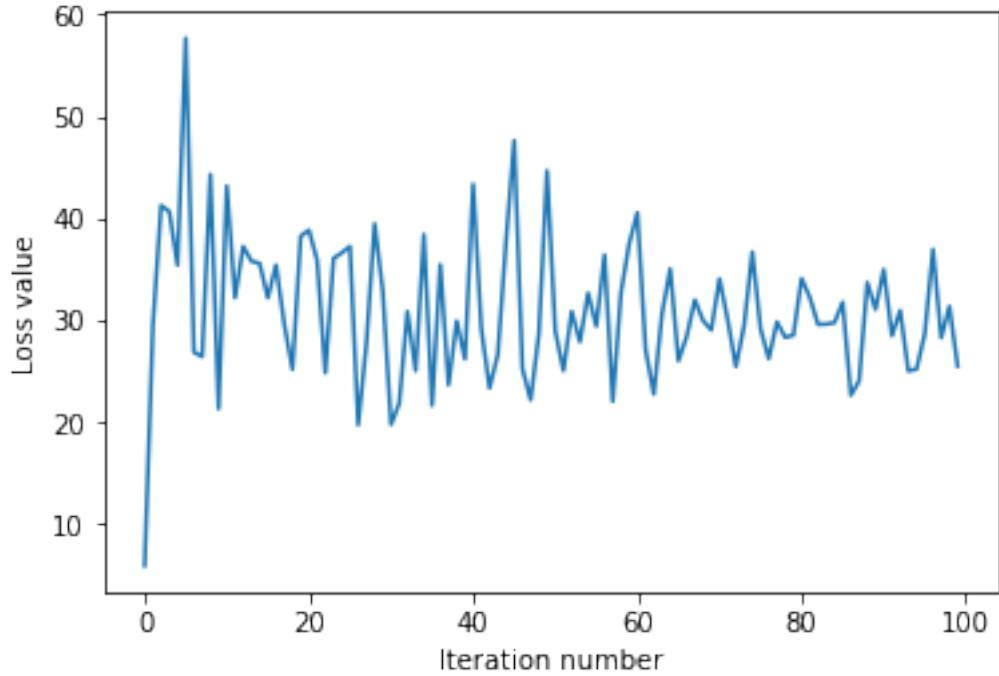
Code for this section is to be written in `cs231n/classifiers/linear_classifier.py`

```
In [65]: # Now that efficient implementations to calculate loss function and gradient of the s
# use it to train the classifier on the cifar-10 data
# Complete the `train` function in cs231n/classifiers/linear_classifier.py

from cs231n.classifiers.linear_classifier import Softmax

classifier = Softmax()
loss_hist = classifier.train(
    X_train,
    y_train,
    learning_rate=1e-4,
    reg=1e-4,
    num_iters=100,
    batch_size=2000,
    verbose=False,
)
# Plot loss vs. iterations
plt.plot(loss_hist)
plt.xlabel("Iteration number")
plt.ylabel("Loss value")
```

Out[65]: `Text(0,0.5,'Loss value')`



```
In [66]: # Complete the `predict` function in cs231n/classifiers/linear_classifier.py
# Evaluate on test set
y_test_pred = classifier.predict(X_test)
test_accuracy = np.mean(y_test == y_test_pred)
print("softmax on raw pixels final test set accuracy: %f" % (test_accuracy,))
```

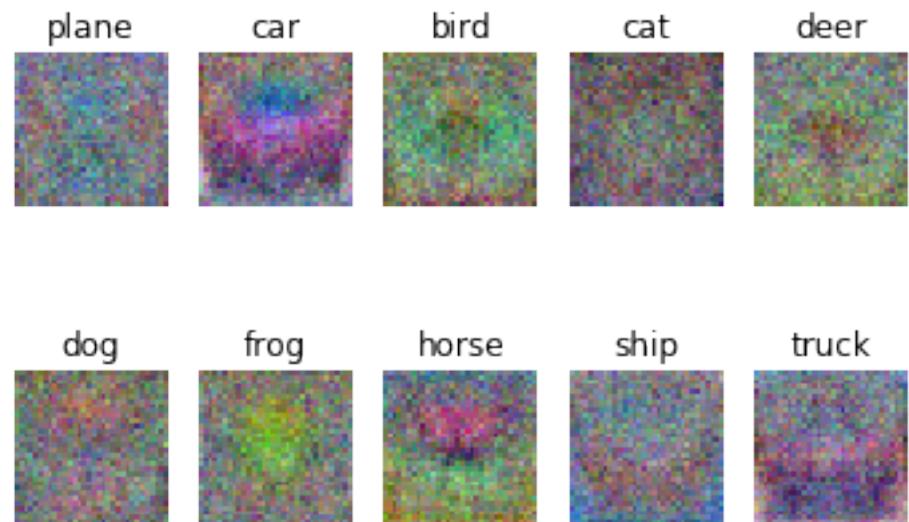
```
softmax on raw pixels final test set accuracy: 0.268000
```

```
In [67]: # Visualize the learned weights for each class
w = classifier.W[:, :-1] # strip out the bias
w = w.reshape(10, 32, 32, 3)

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])
```



In []:

two_layer_net

September 26, 2019

1 Implementing a Neural Network

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

In [47]: # A bit of setup

```
import numpy as np
import matplotlib.pyplot as plt

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

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

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))))
```

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

```
%reload_ext autoreload
```

The neural network parameters will be stored in a dictionary (`model` below), where the keys are the parameter names and the values are numpy arrays. Below, we initialize toy data and a toy model that we will use to verify your implementations.

In [48]: # Create some toy data to check your implementations

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

```

def init_toy_model():
    model = {}
    model['W1'] = np.linspace(-0.2, 0.6, num=input_size*hidden_size).reshape(input_size,
    model['b1'] = np.linspace(-0.3, 0.7, num=hidden_size)
    model['W2'] = np.linspace(-0.4, 0.1, num=hidden_size*num_classes).reshape(hidden_size,
    model['b2'] = np.linspace(-0.5, 0.9, num=num_classes)
    return model

def init_toy_data():
    X = np.linspace(-0.2, 0.5, num=num_inputs*input_size).reshape(num_inputs, input_size)
    y = np.array([0, 1, 2, 2, 1])
    return X, y

model = init_toy_model()
X, y = init_toy_data()

```

2 Forward pass: compute scores

Open the file cs231n/classifiers/neural_net.py and look at the function two_layer_net. This function is very similar to the loss functions you have written for the Softmax exercise in HW0: It takes the data and weights and computes the class scores, the loss, and the gradients on the parameters.

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

```
In [49]: from cs231n.classifiers.neural_net import two_layer_net

scores = two_layer_net(X, model)
print(scores)
correct_scores = [[-0.5328368, 0.20031504, 0.93346689],
                  [-0.59412164, 0.15498488, 0.9040914 ],
                  [-0.67658362, 0.08978957, 0.85616275],
                  [-0.77092643, 0.01339997, 0.79772637],
                  [-0.89110401, -0.08754544, 0.71601312]]

# the difference should be very small. We get 3e-8
print('Difference between your scores and correct scores:')
print(np.sum(np.abs(scores - correct_scores)))

[[-0.5328368  0.20031504  0.93346689]
 [-0.59412164  0.15498488  0.9040914 ]
 [-0.67658362  0.08978957  0.85616275]
 [-0.77092643  0.01339997  0.79772637]
 [-0.89110401 -0.08754544  0.71601312]]
Difference between your scores and correct scores:
3.848682303062012e-08
```

3 Forward pass: compute loss

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

```
In [50]: reg = 0.1
        loss, _ = two_layer_net(X, model, y, reg)
        correct_loss = 1.38191946092

        # should be very small, we get 5e-12
        print('Difference between your loss and correct loss:')
        print(np.sum(np.abs(loss - correct_loss)))
```

Difference between your loss and correct loss:

4.676925513535934e-12

4 Backward pass

Implement the rest of the function. This will compute the gradient of the loss with respect to the variables W_1 , b_1 , W_2 , and b_2 . Now that you (hopefully!) have a correctly implemented forward pass, you can debug your backward pass using a numeric gradient check:

```
In [54]: from cs231n.gradient_check import eval_numerical_gradient

        # Use numeric gradient checking to check your implementation of the backward pass.
        # If your implementation is correct, the difference between the numeric and
        # analytic gradients should be less than 1e-8 for each of  $W_1$ ,  $W_2$ ,  $b_1$ , and  $b_2$ .

        loss, grads = two_layer_net(X, model, y, reg)

        # these should all be less than 1e-8 or so
        for param_name in grads:
            param_grad_num = eval_numerical_gradient(lambda W: two_layer_net(X, model, y, reg)[0], W, -grads[param_name])
            print('%s max relative error: %e' % (param_name, rel_error(param_grad_num, grads[param_name])))
```

W2 max relative error: 1.401431e-09
b2 max relative error: 7.311059e-11
W1 max relative error: 4.426512e-09
b1 max relative error: 2.746125e-08

5 Train the network

To train the network we will use SGD with Momentum. Last assignment you implemented vanilla SGD. You will now implement the momentum update and the RMSProp update. Open the file `classifier_trainer.py` and familiarize yourself with the `ClassifierTrainer` class. It performs optimization given an arbitrary cost function `data`, and `model`. By default it uses vanilla SGD,

which we have already implemented for you. First, run the optimization below using Vanilla SGD:

```
In [55]: from cs231n.classifier_trainer import ClassifierTrainer

model = init_toy_model()
trainer = ClassifierTrainer()
# call the trainer to optimize the loss
# Notice that we're using sample_batches=False, so we're performing Gradient Descent
best_model, loss_history, _, _ = trainer.train(X, y, X, y,
                                                model, two_layer_net,
                                                reg=0.001,
                                                learning_rate=1e-1, momentum=0.0, learning_rate=1e-1,
                                                update='sgd', sample_batches=False,
                                                num_epochs=100,
                                                verbose=False)
print('Final loss with vanilla SGD: %f' % (loss_history[-1],))

starting iteration 0
starting iteration 10
starting iteration 20
starting iteration 30
starting iteration 40
starting iteration 50
starting iteration 60
starting iteration 70
starting iteration 80
starting iteration 90
Final loss with vanilla SGD: 0.940686
```

Now fill in the **momentum update** in the first missing code block inside the train function, and run the same optimization as above but with the momentum update. You should see a much better result in the final obtained loss:

```
In [56]: model = init_toy_model()
trainer = ClassifierTrainer()
# call the trainer to optimize the loss
# Notice that we're using sample_batches=False, so we're performing Gradient Descent
best_model, loss_history, _, _ = trainer.train(X, y, X, y,
                                                model, two_layer_net,
                                                reg=0.001,
                                                learning_rate=1e-1, momentum=0.9, learning_rate=1e-1,
                                                update='momentum', sample_batches=False,
                                                num_epochs=100,
                                                verbose=False)

correct_loss = 0.494394
print('Final loss with momentum SGD: %f. We get: %f' % (loss_history[-1], correct_loss))
```

```

starting iteration 0
starting iteration 10
starting iteration 20
starting iteration 30
starting iteration 40
starting iteration 50
starting iteration 60
starting iteration 70
starting iteration 80
starting iteration 90
Final loss with momentum SGD: 0.494394. We get: 0.494394

```

The **RMSProp** update step is given as follows:

```

cache = decay_rate * cache + (1 - decay_rate) * dx**2
x += - learning_rate * dx / np.sqrt(cache + 1e-8)

```

Here, `decay_rate` is a hyperparameter and typical values are [0.9, 0.99, 0.999].

