

Homework 3

The homework consists of two parts: theoretical part (5 pts) and coding part (25 pts).

- All theoretical questions must be answered in your own words, do not copy-paste text from the internet. Points can be deducted for terrible formatting or incomprehensible English.
- Code must be commented. If you use code you found online, you have to add the link to the source you used. There is no penalty for using outside sources as long as you convince us you understand the code.
- If you have hard time writing equations in markdown/latex in this notebook, write them on paper, take a good quality photo and add that photo to the notebook (make sure the file is included in your upload)

You can earn up to 5 bonus points by tuning the hyperparameters and reaching test accuracies above 52%

Once completed zip the entire directory containing this exercise and upload it to <https://courses.cs.ut.ee/2021/spring/Main/Practices>.

For background reading see <http://cs231n.github.io/optimization-2/>.

For more information about how to derive rules for matrix derivatives (you are not required to work through these, this is more of an extra information for those who are interested in how these are derived):

- https://courses.cs.ut.ee/LAT.02.001/2022_spring/uploads/Main/matrix_calc.pdf
- <http://vision.stanford.edu/teaching/cs231n/handouts/derivatives.pdf>
- <https://web.stanford.edu/class/cs24n/readings/gradients-notes.pdf>

If you did this homework together with classmates, please write here their names (answers still have to be your own!).

Names(s): Soni

Part 1: Lecture Materials (5 pts)

These theoretical questions are about the material covered in the lecture about "Back-propagation"

Task 1.1: Derivatives of scalars, vectors and matrices (1pt)

In this task we ask you what are the dimensions of the gradient vector/matrix when taking the derivative of the object noted in the topmost row with respect to the object given in the first column. You can find help from <http://cs231n.stanford.edu/handouts/derivatives.pdf>, more precisely the first equation (the Jacobian) on page 3 (section 1.3) describing $\frac{\partial}{\partial x}$

Please fill in the table below:

$\frac{\partial}{\partial}$	scalar	vector $\in \mathbb{R}^n$
scalar	M-Dimension	MxN
vector $\in \mathbb{R}^n$	KxM	MxN

Task 1.2: General rule of backprop (1pt)

You are given a feed-forward network with very many hidden layers. For the simplicity of this exercise, consider that there is no non-linearity (no activation function) applied in the nodes. We note as W_i the weight matrix that is between layers i and $i+1$. We note as h_i the hidden node activations in layer i . So $h_{i+1} = W_i h_i$.

Considering that you already know the derivatives with respect to the hidden node activations in layer $N+1$, $\frac{\partial L}{\partial h_{N+1}}$, find the the following derivatives:

$$\begin{aligned}\frac{\partial L}{\partial W_N} &= \frac{\partial L}{\partial h_{N+1}} \frac{\partial h_{N+1}}{\partial W_N} = \frac{\partial L}{\partial h_{N+1}} W_N^T \frac{\partial h_{N+1}}{\partial W_N} \\ \frac{\partial L}{\partial h_N} &= \frac{\partial L}{\partial h_{N+1}} \frac{\partial h_{N+1}}{\partial h_N} = \frac{\partial L}{\partial h_{N+1}} W_N^T \\ \frac{\partial L}{\partial W_{N-1}} &= \frac{\partial L}{\partial h_N} \frac{\partial h_N}{\partial W_{N-1}} = \frac{\partial L}{\partial h_N} W_{N-1}^T \frac{\partial h_N}{\partial W_{N-1}} \\ \frac{\partial L}{\partial h_{N-1}} &= \frac{\partial L}{\partial h_N} \frac{\partial h_N}{\partial h_{N-1}} = \frac{\partial L}{\partial h_N} W_{N-1}^T \\ \frac{\partial L}{\partial W_{N-2}} &= \frac{\partial L}{\partial h_{N-1}} \frac{\partial h_{N-1}}{\partial W_{N-2}} = \frac{\partial L}{\partial h_{N-1}} W_{N-2}^T \frac{\partial h_{N-1}}{\partial W_{N-2}}\end{aligned}$$

Express the derivatives in function of the known derivative $\frac{\partial L}{\partial h_{N+1}}$, weight matrices W_i and hidden node activations h_i .

