

Due Friday, 7 Feb 2025, by 11:59pm to Gradescope.
 100 points total.

1. (15 points) **Backpropagation for autoencoders.** In an autoencoder, we seek to reconstruct the original data after some operation that reduces the data's dimensionality. We may be interested in reducing the data's dimensionality to gain a more compact representation of the data.

For example, consider $\mathbf{x} \in \mathbb{R}^n$. Further, consider $\mathbf{W} \in \mathbb{R}^{m \times n}$ where $m < n$. Then $\mathbf{W}\mathbf{x}$ is of lower dimensionality than \mathbf{x} . One way to design \mathbf{W} so that $\mathbf{W}\mathbf{x}$ still contains key features of \mathbf{x} is to minimize the following expression

$$\mathcal{L} = \frac{1}{2} \|\mathbf{W}^T \mathbf{W}\mathbf{x} - \mathbf{x}\|^2$$

with respect to \mathbf{W} . (To be complete, autoencoders also have a nonlinearity in each layer, i.e., the loss is $\frac{1}{2} \|f(\mathbf{W}^T f(\mathbf{W}\mathbf{x})) - \mathbf{x}\|^2$. However, we'll work with the linear example.)

- (a) (3 points) In words, describe why this minimization finds a $\boxed{\mathbf{W}}$ that ought to preserve information about \mathbf{x} .
 - (b) (3 points) Draw the computational graph for \mathcal{L} . **Hint:** You can set up the computational graph to this problem in a way that will allow you to solve for part (d) without taking 4D tensor derivative.
 - (c) (3 points) In the computational graph, there should be two paths to \mathbf{W} . How do we account for these two paths when calculating $\nabla_{\mathbf{W}} \mathcal{L}$? Your answer should include a mathematical argument.
 - (d) (6 points) Calculate the gradient: $\nabla_{\mathbf{W}} \mathcal{L}$.
2. (20 points) **Backpropagation for Gaussian-process latent variable model.** (Optional for students in C147: Please write 'I am a C147 student' in the solution and you will get full credit for this problem). An important component of unsupervised learning is visualizing high-dimensional data in low-dimensional spaces. One such nonlinear algorithm to do so is from Lawrence, NIPS 2004, called GP-LVM. GP-LVM optimizes the maximum-likelihood of a probabilistic model. We won't get into the details here, but rather to the bottom line: in this paper, a log-likelihood has to be differentiated with respect to a matrix to derive the optimal parameters.

To do so, we will apply the chain rule for multivariate derivatives via backpropagation. The log-likelihood is:

$$\mathcal{L} = -c - \frac{D}{2} \log |\mathbf{K}| - \frac{1}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y} \mathbf{Y}^T)$$

where $\mathbf{K} = \alpha \mathbf{XX}^T + \beta^{-1} \mathbf{I}$ and c is a constant. The $|.|$ symbol in this context refers to the determinant of a matrix. To solve this, we'll take the derivatives with respect to the two terms with dependencies on \mathbf{X} :

$$\begin{aligned}\mathcal{L}_1 &= -\frac{D}{2} \log |\alpha \mathbf{XX}^T + \beta^{-1} \mathbf{I}| \\ \mathcal{L}_2 &= -\frac{1}{2} \text{tr} ((\alpha \mathbf{XX}^T + \beta^{-1} \mathbf{I})^{-1} \mathbf{YY}^T)\end{aligned}$$

Hint: To receive full credit, you will be required to show all work. You may use the following matrix derivative without proof:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{K}} = -\mathbf{K}^{-T} \frac{\partial \mathcal{L}}{\partial \mathbf{K}^{-1}} \mathbf{K}^{-T}.$$