Implement the **RMSProp** update rule inside the `train` function and rerun the optimization:

```

In [57]: model = init_toy_model()
          trainer = ClassifierTrainer()
          # call the trainer to optimize the loss
          # Notice that we're using sample_batches=False, so we're performing Gradient Descent
          best_model, loss_history, _, _ = trainer.train(X, y, X, y,
                                                       model, two_layer_net,
                                                       reg=0.001,
                                                       learning_rate=1e-1, momentum=0.9, learning_rate=1e-1,
                                                       update='rmsprop', sample_batches=False,
                                                       num_epochs=100,
                                                       verbose=False)

          correct_loss = 0.439368
          print('Final loss with RMSProp: %f. We get: %f' % (loss_history[-1], correct_loss))

starting iteration 0
starting iteration 10
starting iteration 20
starting iteration 30
starting iteration 40
starting iteration 50
starting iteration 60
starting iteration 70
starting iteration 80
starting iteration 90
Final loss with RMSProp: 0.439368. We get: 0.439368

```

6 Load the data

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

```
In [58]: from cs231n.data_utils import load_CIFAR10

def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000):
    """
    Load the CIFAR-10 dataset from disk and perform preprocessing to prepare
    it for the two-layer neural net classifier.
    """
    # Load the raw CIFAR-10 data
    cifar10_dir = 'cs231n/datasets/cifar-10-batches-py'
    X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

    # Subsample the data
    mask = range(num_training, num_training + num_validation)
    X_val = X_train[mask]
    y_val = y_train[mask]
    mask = range(num_training)
    X_train = X_train[mask]
    y_train = y_train[mask]
    mask = range(num_test)
    X_test = X_test[mask]
    y_test = y_test[mask]

    # Normalize the data: subtract the mean image
    mean_image = np.mean(X_train, axis=0)
    X_train -= mean_image
    X_val -= mean_image
    X_test -= mean_image

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

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

# Invoke the above function to get our data.
X_train, y_train, X_val, y_val, X_test, y_test = get_CIFAR10_data()
print('Train data shape: ', X_train.shape)
print('Train labels shape: ', y_train.shape)
print('Validation data shape: ', X_val.shape)
print('Validation labels shape: ', y_val.shape)
print('Test data shape: ', X_test.shape)
```

```

print('Test labels shape: ', y_test.shape)

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

```

7 Train a network

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

In [59]: `from cs231n.classifiers.neural_net import init_two_layer_model`

```

model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
best_model, loss_history, train_acc, val_acc = trainer.train(X_train, y_train, X_val,
                                                             model, two_layer_net,
                                                             num_epochs=5, reg=1.0,
                                                             momentum=0.9, learning_rate_decay = 0.95
                                                             learning_rate=1e-5, verbose=True)

starting iteration 0
Finished epoch 0 / 5: cost 2.302593, train: 0.086000, val 0.088000, lr 1.000000e-05
starting iteration 10
starting iteration 20
starting iteration 30
starting iteration 40
starting iteration 50
starting iteration 60
starting iteration 70
starting iteration 80
starting iteration 90
starting iteration 100
starting iteration 110
starting iteration 120
starting iteration 130
starting iteration 140
starting iteration 150
starting iteration 160
starting iteration 170
starting iteration 180
starting iteration 190
starting iteration 200

```

```
starting iteration 210
starting iteration 220
starting iteration 230
starting iteration 240
starting iteration 250
starting iteration 260
starting iteration 270
starting iteration 280
starting iteration 290
starting iteration 300
starting iteration 310
starting iteration 320
starting iteration 330
starting iteration 340
starting iteration 350
starting iteration 360
starting iteration 370
starting iteration 380
starting iteration 390
starting iteration 400
starting iteration 410
starting iteration 420
starting iteration 430
starting iteration 440
starting iteration 450
starting iteration 460
starting iteration 470
starting iteration 480
Finished epoch 1 / 5: cost 2.279186, train: 0.172000, val 0.185000, lr 9.500000e-06
starting iteration 490
starting iteration 500
starting iteration 510
starting iteration 520
starting iteration 530
starting iteration 540
starting iteration 550
starting iteration 560
starting iteration 570
starting iteration 580
starting iteration 590
starting iteration 600
starting iteration 610
starting iteration 620
starting iteration 630
starting iteration 640
starting iteration 650
starting iteration 660
starting iteration 670
```

```
starting iteration 680
starting iteration 690
starting iteration 700
starting iteration 710
starting iteration 720
starting iteration 730
starting iteration 740
starting iteration 750
starting iteration 760
starting iteration 770
starting iteration 780
starting iteration 790
starting iteration 800
starting iteration 810
starting iteration 820
starting iteration 830
starting iteration 840
starting iteration 850
starting iteration 860
starting iteration 870
starting iteration 880
starting iteration 890
starting iteration 900
starting iteration 910
starting iteration 920
starting iteration 930
starting iteration 940
starting iteration 950
starting iteration 960
starting iteration 970
Finished epoch 2 / 5: cost 2.122511, train: 0.274000, val 0.237000, lr 9.025000e-06
starting iteration 980
starting iteration 990
starting iteration 1000
starting iteration 1010
starting iteration 1020
starting iteration 1030
starting iteration 1040
starting iteration 1050
starting iteration 1060
starting iteration 1070
starting iteration 1080
starting iteration 1090
starting iteration 1100
starting iteration 1110
starting iteration 1120
starting iteration 1130
starting iteration 1140
```

```
starting iteration 1150
starting iteration 1160
starting iteration 1170
starting iteration 1180
starting iteration 1190
starting iteration 1200
starting iteration 1210
starting iteration 1220
starting iteration 1230
starting iteration 1240
starting iteration 1250
starting iteration 1260
starting iteration 1270
starting iteration 1280
starting iteration 1290
starting iteration 1300
starting iteration 1310
starting iteration 1320
starting iteration 1330
starting iteration 1340
starting iteration 1350
starting iteration 1360
starting iteration 1370
starting iteration 1380
starting iteration 1390
starting iteration 1400
starting iteration 1410
starting iteration 1420
starting iteration 1430
starting iteration 1440
starting iteration 1450
starting iteration 1460
Finished epoch 3 / 5: cost 1.995343, train: 0.282000, val 0.292000, lr 8.573750e-06
starting iteration 1470
starting iteration 1480
starting iteration 1490
starting iteration 1500
starting iteration 1510
starting iteration 1520
starting iteration 1530
starting iteration 1540
starting iteration 1550
starting iteration 1560
starting iteration 1570
starting iteration 1580
starting iteration 1590
starting iteration 1600
starting iteration 1610
```

```
starting iteration 1620
starting iteration 1630
starting iteration 1640
starting iteration 1650
starting iteration 1660
starting iteration 1670
starting iteration 1680
starting iteration 1690
starting iteration 1700
starting iteration 1710
starting iteration 1720
starting iteration 1730
starting iteration 1740
starting iteration 1750
starting iteration 1760
starting iteration 1770
starting iteration 1780
starting iteration 1790
starting iteration 1800
starting iteration 1810
starting iteration 1820
starting iteration 1830
starting iteration 1840
starting iteration 1850
starting iteration 1860
starting iteration 1870
starting iteration 1880
starting iteration 1890
starting iteration 1900
starting iteration 1910
starting iteration 1920
starting iteration 1930
starting iteration 1940
starting iteration 1950
Finished epoch 4 / 5: cost 1.783637, train: 0.339000, val 0.333000, lr 8.145063e-06
starting iteration 1960
starting iteration 1970
starting iteration 1980
starting iteration 1990
starting iteration 2000
starting iteration 2010
starting iteration 2020
starting iteration 2030
starting iteration 2040
starting iteration 2050
starting iteration 2060
starting iteration 2070
starting iteration 2080
```

```
starting iteration 2090
starting iteration 2100
starting iteration 2110
starting iteration 2120
starting iteration 2130
starting iteration 2140
starting iteration 2150
starting iteration 2160
starting iteration 2170
starting iteration 2180
starting iteration 2190
starting iteration 2200
starting iteration 2210
starting iteration 2220
starting iteration 2230
starting iteration 2240
starting iteration 2250
starting iteration 2260
starting iteration 2270
starting iteration 2280
starting iteration 2290
starting iteration 2300
starting iteration 2310
starting iteration 2320
starting iteration 2330
starting iteration 2340
starting iteration 2350
starting iteration 2360
starting iteration 2370
starting iteration 2380
starting iteration 2390
starting iteration 2400
starting iteration 2410
starting iteration 2420
starting iteration 2430
starting iteration 2440
Finished epoch 5 / 5: cost 1.874293, train: 0.360000, val 0.363000, lr 7.737809e-06
finished optimization. best validation accuracy: 0.363000
```

8 Debug the training

With the default parameters we provided above, you should get a validation accuracy of about 0.37 on the validation set. This isn't very good.