Knowing these rules will help you to solve the task:

Suppose

$$L = f(Y), \quad Y = XW.$$

Then

$$\begin{aligned}\frac{\partial L}{\partial X} &= \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial X} = \frac{\partial L}{\partial Y} \frac{\partial XW}{\partial X} = \frac{\partial L}{\partial Y} W^T \\ \frac{\partial L}{\partial W} &= \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial W} = \frac{\partial L}{\partial Y} \frac{\partial XW}{\partial W} = \frac{\partial L}{\partial Y} X \frac{\partial W}{\partial W} = \frac{\partial L}{\partial Y} X\end{aligned}$$

Your Answer:

$$\begin{aligned}\frac{\partial L}{\partial W_N} &= \frac{\partial L}{\partial h_{N+1}} \frac{\partial h_{N+1}}{\partial W_N} = \frac{\partial L}{\partial h_{N+1}} W_N^T \frac{\partial h_{N+1}}{\partial W_N} \\ \frac{\partial L}{\partial h_N} &= \frac{\partial L}{\partial h_{N+1}} \frac{\partial h_{N+1}}{\partial h_N} = \frac{\partial L}{\partial h_{N+1}} W_N^T \\ \frac{\partial L}{\partial W_{N-1}} &= \frac{\partial L}{\partial h_N} \frac{\partial h_N}{\partial W_{N-1}} = \frac{\partial L}{\partial h_N} W_{N-1}^T \frac{\partial h_N}{\partial W_{N-1}} \\ \frac{\partial L}{\partial h_{N-1}} &= \frac{\partial L}{\partial h_N} \frac{\partial h_N}{\partial h_{N-1}} = \frac{\partial L}{\partial h_N} W_{N-1}^T \\ \frac{\partial L}{\partial W_{N-2}} &= \frac{\partial L}{\partial h_{N-1}} \frac{\partial h_{N-1}}{\partial W_{N-2}} = \frac{\partial L}{\partial h_{N-1}} W_{N-2}^T \frac{\partial h_{N-1}}{\partial W_{N-2}}\end{aligned}$$

Task 1.3: Complexity task 1 (1pt)

How many multiplications and how many additions/subtractions does it take to calculate:

- the product of a vector $v \in \mathbb{R}^n$ and a matrix $M \in \mathbb{R}^{n \times m}$ (i.e. complexity of doing $v^T M$)?

Your Answer: fill this in.

Task 1.4: Complexity task 2 (2pt)

Imagine that you have a neural network with one hidden layer with 100 nodes (a "two-layer network", because it has two weight matrices), and you are using ReLU activation function at the nodes (for both layers). The inputs to the network are RGB images of size $32 \times 32 \times 3$. The output is 10 probability values (obtained by passing the results through Softmax function).

How many scalar operations (additions, subtractions, multiplications, divisions, exponentiations, logarithms and comparisons) does it take to do:

- perform a forward pass on one data point (to calculate the cross-entropy loss)

NB!

- Don't forget the bias (you can think of it as an extra dimension in the input).
- The number of operations for the softmax depends on how you would "implement" it. Is the denominator calculated once (efficient way) or separately for each value (inefficient way). You can decide yourself which way you use it but please state it when reporting the operation count.

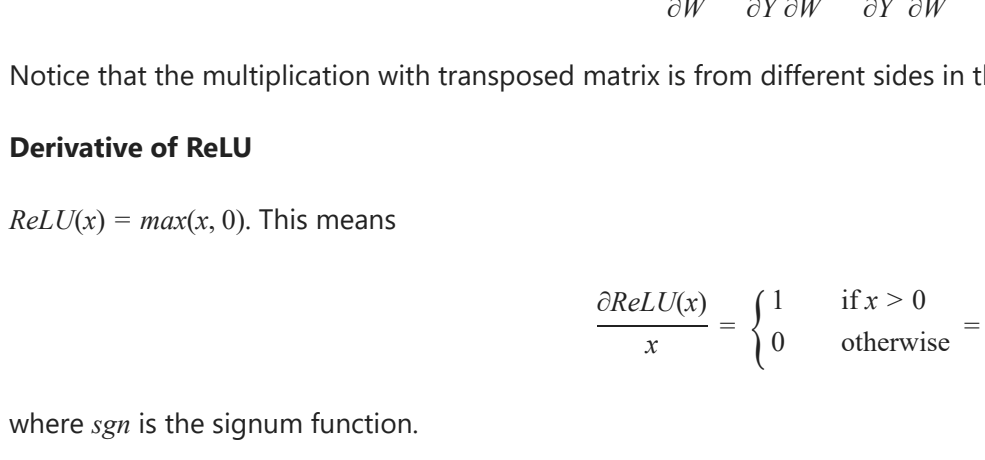
Your Answer: fill this in.