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

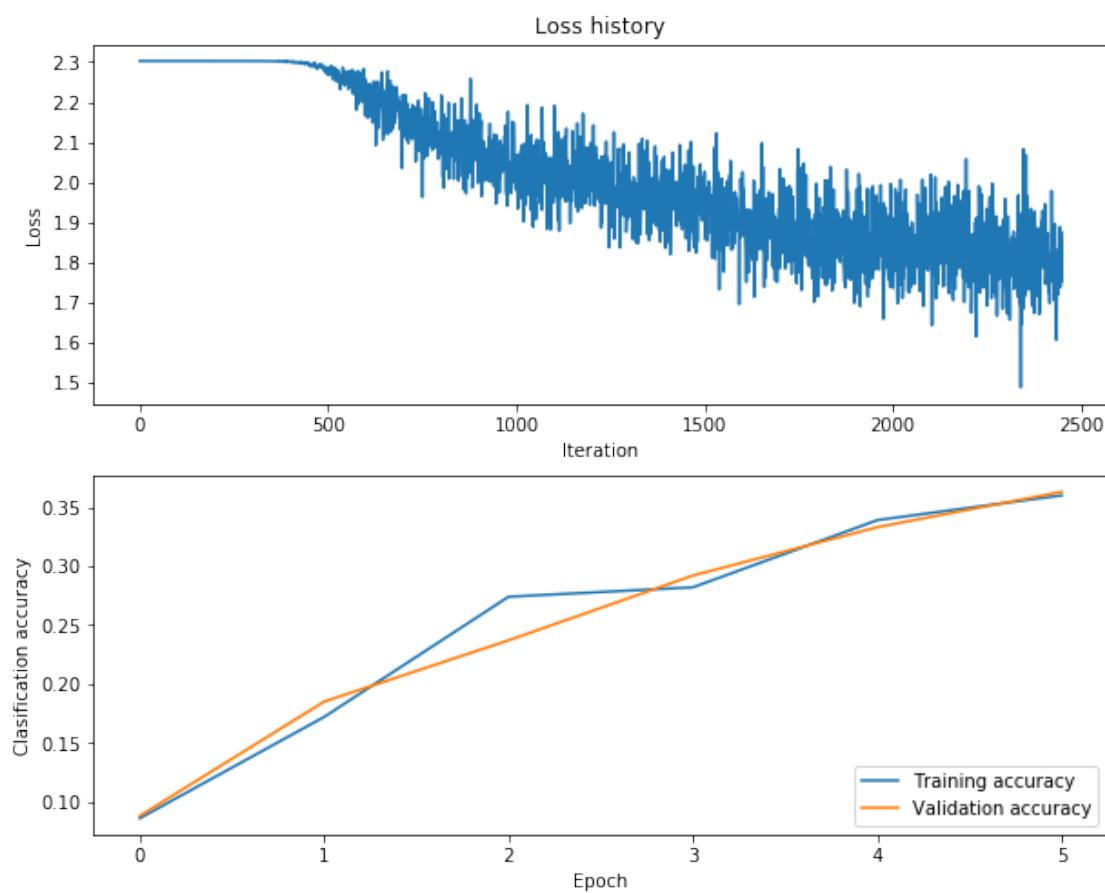
Another strategy is to visualize the weights that were learned in the first layer of the network.

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

```
In [60]: # Plot the loss function and train / validation accuracies
    plt.subplot(2, 1, 1)
    plt.plot(loss_history)
    plt.title('Loss history')
    plt.xlabel('Iteration')
    plt.ylabel('Loss')

    plt.subplot(2, 1, 2)
    plt.plot(train_acc)
    plt.plot(val_acc)
    plt.legend(['Training accuracy', 'Validation accuracy'], loc='lower right')
    plt.xlabel('Epoch')
    plt.ylabel('Classification accuracy')

Out[60]: Text(0,0.5,'Classification accuracy')
```



```
In [61]: from cs231n.vis_utils import visualize_grid
```

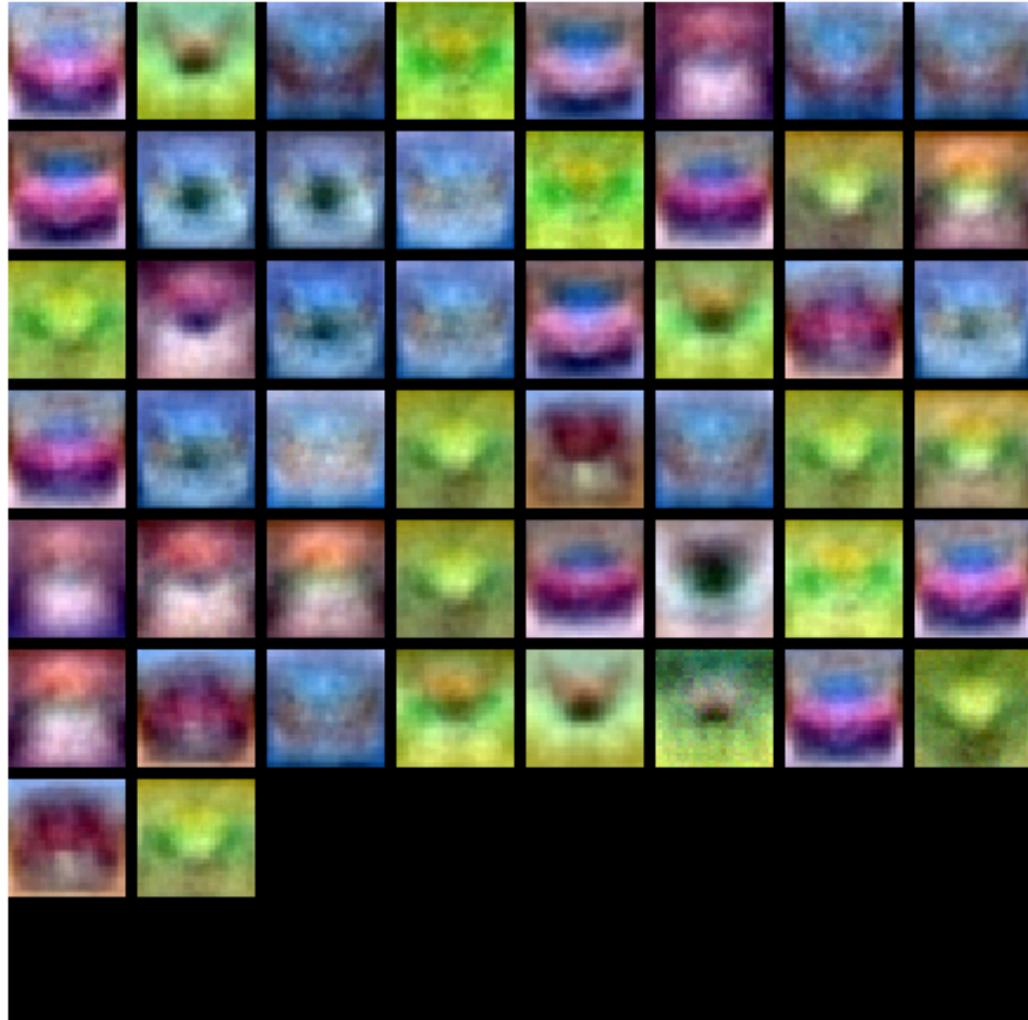
```

# Visualize the weights of the network

def show_net_weights(model):
    plt.imshow(visualize_grid(model['W1']).T.reshape(-1, 32, 32, 3), padding=3).astype('uint8')
    plt.gca().axis('off')
    plt.show()

show_net_weights(model)

```



9 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, numer of training epochs, and regularization strength. You might also consider tuning the momentum and learning rate decay parameters, but you should be able to get good performance using the default values.

Approximate results. You should be aim to achieve a classification accuracy of greater than 50% on the validation set. Our best network gets over 56% on the validation set.

Experiment: Your goal in this exercise is to get as good of a result on CIFAR-10 as you can, with a fully-connected Neural Network. For every 1% above 56% on the Test set we will award you with one extra bonus point. Feel free implement your own techniques (e.g. PCA to reduce dimensionality, or adding dropout, or adding features to the solver, etc.).

In [1]: `best_model = None # store the best model into this`

```
#####
# 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 the previous assignment.                      #
#####
# input size, hidden size, number of classes
model = init_two_layer_model(32*32*3, 1000, 10)
trainer = ClassifierTrainer()
best_model, loss_history, train_acc, val_acc = trainer.train(X_train, y_train,
                                                             X_val, y_val,
                                                             model, two_layer_net,
                                                             num_epochs=10, reg=1e-3,
                                                             momentum=0.99,
                                                             learning_rate_decay=0.95,
                                                             learning_rate=3e-6, verbose=False)
#####
#                                              END OF YOUR CODE                  #
#####
```

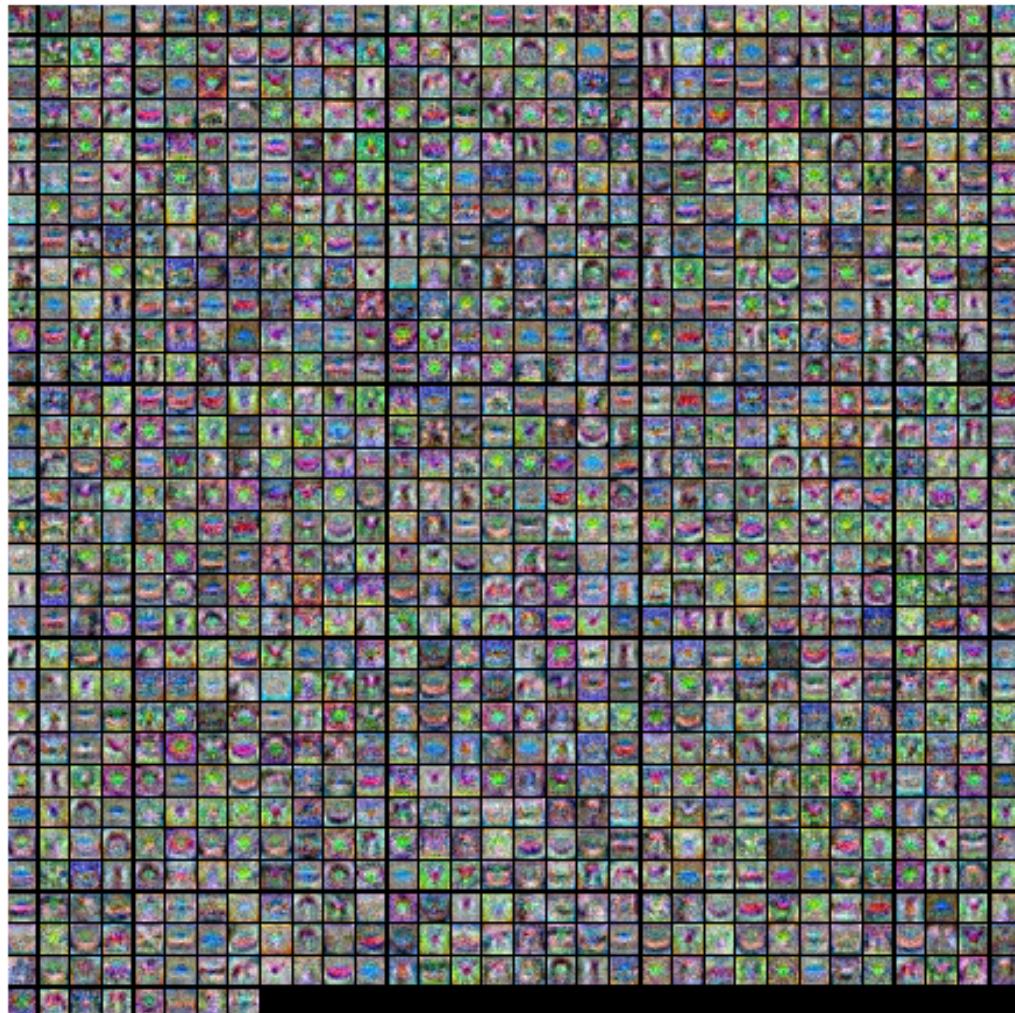
```
NameError
```

```
Traceback (most recent call last)
```

```
<ipython-input-1-97978f118412> in <module>()
 14 #####
 15 # input size, hidden size, number of classes
---> 16 model = init_two_layer_model(32*32*3, 1000, 10)
 17 trainer = ClassifierTrainer()
 18 best_model, loss_history, train_acc, val_acc = trainer.train(X_train, y_train,
```

```
NameError: name 'init_two_layer_model' is not defined
```

```
In [76]: # visualize the weights
show_net_weights(best_model)
```