Part 2: Practical Tasks (25 pts)

Applied Theory

In this exercise we have a neural network with one hidden layer. This time we derive gradients slightly differently.

- we use batch of samples instead of one sample.
- we use matrix calculus for gradients.
- we add L2 regularization to the loss.



Notation:

- X is $N \times D$ input matrix, where N is number of samples in batch and D is number of features.
- Y is $N \times C$ one-hot coded target matrix, where N is number of samples in batch and C is number of classes.
- c is N -dimensional vector of correct classes for all samples, $c_i \in \{1, \dots, C\}$.
- $W^{(1)}$ is $D \times M$ weight matrix of the first layer, where D is the number of features and M is number of hidden nodes.
- $W^{(2)}$ is $M \times C$ weight matrix of the second layer, where M is the number of hidden nodes and C is the number of classes.
- $b^{(1)}$ is bias vector of the first layer with dimension of $1 \times M$.
- $b^{(2)}$ is bias vector of the second layer with dimension of $1 \times C$.
- 1_N is a column vector of ones of length N with dimension $N \times 1$.
- $\text{ReLU}(x) = \max(x, 0)$ is rectified linear unit activation function.
- $\text{softmax}(A)$ converts activations into probabilities row by row $P_j = \frac{e^{A_{ij}}}{\sum_k e^{A_{ik}}}$.

Feed-forward pass:

$$\begin{aligned}A^{(1)} &= XW^{(1)} + 1_N b^{(1)} & (N \times M) \\ H^{(1)} &= \text{ReLU}(A^{(1)}) & (N \times M) \\ A^{(2)} &= H^{(1)}W^{(2)} + 1_N b^{(2)} & (N \times C) \\ P &= \text{softmax}(A^{(2)}) & (N \times C) \\ L_{CE} &= -\frac{1}{N} \sum_{i=1}^N Y_i \odot \log P_i = -\frac{1}{N} \sum_{i=1}^N \log P_{c_i} & (\text{scalar})\end{aligned}$$

Here AB is matrix multiplication and $A \odot B$ is element-wise multiplication.

$\frac{1}{N}$ in the loss function produces mean instance loss, which is good because then loss value does not depend on batch size. If we had used sum, both loss value and gradients would have different magnitude for different batch sizes and we would need to adapt learning rate for each batch size.

Backward pass:

Derivatives of matrix multiplication

Derivatives of matrix multiplication work very similarly to the scalar case. However, as the order of elements in product matters ($AB \neq BA$), pay attention to transposes and also on which side of the overall expression the derivative ends up. Suppose

$$L = f(Y), \quad Y = XW.$$

Then

$$\begin{aligned}\frac{\partial L}{\partial X} &= \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial X} = \frac{\partial L}{\partial Y} \frac{\partial XW}{\partial X} = \frac{\partial L}{\partial Y} W^T \\ \frac{\partial L}{\partial W} &= \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial W} = \frac{\partial L}{\partial Y} \frac{\partial XW}{\partial W} = \frac{\partial L}{\partial Y} X \frac{\partial W}{\partial W} = \frac{\partial L}{\partial Y} X\end{aligned}$$

Notice that the multiplication with transposed matrix is from different sides in these two cases.

Derivative of ReLU

$\text{ReLU}(x) = \max(x, 0)$. This means

$$\frac{\partial \text{ReLU}(x)}{\partial x} = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases} = \text{sgn}(\text{ReLU}(x))$$

where sgn is the signum function.