10 Run on the test set

When you are done experimenting, you should evaluate your final trained network on the test set.

```
In [77]: scores_test = two_layer_net(X_test, best_model)
print('Test accuracy: ', np.mean(np.argmax(scores_test, axis=1) == y_test))
```

```
Test accuracy:  0.516
```

```
In [ ]:
```

layers

September 26, 2019

1 Modular neural nets

In the previous exercise, we computed the loss and gradient for a two-layer neural network in a single monolithic function. This isn't very difficult for a small two-layer network, but would be tedious and error-prone for larger networks. Ideally we want to build networks using a more modular design so that we can snap together different types of layers and loss functions in order to quickly experiment with different architectures.

In this exercise we will implement this approach, and develop a number of different layer types in isolation that can then be easily plugged together. For each layer we will implement `forward` and `backward` functions. The `forward` function will receive data, weights, and other parameters, and will return both an output and a `cache` object that stores data needed for the `backward` pass. The `backward` function will receive upstream derivatives and the `cache` object, and will return gradients with respect to the data and all of the weights. This will allow us to write code that looks like this:

```
def two_layer_net(X, W1, b1, W2, b2, reg):
    # Forward pass; compute scores
    s1, fc1_cache = affine_forward(X, W1, b1)
    a1, relu_cache = relu_forward(s1)
    scores, fc2_cache = affine_forward(a1, W2, b2)

    # Loss functions return data loss and gradients on scores
    data_loss, dscores = svm_loss(scores, y)

    # Compute backward pass
    da1, dW2, db2 = affine_backward(dscores, fc2_cache)
    ds1 = relu_backward(da1, relu_cache)
    dX, dW1, db1 = affine_backward(ds1, fc1_cache)

    # A real network would add regularization here

    # Return loss and gradients
    return loss, dW1, db1, dW2, db2
```

In [1]: # As usual, a bit of setup

```
import numpy as np
```

```

import matplotlib.pyplot as plt
from cs231n.gradient_check import eval_numerical_gradient_array, eval_numerical_gradient
from cs231n.layers import *

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

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

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

```

2 Affine layer: forward

Open the file `cs231n/layers.py` and implement the `affine_forward` function.

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

In [2]: `# 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 1e-9.
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. You can test your implementation using numeric gradient checking.

In [7]: # Test the `affine_backward` function

```
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 less than 1e-10
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:  3.3648272856060417e-10
dw error:  2.3867686814528706e-10
db error:  1.4081818187780606e-11
```

4 ReLU layer: forward

Implement the `relu_forward` function and test your implementation by running the following:

In [8]: # 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 around 1e-8
print('Testing relu_forward function:')
print('difference: ', rel_error(out, correct_out))
```

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

5 ReLU layer: backward

Implement the `relu_backward` function and test your implementation using numeric gradient checking:

```
In [9]: 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 around 1e-12  
print('Testing relu_backward function:')  
print('dx error: ', rel_error(dx_num, dx))  
  
Testing relu_backward function:  
dx error: 3.2755825341276888e-12
```

6 Loss layers: Softmax and SVM

You implemented these loss functions in the last assignment, so we'll give them to you for free here. It's still a good idea to test them to make sure they work correctly.

```
In [10]: 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: svm_loss(x, y)[0], x, verbose=False)  
loss, dx = svm_loss(x, y)  
  
# Test sum_loss function. Loss should be around 9 and dx error should be 1e-9  
print('Testing svm_loss:')  
print('loss: ', loss)  
print('dx error: ', rel_error(dx_num, dx))  
  
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 2.3 and dx error should be 1e-8  
print('\nTesting softmax_loss:')
```

```

print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))

Testing svm_loss:
loss: 9.000413245081784
dx error: 8.182894472887002e-10

Testing softmax_loss:
loss: 2.302626876996529
dx error: 9.558381239354783e-09

```

7 Convolution layer: forward naive

We are now ready to implement the forward pass for a convolutional layer. Implement the function `conv_forward_naive` in the file `cs231n/layers.py`.

You don't have to worry too much about efficiency at this point; just write the code in whatever way you find most clear.

You can test your implementation by running the following:

```

In [24]: x_shape = (2, 3, 4, 4)
w_shape = (3, 3, 4, 4)
x = np.linspace(-0.1, 0.5, num=np.prod(x_shape)).reshape(x_shape)
w = np.linspace(-0.2, 0.3, num=np.prod(w_shape)).reshape(w_shape)
b = np.linspace(-0.1, 0.2, num=3)

conv_param = {'stride': 2, 'pad': 1}
out, _ = conv_forward_naive(x, w, b, conv_param)
correct_out = np.array([[[[-0.08759809, -0.10987781],
                         [-0.18387192, -0.2109216 ]],,
                        [[ 0.21027089,  0.21661097],
                         [ 0.22847626,  0.23004637]],,
                        [[ 0.50813986,  0.54309974],
                         [ 0.64082444,  0.67101435]]],,
                       [[[ -0.98053589, -1.03143541],
                         [-1.19128892, -1.24695841]],,
                        [[ 0.69108355,  0.66880383],
                         [ 0.59480972,  0.56776003]],,
                        [[ 2.36270298,  2.36904306],
                         [ 2.38090835,  2.38247847]]]]))

# Compare your output to ours; difference should be around 1e-8
print('Testing conv_forward_naive')
print('difference: ', rel_error(out, correct_out))

Testing conv_forward_naive
difference: 2.2121476417505994e-08

```

8 Aside: Image processing via convolutions

As fun way to both check your implementation and gain a better understanding of the type of operation that convolutional layers can perform, we will set up an input containing two images and manually set up filters that perform common image processing operations (grayscale conversion and edge detection). The convolution forward pass will apply these operations to each of the input images. We can then visualize the results as a sanity check.

```
In [25]: from scipy.misc import imread, imresize

kitten, puppy = imread('kitten.jpg'), imread('puppy.jpg')
# kitten is wide, and puppy is already square
d = kitten.shape[1] - kitten.shape[0]
kitten_cropped = kitten[:, d//2:-d//2, :]

img_size = 200    # Make this smaller if it runs too slow
x = np.zeros((2, 3, img_size, img_size))
x[0, :, :, :] = imresize(puppy, (img_size, img_size)).transpose((2, 0, 1))
x[1, :, :, :] = imresize(kitten_cropped, (img_size, img_size)).transpose((2, 0, 1))

# Set up a convolutional weights holding 2 filters, each 3x3
w = np.zeros((2, 3, 3, 3))

# The first filter converts the image to grayscale.
# Set up the red, green, and blue channels of the filter.
w[0, 0, :, :] = [[0, 0, 0], [0, 0.3, 0], [0, 0, 0]]
w[0, 1, :, :] = [[0, 0, 0], [0, 0.6, 0], [0, 0, 0]]
w[0, 2, :, :] = [[0, 0, 0], [0, 0.1, 0], [0, 0, 0]]

# Second filter detects horizontal edges in the blue channel.
w[1, 2, :, :] = [[1, 2, 1], [0, 0, 0], [-1, -2, -1]]

# Vector of biases. We don't need any bias for the grayscale
# filter, but for the edge detection filter we want to add 128
# to each output so that nothing is negative.
b = np.array([0, 128])

# Compute the result of convolving each input in x with each filter in w,
# offsetting by b, and storing the results in out.
out, _ = conv_forward_naive(x, w, b, {'stride': 1, 'pad': 1})

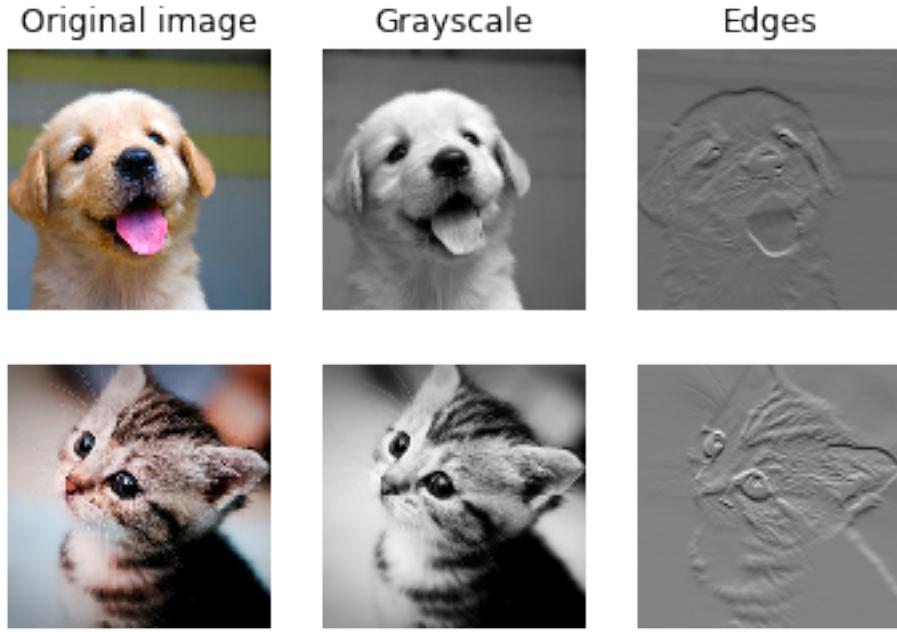
def imshow_noax(img, normalize=True):
    """ Tiny helper to show images as uint8 and remove axis labels """
    if normalize:
        img_max, img_min = np.max(img), np.min(img)
        img = 255.0 * (img - img_min) / (img_max - img_min)
    plt.imshow(img.astype('uint8'))
    plt.gca().axis('off')
```

```

# Show the original images and the results of the conv operation
plt.subplot(2, 3, 1)
imshow_noax(puppy, normalize=False)
plt.title('Original image')
plt.subplot(2, 3, 2)
imshow_noax(out[0, 0])
plt.title('Grayscale')
plt.subplot(2, 3, 3)
imshow_noax(out[0, 1])
plt.title('Edges')
plt.subplot(2, 3, 4)
imshow_noax(kitten_cropped, normalize=False)
plt.subplot(2, 3, 5)
imshow_noax(out[1, 0])
plt.subplot(2, 3, 6)
imshow_noax(out[1, 1])
plt.show()

c:\users\kingsman142\appdata\local\programs\python\python36\lib\site-packages\ipykernel_launcher
`imread` is deprecated in SciPy 1.0.0, and will be removed in 1.2.0.
Use ``imageio.imread`` instead.
This is separate from the ipykernel package so we can avoid doing imports until
c:\users\kingsman142\appdata\local\programs\python\python36\lib\site-packages\ipykernel_launcher
`imresize` is deprecated in SciPy 1.0.0, and will be removed in 1.2.0.
Use ``skimage.transform.resize`` instead.
# Remove the CWD from sys.path while we load stuff.
c:\users\kingsman142\appdata\local\programs\python\python36\lib\site-packages\ipykernel_launcher
`imresize` is deprecated in SciPy 1.0.0, and will be removed in 1.2.0.
Use ``skimage.transform.resize`` instead.
# This is added back by InteractiveShellApp.init_path()

```



9 Convolution layer: backward naive

Next you need to implement the function `conv_backward_naive` in the file `cs231n/layers.py`. As usual, we will check your implementation with numeric gradient checking.

```
In [75]: x = np.random.randn(4, 3, 5, 5)
w = np.random.randn(2, 3, 3, 3)
b = np.random.randn(2,)
dout = np.random.randn(4, 2, 5, 5)
conv_param = {'stride': 1, 'pad': 1}

dx_num = eval_numerical_gradient_array(lambda x: conv_forward_naive(x, w, b, conv_param),
dw_num = eval_numerical_gradient_array(lambda w: conv_forward_naive(x, w, b, conv_param),
db_num = eval_numerical_gradient_array(lambda b: conv_forward_naive(x, w, b, conv_param))

out, cache = conv_forward_naive(x, w, b, conv_param)
dx, dw, db = conv_backward_naive(dout, cache)

# Your errors should be around 1e-9'
print('Testing conv_backward_naive function')
print('dx error: ', rel_error(dx, dx_num))
print('dw error: ', rel_error(dw, dw_num))
print('db error: ', rel_error(db, db_num))
```

```
Testing conv_backward_naive function
dx error:  8.154152208997038e-10
dw error:  8.668200116734592e-10
db error:  2.412540530821696e-10
```

10 Max pooling layer: forward naive

The last layer we need for a basic convolutional neural network is the max pooling layer. First implement the forward pass in the function `max_pool_forward_naive` in the file `cs231n/layers.py`.

```
In [76]: x_shape = (2, 3, 4, 4)
x = np.linspace(-0.3, 0.4, num=np.prod(x_shape)).reshape(x_shape)
pool_param = {'pool_width': 2, 'pool_height': 2, 'stride': 2}

out, _ = max_pool_forward_naive(x, pool_param)

correct_out = np.array([[[[-0.26315789, -0.24842105],
                         [-0.20421053, -0.18947368]],
                        [[-0.14526316, -0.13052632],
                         [-0.08631579, -0.07157895]],
                        [[-0.02736842, -0.01263158],
                         [ 0.03157895,  0.04631579]]],
                       [[[ 0.09052632,  0.10526316],
                         [ 0.14947368,  0.16421053]],
                        [[ 0.20842105,  0.22315789],
                         [ 0.26736842,  0.28210526]],
                        [[ 0.32631579,  0.34105263],
                         [ 0.38526316,  0.4        ]]]])

# Compare your output with ours. Difference should be around 1e-8.
print('Testing max_pool_forward_naive function:')
print('difference: ', rel_error(out, correct_out))

Testing max_pool_forward_naive function:
difference:  4.1666665157267834e-08
```

11 Max pooling layer: backward naive

Implement the backward pass for a max pooling layer in the function `max_pool_backward_naive` in the file `cs231n/layers.py`. As always we check the correctness of the backward pass using numerical gradient checking.

```
In [77]: x = np.random.randn(3, 2, 8, 8)
dout = np.random.randn(3, 2, 4, 4)
pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}
```

```

dx_num = eval_numerical_gradient_array(lambda x: max_pool_forward_naive(x, pool_param),
                                         x, pool_param)

out, cache = max_pool_forward_naive(x, pool_param)
dx = max_pool_backward_naive(dout, cache)

# Your error should be around 1e-12
print('Testing max_pool_backward_naive function:')
print('dx error: ', rel_error(dx, dx_num))

Testing max_pool_backward_naive function:
dx error:  3.2756125356666546e-12

```

12 Fast layers

Making convolution and pooling layers fast can be challenging. To spare you the pain, we've provided fast implementations of the forward and backward passes for convolution and pooling layers in the file `cs231n/fast_layers.py`.

The fast convolution implementation depends on a Cython extension; to compile it you need to run the following from the `cs231n` directory:

```
python setup.py build_ext --inplace
```

The API for the fast versions of the convolution and pooling layers is exactly the same as the naive versions that you implemented above: the forward pass receives data, weights, and parameters and produces outputs and a cache object; the backward pass receives upstream derivatives and the cache object and produces gradients with respect to the data and weights.

NOTE: The fast implementation for pooling will only perform optimally if the pooling regions are non-overlapping and tile the input. If these conditions are not met then the fast pooling implementation will not be much faster than the naive implementation.

You can compare the performance of the naive and fast versions of these layers by running the following:

```
In [78]: from cs231n.fast_layers import conv_forward_fast, conv_backward_fast
         from time import time

         x = np.random.randn(100, 3, 31, 31)
         w = np.random.randn(25, 3, 3, 3)
         b = np.random.randn(25,)
         dout = np.random.randn(100, 25, 16, 16)
         conv_param = {'stride': 2, 'pad': 1}

         t0 = time()
         out_naive, cache_naive = conv_forward_naive(x, w, b, conv_param)
         t1 = time()
         out_fast, cache_fast = conv_forward_fast(x, w, b, conv_param)
         t2 = time()
```

```

print('Testing conv_forward_fast:')
print('Naive: %fs' % (t1 - t0))
print('Fast: %fs' % (t2 - t1))
print('Speedup: %fx' % ((t1 - t0) / (t2 - t1)))
print('Difference: ', rel_error(out_naive, out_fast))

t0 = time()
dx_naive, dw_naive, db_naive = conv_backward_naive(dout, cache_naive)
t1 = time()
dx_fast, dw_fast, db_fast = conv_backward_fast(dout, cache_fast)
t2 = time()

print('\nTesting conv_backward_fast:')
print('Naive: %fs' % (t1 - t0))
print('Fast: %fs' % (t2 - t1))
print('Speedup: %fx' % ((t1 - t0) / (t2 - t1)))
print('dx difference: ', rel_error(dx_naive, dx_fast))
print('dw difference: ', rel_error(dw_naive, dw_fast))
print('db difference: ', rel_error(db_naive, db_fast))

Testing conv_forward_fast:
Naive: 3.896232s
Fast: 0.020979s
Speedup: 185.716992x
Difference: 2.7676579517961165e-09

Testing conv_backward_fast:
Naive: 6.921198s
Fast: 0.017982s
Speedup: 384.895747x
dx difference: 1.4424397962353693e-11
dw difference: 2.552069923833236e-12
db difference: 0.0

```

```

In [79]: from cs231n.fast_layers import max_pool_forward_fast, max_pool_backward_fast

x = np.random.randn(100, 3, 32, 32)
dout = np.random.randn(100, 3, 16, 16)
pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}

t0 = time()
out_naive, cache_naive = max_pool_forward_naive(x, pool_param)
t1 = time()
out_fast, cache_fast = max_pool_forward_fast(x, pool_param)
t2 = time()

```

```

print('Testing pool_forward_fast:')
print('Naive: %fs' % (t1 - t0))
print('fast: %fs' % (t2 - t1))
print('speedup: %fx' % ((t1 - t0) / (t2 - t1)))
print('difference: ', rel_error(out_naive, out_fast))

t0 = time()
dx_naive = max_pool_backward_naive(dout, cache_naive)
t1 = time()
dx_fast = max_pool_backward_fast(dout, cache_fast)
t2 = time()

print('\nTesting pool_backward_fast:')
print('Naive: %fs' % (t1 - t0))
print('speedup: %fx' % ((t1 - t0) / (t2 - t1)))
print('dx difference: ', rel_error(dx_naive, dx_fast))

Testing pool_forward_fast:
Naive: 0.299684s
fast: 0.003988s
speedup: 75.154918x
difference:  0.0

Testing pool_backward_fast:
Naive: 0.355653s
speedup: 25.430301x
dx difference:  0.0

```

13 Sandwich layers

There are a couple common layer “sandwiches” that frequently appear in ConvNets. For example convolutional layers are frequently followed by ReLU and pooling, and affine layers are frequently followed by ReLU. To make it more convenient to use these common patterns, we have defined several convenience layers in the file cs231n/layer_utils.py. Lets grad-check them to make sure that they work correctly:

```
In [80]: from cs231n.layer_utils import conv_relu_pool_forward, conv_relu_pool_backward

x = np.random.randn(2, 3, 16, 16)
w = np.random.randn(3, 3, 3, 3)
b = np.random.randn(3,)
dout = np.random.randn(2, 3, 8, 8)
conv_param = {'stride': 1, 'pad': 1}
pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}

out, cache = conv_relu_pool_forward(x, w, b, conv_param, pool_param)
dx, dw, db = conv_relu_pool_backward(dout, cache)
```

```

dx_num = eval_numerical_gradient_array(lambda x: conv_relu_pool_forward(x, w, b, conv_param),
dw_num = eval_numerical_gradient_array(lambda w: conv_relu_pool_forward(x, w, b, conv_param),
db_num = eval_numerical_gradient_array(lambda b: conv_relu_pool_forward(x, w, b, conv_param))

print('Testing conv_relu_pool_forward:')
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 conv_relu_pool_forward:
dx error:  1.2684020766235086e-08
dw error:  4.3494413138337085e-10
db error:  9.731066997072821e-11

```

In [81]: `from cs231n.layer_utils import conv_relu_forward, conv_relu_backward`

```

x = np.random.randn(2, 3, 8, 8)
w = np.random.randn(3, 3, 3, 3)
b = np.random.randn(3,)
dout = np.random.randn(2, 3, 8, 8)
conv_param = {'stride': 1, 'pad': 1}

out, cache = conv_relu_forward(x, w, b, conv_param)
dx, dw, db = conv_relu_backward(dout, cache)

dx_num = eval_numerical_gradient_array(lambda x: conv_relu_forward(x, w, b, conv_param),
dw_num = eval_numerical_gradient_array(lambda w: conv_relu_forward(x, w, b, conv_param),
db_num = eval_numerical_gradient_array(lambda b: conv_relu_forward(x, w, b, conv_param))

print('Testing conv_relu_forward:')
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 conv_relu_forward:
dx error:  7.157175523005417e-10
dw error:  2.523250094484491e-10
db error:  4.754190955641717e-12

```

In [82]: `from cs231n.layer_utils import affine_relu_forward, affine_relu_backward`

```

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, dx)
dw_num = eval_numerical_gradient_array(lambda w: affine_relu_forward(x, w, b)[0], w, dw)
db_num = eval_numerical_gradient_array(lambda b: affine_relu_forward(x, w, b)[0], b, db)

print('Testing affine_relu_forward:')
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:
dx error:  1.0549467212422243e-09
dw error:  6.433770438106674e-10
db error:  2.5479984708360518e-11
```

In []:

convnet

September 26, 2019

1 Train a ConvNet!

We now have a generic solver and a bunch of modularized layers. It's time to put it all together, and train a ConvNet to recognize the classes in CIFAR-10. In this notebook we will walk you through training a simple two-layer ConvNet and then set you free to build the best net that you can to perform well on CIFAR-10.

Open up the file `cs231n/classifiers/convnet.py`; you will see that the `two_layer_convnet` function computes the loss and gradients for a two-layer ConvNet. Note that this function uses the "sandwich" layers defined in `cs231n/layer_utils.py`.

In [1]: # As usual, a bit of setup

```
import numpy as np
import matplotlib.pyplot as plt
from cs231n.classifier_trainer import ClassifierTrainer
from cs231n.gradient_check import eval_numerical_gradient
from cs231n.classifiers.convnet import *

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

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

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

In [2]: from cs231n.data_utils import load_CIFAR10

```
def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000):
    """
    Load the CIFAR-10 dataset from disk and perform preprocessing to prepare
    it for the two-layer neural net classifier. These are the same steps as
    
```

```

we used for the SVM, but condensed to a single function.
"""

# Load the raw CIFAR-10 data
cifar10_dir = 'cs231n/datasets/cifar-10-batches-py'
X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

# Subsample the data
mask = range(num_training, num_training + num_validation)
X_val = X_train[mask]
y_val = y_train[mask]
mask = range(num_training)
X_train = X_train[mask]
y_train = y_train[mask]
mask = range(num_test)
X_test = X_test[mask]
y_test = y_test[mask]

# Normalize the data: subtract the mean image
mean_image = np.mean(X_train, axis=0)
X_train -= mean_image
X_val -= mean_image
X_test -= mean_image

# Transpose so that channels come first
X_train = X_train.transpose(0, 3, 1, 2).copy()
X_val = X_val.transpose(0, 3, 1, 2).copy()
X_test = X_test.transpose(0, 3, 1, 2).copy()

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

# Invoke the above function to get our data.
X_train, y_train, X_val, y_val, X_test, y_test = get_CIFAR10_data()
print('Train data shape: ', X_train.shape)
print('Train labels shape: ', y_train.shape)
print('Validation data shape: ', X_val.shape)
print('Validation labels shape: ', y_val.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)

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

```

2 Sanity check loss

After you build a new network, one of the first things you should do is sanity check the loss. When we use the softmax loss, we expect the loss for random weights (and no regularization) to be about $\log(C)$ for C classes. When we add regularization this should go up.

```
In [3]: model = init_two_layer_convnet()

X = np.random.randn(100, 3, 32, 32)
y = np.random.randint(10, size=100)

loss, _ = two_layer_convnet(X, model, y, reg=0)

# Sanity check: Loss should be about log(10) = 2.3026
print('Sanity check loss (no regularization): ', loss)

# Sanity check: Loss should go up when you add regularization
loss, _ = two_layer_convnet(X, model, y, reg=1)
print('Sanity check loss (with regularization): ', loss)

Sanity check loss (no regularization):  2.3026068267795177
Sanity check loss (with regularization):  2.3447459205935943
```

3 Gradient check

After the loss looks reasonable, you should always use numeric gradient checking to make sure that your backward pass is correct. When you use numeric gradient checking you should use a small amount of artificial data and a small number of neurons at each layer.

```
In [4]: num_inputs = 2
        input_shape = (3, 16, 16)
        reg = 0.0
        num_classes = 10
        X = np.random.randn(num_inputs, *input_shape)
        y = np.random.randint(num_classes, size=num_inputs)

        model = init_two_layer_convnet(num_filters=3, filter_size=3, input_shape=input_shape)
        loss, grads = two_layer_convnet(X, model, y)
        for param_name in sorted(grads):
            f = lambda _: two_layer_convnet(X, model, y)[0]
            param_grad_num = eval_numerical_gradient(f, model[param_name], verbose=False, h=1e-05)
            e = rel_error(param_grad_num, grads[param_name])
            print('%s max relative error: %e' % (param_name, rel_error(param_grad_num, grads[p
```

W1 max relative error: 3.090184e-06
W2 max relative error: 3.283356e-05
b1 max relative error: 6.233016e-08

```
b2 max relative error: 7.612983e-10
```

4 Overfit small data

A nice trick is to train your model with just a few training samples. You should be able to overfit small datasets, which will result in very high training accuracy and comparatively low validation accuracy.

```
In [5]: # Use a two-layer ConvNet to overfit 50 training examples.
```

```
model = init_two_layer_convnet()
trainer = ClassifierTrainer()
best_model, loss_history, train_acc_history, val_acc_history = trainer.train(
    X_train[:50], y_train[:50], X_val, y_val, model, two_layer_convnet,
    reg=0.001, momentum=0.9, learning_rate=0.0001, batch_size=10, num_epochs=10,
    verbose=True)

starting iteration  0
Finished epoch 0 / 10: cost 2.321363, train: 0.100000, val 0.083000, lr 1.000000e-04
Finished epoch 1 / 10: cost 2.295679, train: 0.240000, val 0.089000, lr 9.500000e-05
Finished epoch 2 / 10: cost 1.801747, train: 0.240000, val 0.106000, lr 9.025000e-05
starting iteration  10
Finished epoch 3 / 10: cost 1.231421, train: 0.340000, val 0.136000, lr 8.573750e-05
Finished epoch 4 / 10: cost 1.080222, train: 0.500000, val 0.159000, lr 8.145062e-05
starting iteration  20
Finished epoch 5 / 10: cost 1.322249, train: 0.640000, val 0.188000, lr 7.737809e-05
Finished epoch 6 / 10: cost 0.554413, train: 0.620000, val 0.145000, lr 7.350919e-05
starting iteration  30
Finished epoch 7 / 10: cost 0.648387, train: 0.760000, val 0.122000, lr 6.983373e-05
Finished epoch 8 / 10: cost 0.370819, train: 0.800000, val 0.123000, lr 6.634204e-05
starting iteration  40
Finished epoch 9 / 10: cost 0.213383, train: 0.880000, val 0.162000, lr 6.302494e-05
Finished epoch 10 / 10: cost 0.347561, train: 0.940000, val 0.180000, lr 5.987369e-05
finished optimization. best validation accuracy: 0.188000
```

Plotting the loss, training accuracy, and validation accuracy should show clear overfitting:

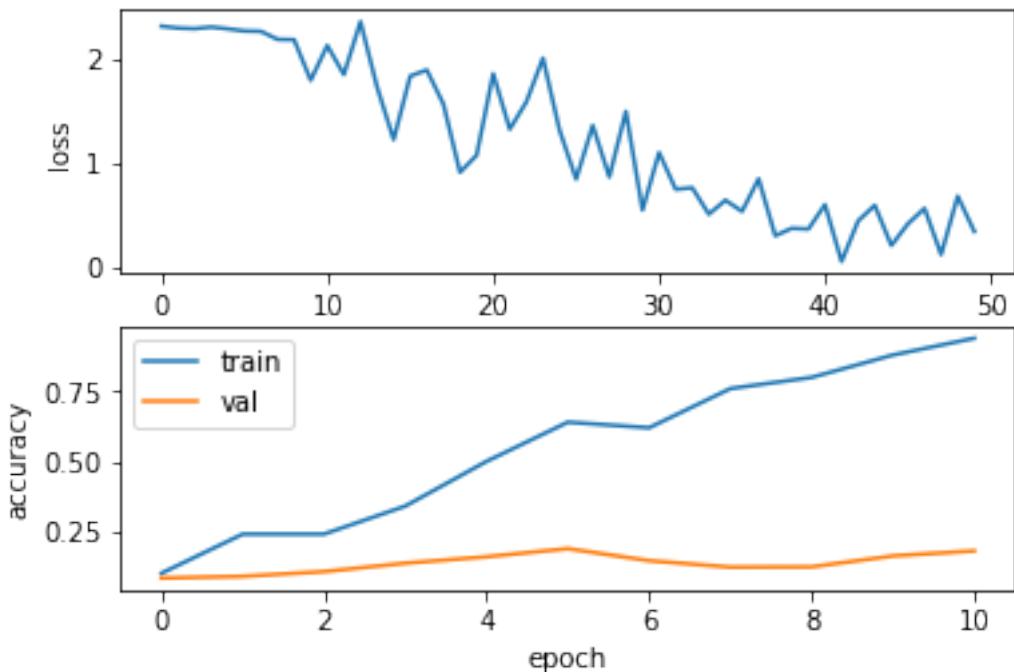
```
In [6]: plt.subplot(2, 1, 1)
plt.plot(loss_history)
plt.xlabel('iteration')
plt.ylabel('loss')

plt.subplot(2, 1, 2)
plt.plot(train_acc_history)
plt.plot(val_acc_history)
plt.legend(['train', 'val'], loc='upper left')
```

```

plt.xlabel('epoch')
plt.ylabel('accuracy')
plt.show()

```



5 Train the net

Once the above works, training the net is the next thing to try. You can set the `acc_frequency` parameter to change the frequency at which the training and validation set accuracies are tested. If your parameters are set properly, you should see the training and validation accuracy start to improve within a hundred iterations, and you should be able to train a reasonable model with just one epoch.