Derivative of softmax

If $P = \text{softmax}(A)$ and $L = CE(P, Y)$ and N is the number of samples in the batch, then

$$\frac{\partial L}{\partial A} = \frac{1}{N}(P - Y)$$

Backward pass for $W^{(1)}$ using chain rule:

$$\begin{aligned}\frac{\partial L_{CE}}{\partial W^{(1)}} &= \frac{\partial L_{CE}}{\partial A^{(1)}} \frac{\partial A^{(1)}}{\partial W^{(1)}} & (\text{chain rule}) \\ &= X^T \frac{\partial L_{CE}}{\partial A^{(1)}} & (\text{matrix multiplication formula}) \\ &= X^T \frac{\partial L_{CE}}{\partial H^{(1)}} \frac{\partial H^{(1)}}{\partial W^{(1)}} & (\text{chain rule}) \\ &= X^T \left(\frac{\partial L_{CE}}{\partial H^{(1)}} \odot (A^{(1)} > 0) \right) & (\text{ReLU formula}) \\ &= X^T \left(\frac{\partial L_{CE}}{\partial A^{(2)}} \frac{\partial A^{(2)}}{\partial H^{(1)}} \odot (A^{(1)} > 0) \right) & (\text{chain rule}) \\ &= X^T \left(\left(-\frac{\partial L_{CE}}{\partial A^{(2)}} \right) (W^{(2)})^T \odot (A^{(1)} > 0) \right) & (\text{matrix multiplication formula}) \\ &= \frac{1}{N} X^T \left((P - Y) W^{(2)T} \right) \odot (A^{(1)} > 0) & (\text{softmax formula})\end{aligned}$$

Here $(A^{(1)} > 0)$ denotes a matrix whose value at row i and column j is one, if the corresponding value in matrix A is greater than 0, and 0 otherwise.

Task 2.1 (4pts)

Write down the partial derivatives of classification loss function with respect to other weights and biases. Verify that the dimensions match!

$$\begin{aligned}\frac{\partial L_{CE}}{\partial W^{(2)}} &= \dots & (M \times C) \\ \frac{\partial L_{CE}}{\partial b^{(2)}} &= \dots & (1 \times M) \\ \frac{\partial L_{CE}}{\partial b^{(1)}} &= \dots & (1 \times C)\end{aligned}$$

Hint:

- You can use the matrix multiplication derivative rules for the partial derivative for the bias (since it is also multiplication of two matrices, identity matrix and bias itself (matrix with 1 row)).
- You can reuse parts of the calculation of $\frac{\partial L_{CE}}{\partial W^{(1)}}$ (the example above).

L2 regularization

L2 regularization drives the weights to be small.

$$L_r = \sum_i (W_i)^2 + \sum_j (W_j)^2$$

In the final loss the classification loss and regularization loss are added together.

$$L = L_{CE} + \lambda L_r$$

The strenght of regularization is determined by regularization coefficient λ .

Task 2.2 (2pts):

Write down the partial derivatives of regularization loss with respect to weights:

$$\begin{aligned}\frac{\partial L_r}{\partial W^{(1)}} &= \dots & (D \times M) \\ \frac{\partial L_r}{\partial W^{(2)}} &= \dots & (M \times C)\end{aligned}$$

In [1]:

```
# A bit of setup
from future import print_function
import numpy as np
import matplotlib.pyplot as plt

from neural_net import TwoLayerNet

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

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

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

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

In [2]:

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

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

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

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

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

Forward pass: compute scores

Task 2.3 (2pts)

Open the file 'neural_net.py' and look at the method TwoLayerNet.loss. This function is very similar to the loss functions you have written for the Softmax exercise: 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 (activations in last layer, before softmax) for all inputs.

Hint: When implementing ReLU, use np.maximum, not np.max (np.max finds max along an axis of the matrix, but you need element-wise maximum).

In [3]:

```
scores = net.loss(X)
print('Your scores:')
print(scores)

correct_scores = np.asarray([
    [-0.81233741, -1.27654624, -0.70335995],
    [-0.17129677, -1.18803311, -0.47310444],
    [-0.51590475, -1.01354314, -0.8504215 ],
    [-0.15419291, -0.48629638, -0.52901952],
    [-0.00618733, -0.12435261, -0.15226949]])
print(correct_scores)
print()

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

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

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

Difference between your scores and correct scores:
3.660270745909845e-08
```

Forward pass: compute loss

Task 2.4 (2pts)

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

In [4]:

```
loss, _ = net.loss(X, y, reg=0.05)
correct_loss = 1.30378789135

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

Your loss and correct loss:
0.01896541560606335
```

Backward pass

Task 2.5 (6pts)

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

In [5]:

```
from gradient_check import eval_numerical_gradient

# Use numeric gradient checking to check your implementation of the backward pass.
# If your implementation is correct, the difference between the numeric and
# analytic gradients should be less than 1e-8 for each of W1, b1, W2, b2.

loss, grads = net.loss(X, y, reg=0.05)

# these should all be less than 1e-8 or so
for param_name in grads:
    f = lambda W: net.loss(X, y, reg=0.05, dtype='float64')
    param_grad_num = eval_numerical_gradient(f, [param_name], verbose=False)
    print('max relative error: %e' % (param_name, rel_error(param_grad_num, grads[param_name])))

W2 max relative error: 3.440708e-09
b2 max relative error: 4.447646e-11
W1 max relative error: 6.890046e-09
b1 max relative error: 2.738421e-09
```

Train the network

Task 2.6 (3pts)

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

Once you have implemented the method, run the code below to train a two-layer network on toy data. You should achieve a training loss less than 0.2.

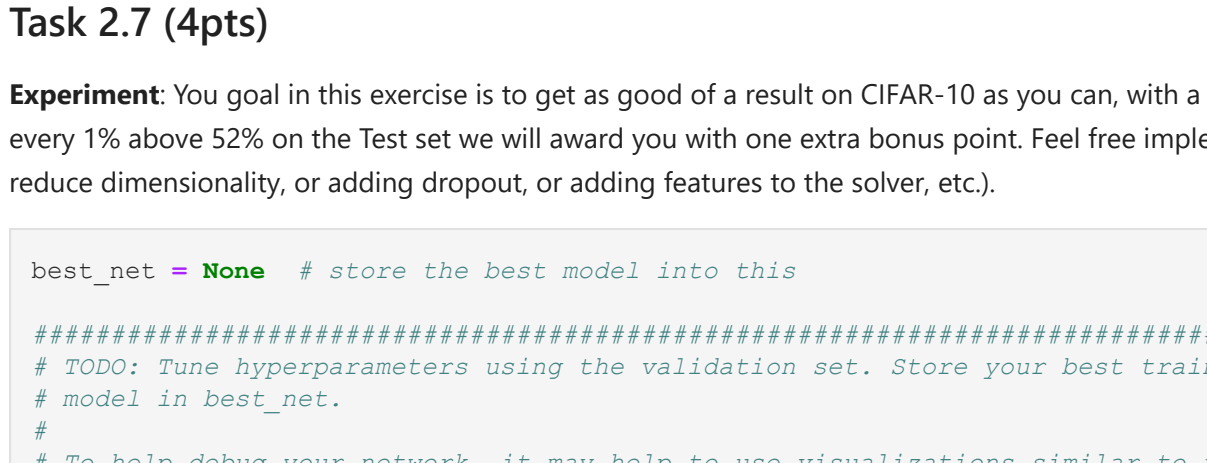
In [6]:

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

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

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

Final training loss: 0.017143643532923757
```



Load the data

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

In [7]:

```
from 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
    used for the SVM, but condensed to a single function.
    """
    # Load the raw CIFAR-10 data
    cifar10_dir = './datasets/cifar-10-batches-py'
    X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

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

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

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

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

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

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

Train a network

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

In [8]:

```
input_size = 32 * 32 * 3
hidden_size = 50
num_classes = 10
net = TwoLayerNet(input_size, hidden_size, num_classes)

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

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

iteration 0 / 1000: loss 2.302762
iteration 100 / 1000: loss 2.302358
iteration 200 / 1000: loss 2.297404
iteration 300 / 1000: loss 2.298971
iteration 400 / 1000: loss 2.202975
iteration 500 / 1000: loss 2.116816
iteration 600 / 1000: loss 2.040789
iteration 700 / 1000: loss 1.985711
iteration 800 / 1000: loss 2.003726
iteration 900 / 1000: loss 1.948076
Validation accuracy: 0.287
```

Debug the training

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

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

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

In [9]:

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

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


In [10]:

```
from vis_utils import visualize_grid

# Visualize the weights of the network
W1 = net.params['W1']
W1 = W1.reshape(32, 32, 3, -1).transpose(3, 0, 1, 2)
plt.imshow(visualize_grid(W1, padding=3, stype='uint8'))
plt.gca().axis('off')
plt.show()

show_net_weights(net)
```


Tune your hyperparameters

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

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

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