Using the parameters below you should be able to get around 50% accuracy on the validation set.

```

In [7]: model = init_two_layer_convnet(filter_size=7)
trainer = ClassifierTrainer()
best_model, loss_history, train_acc_history, val_acc_history = trainer.train(
    X_train, y_train, X_val, y_val, model, two_layer_convnet,
    reg=0.001, momentum=0.9, learning_rate=0.0001, batch_size=50, num_epochs=1,
    acc_frequency=50, verbose=True)

starting iteration  0
Finished epoch 0 / 1: cost 2.299287, train: 0.100000, val 0.094000, lr 1.000000e-04
starting iteration  10

```

```
starting iteration 20
starting iteration 30
starting iteration 40
starting iteration 50
Finished epoch 0 / 1: cost 2.113618, train: 0.336000, val 0.333000, lr 1.000000e-04
starting iteration 60
starting iteration 70
starting iteration 80
starting iteration 90
starting iteration 100
Finished epoch 0 / 1: cost 1.491748, train: 0.395000, val 0.397000, lr 1.000000e-04
starting iteration 110
starting iteration 120
starting iteration 130
starting iteration 140
starting iteration 150
Finished epoch 0 / 1: cost 2.198996, train: 0.398000, val 0.422000, lr 1.000000e-04
starting iteration 160
starting iteration 170
starting iteration 180
starting iteration 190
starting iteration 200
Finished epoch 0 / 1: cost 1.582417, train: 0.384000, val 0.415000, lr 1.000000e-04
starting iteration 210
starting iteration 220
starting iteration 230
starting iteration 240
starting iteration 250
Finished epoch 0 / 1: cost 1.117805, train: 0.370000, val 0.389000, lr 1.000000e-04
starting iteration 260
starting iteration 270
starting iteration 280
starting iteration 290
starting iteration 300
Finished epoch 0 / 1: cost 1.755130, train: 0.437000, val 0.450000, lr 1.000000e-04
starting iteration 310
starting iteration 320
starting iteration 330
starting iteration 340
starting iteration 350
Finished epoch 0 / 1: cost 1.773852, train: 0.389000, val 0.368000, lr 1.000000e-04
starting iteration 360
starting iteration 370
starting iteration 380
starting iteration 390
starting iteration 400
Finished epoch 0 / 1: cost 1.933786, train: 0.407000, val 0.394000, lr 1.000000e-04
starting iteration 410
```

```
starting iteration 420
starting iteration 430
starting iteration 440
starting iteration 450
Finished epoch 0 / 1: cost 1.691737, train: 0.448000, val 0.464000, lr 1.000000e-04
starting iteration 460
starting iteration 470
starting iteration 480
starting iteration 490
starting iteration 500
Finished epoch 0 / 1: cost 1.386387, train: 0.444000, val 0.432000, lr 1.000000e-04
starting iteration 510
starting iteration 520
starting iteration 530
starting iteration 540
starting iteration 550
Finished epoch 0 / 1: cost 1.736360, train: 0.471000, val 0.457000, lr 1.000000e-04
starting iteration 560
starting iteration 570
starting iteration 580
starting iteration 590
starting iteration 600
Finished epoch 0 / 1: cost 1.420449, train: 0.457000, val 0.480000, lr 1.000000e-04
starting iteration 610
starting iteration 620
starting iteration 630
starting iteration 640
starting iteration 650
Finished epoch 0 / 1: cost 1.539274, train: 0.433000, val 0.407000, lr 1.000000e-04
starting iteration 660
starting iteration 670
starting iteration 680
starting iteration 690
starting iteration 700
Finished epoch 0 / 1: cost 1.752273, train: 0.514000, val 0.494000, lr 1.000000e-04
starting iteration 710
starting iteration 720
starting iteration 730
starting iteration 740
starting iteration 750
Finished epoch 0 / 1: cost 1.709746, train: 0.477000, val 0.506000, lr 1.000000e-04
starting iteration 760
starting iteration 770
starting iteration 780
starting iteration 790
starting iteration 800
Finished epoch 0 / 1: cost 1.855160, train: 0.440000, val 0.460000, lr 1.000000e-04
starting iteration 810
```

```
starting iteration 820
starting iteration 830
starting iteration 840
starting iteration 850
Finished epoch 0 / 1: cost 1.565836, train: 0.542000, val 0.503000, lr 1.000000e-04
starting iteration 860
starting iteration 870
starting iteration 880
starting iteration 890
starting iteration 900
Finished epoch 0 / 1: cost 1.623981, train: 0.477000, val 0.467000, lr 1.000000e-04
starting iteration 910
starting iteration 920
starting iteration 930
starting iteration 940
starting iteration 950
Finished epoch 0 / 1: cost 1.957127, train: 0.493000, val 0.506000, lr 1.000000e-04
starting iteration 960
starting iteration 970
Finished epoch 1 / 1: cost 1.945312, train: 0.510000, val 0.474000, lr 9.500000e-05
finished optimization. best validation accuracy: 0.506000
```

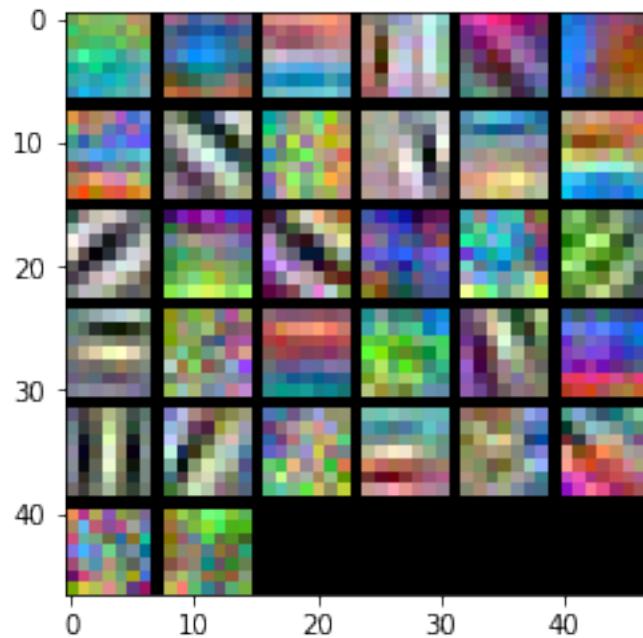
6 Visualize weights

We can visualize the convolutional weights from the first layer. If everything worked properly, these will usually be edges and blobs of various colors and orientations.

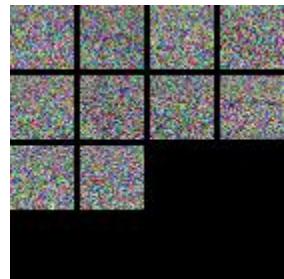
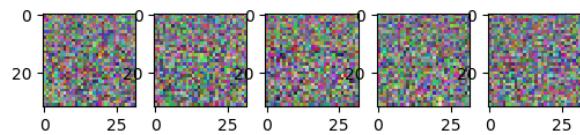
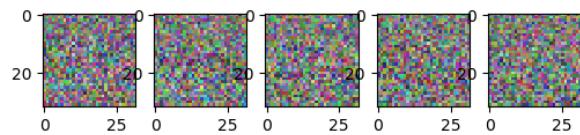
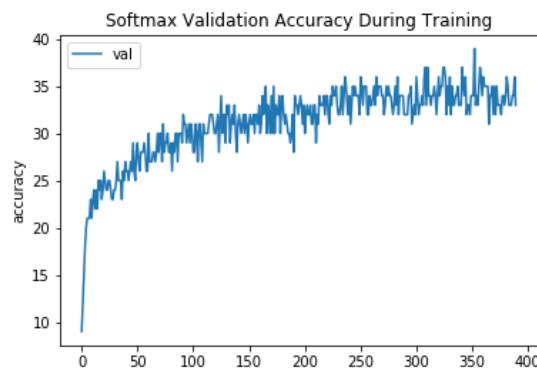
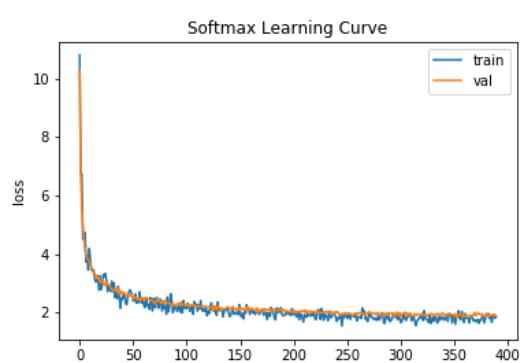
```
In [8]: from cs231n.vis_utils import visualize_grid

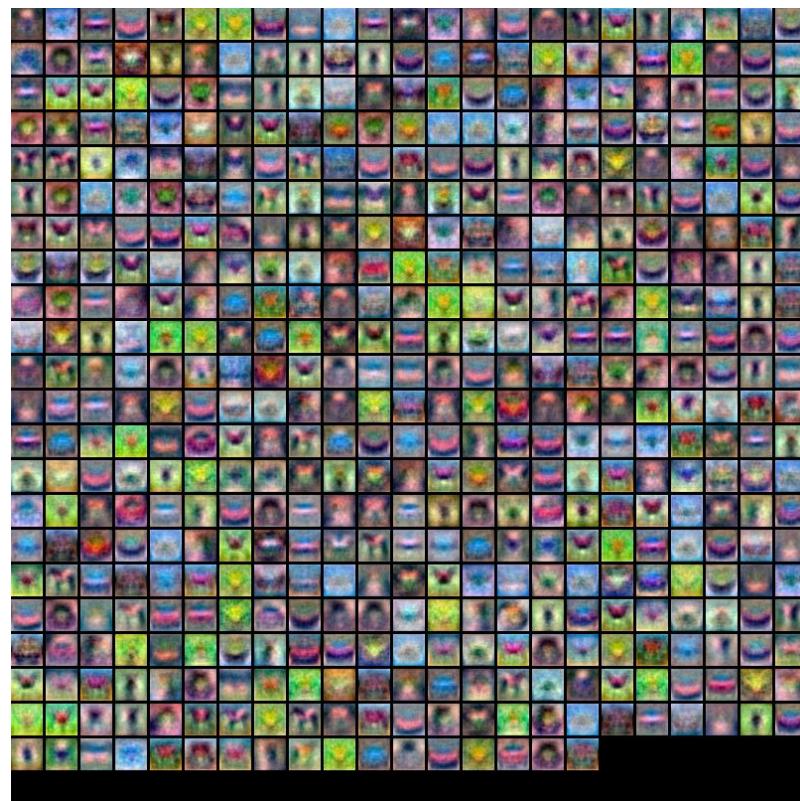
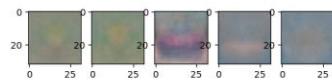
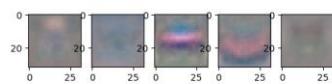
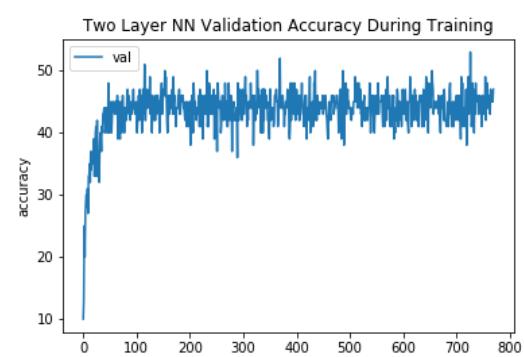
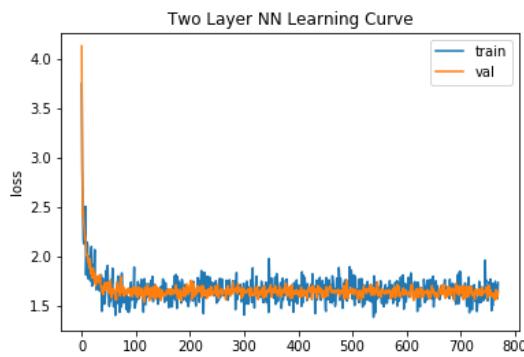
grid = visualize_grid(best_model['W1'].transpose(0, 2, 3, 1))
plt.imshow(grid.astype('uint8'))
```

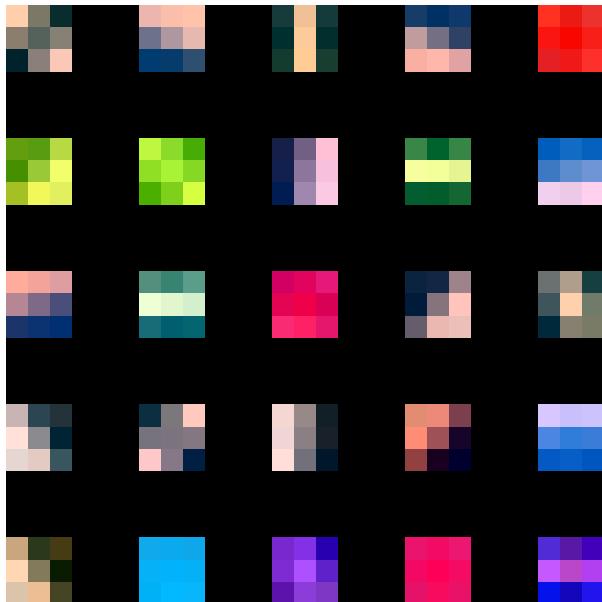
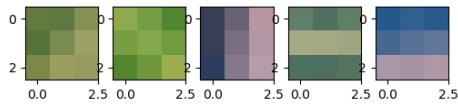
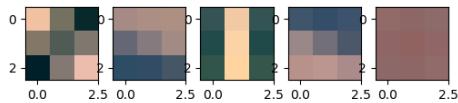
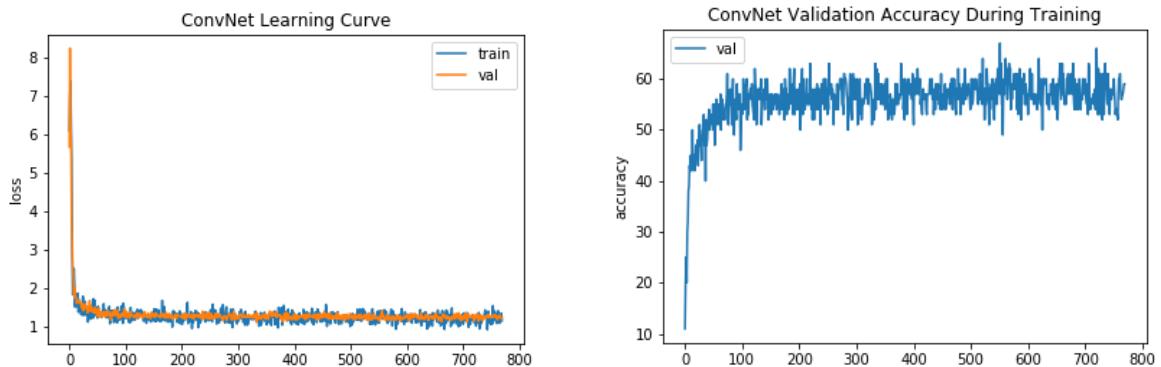
Out[8]: <matplotlib.image.AxesImage at 0x1f800126be0>

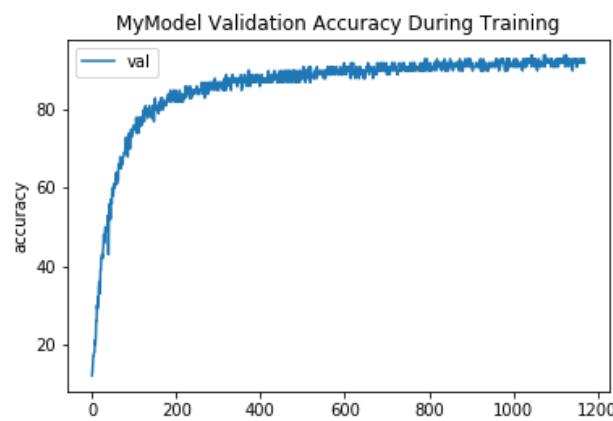
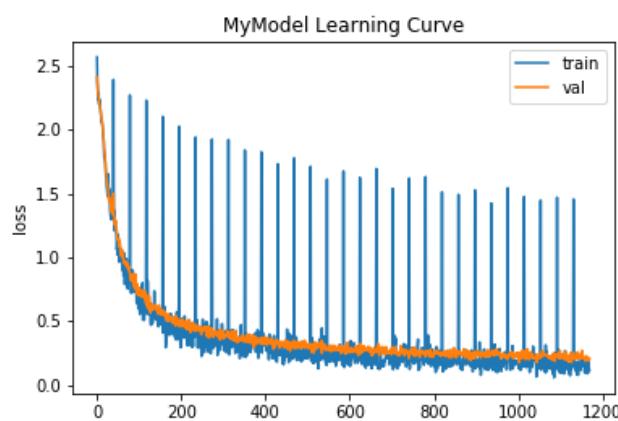


In []:









Name: James Hahn
Email ID: jhahn37@gatech.edu
Best Accuracy: 91%

Procedure:

When I first saw this task, I didn't think twice. I knew this was a perfect task for transfer learning, since the ImageNet task is similar in nature, just with larger images and more classes. The benefit CIFAR10 has over ImageNet is that CIFAR10 has more samples per class. So, I went to PyTorch's website and found the pre-trained models (<https://pytorch.org/docs/stable/torchvision/models.html>). Lucky for me, there is a table with all the top-1 and top-5 accuracies on ImageNet. I basically just went down the list, tried ResNet18, VGG16, AlexNet, DenseNet161, DenseNet121. Initially, I was training this on my desktop and it was extremely slow, so I was training on either 1, 5, or 10 epochs. My accuracies were pretty poor, but I reached 76% accuracy on DenseNet161. It's important to note in this initial phase, I froze all the pre-trained layers and replaced the last FC layer with my own layer to output 10 scores for the CIFAR10 dataset; also, I resized the images from 32x32 to 224x224 for the pre-trained models. I didn't really know how to get up to a higher accuracy, most of my models were stuck around 76%. As such, I did a lot of searches online, specifically for CIFAR10 papers and Github repositories. I found a lot of articles using either their own network or pre-trained networks, but a lot of the articles were pretty annoying to read because they discussed resizing images and data augmentation and I was too lazy to do any of that. Eventually, I went to one of my other machine learning classes, and found a student from the deep learning class that had already beaten the TA's benchmark. I asked him what his approach was, and to my shock, he used pre-trained networks as well, but he was training on a significant higher number of epochs. I was kind of confused, because my computer wasn't even using all of my GPU VRAM, but eventually found out he was training on Google Colab. I went back to my apartment, plugged my code into Colab, used their GPU, and the training was at least 5x faster. I was able to train for 30 or 50 epochs and run through the pre-trained models once again. I eventually found some with 81%, 81%, 82%, and 83% accuracy. The DenseNet161 from earlier was the one with 83% accuracy, so I just stuck with that one. Initially, I was just going to stop here since it accomplished the goal of beating the TA, but I wanted to be as safe as possible when it was evaluated on the full test set. After a lot of debugging, I tested a few more models and hyperparameters and eventually found ResNet18 to work best, accomplishing 91% on the EvalAI set. Additionally, I added data augmentation by randomly performing a horizontal flip of the images with a probability of 0.5. Finally, from before, I froze all the weights, but this time around, I brought in the pre-trained weights and continued to train the entire network without any frozen layers. I think it's important to note I tried three really different approaches to finding my final solution: prior knowledge and experience in DL/transfer learning, online searches of previous results on CIFAR and similar datasets, and person-to-person communication to hear their struggles/experiences. Overall, pretty well-rounded, but ended up with data augmentation + pre-trained model.