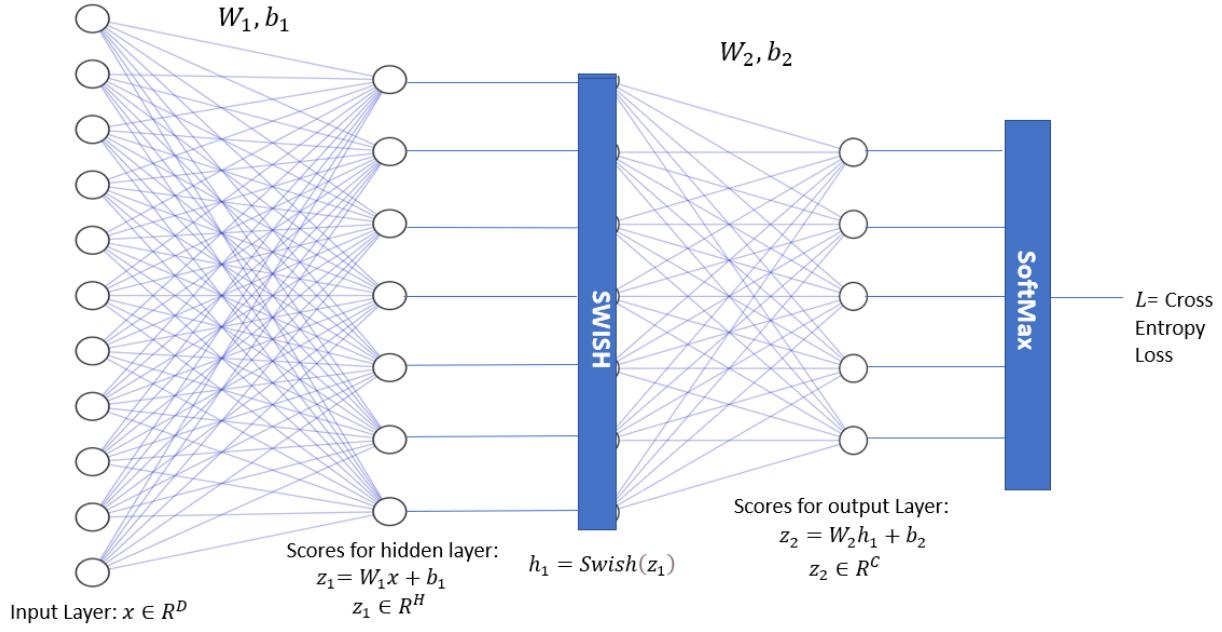
Also, consider the matrix operation, $\mathbf{Z} = \mathbf{XY}$. If we have an upstream derivative, $\partial \mathcal{L} / \partial \mathbf{Z}$, then backpropagate the derivatives in the following way:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mathbf{X}} &= \frac{\partial \mathcal{L}}{\partial \mathbf{Z}} \mathbf{Y}^T \\ \frac{\partial \mathcal{L}}{\partial \mathbf{Y}} &= \mathbf{X}^T \frac{\partial \mathcal{L}}{\partial \mathbf{Z}}\end{aligned}$$

- (a) (3 points) Draw a computational graph for \mathcal{L}_1 .
 - (b) (6 points) Compute $\frac{\partial \mathcal{L}_1}{\partial \mathbf{X}}$.
 - (c) (3 points) Draw a computational graph for \mathcal{L}_2 .
 - (d) (6 points) Compute $\frac{\partial \mathcal{L}_2}{\partial \mathbf{X}}$.
 - (e) (2 points) Compute $\frac{\partial \mathcal{L}}{\partial \mathbf{X}}$.
3. (15 points) **NNDL to the rescue!!**

It looks like a calm Monday morning and you are almost done with NNDL HW for the week (sigh)! but then suddenly (ting ting ...) your phone starts buzzing, you pick the call and the person from the other end sounds tense, the person exclaims ... There is a national Emergency !!: *7 different Pandora creature species (from Avatar) have been spotted in 1000's of numbers across various places in the country . They are having a hard time to adjust to earth's climate and are causing chaos. As a result there has been a Power outage in many cities. Luckily LA is an exception. UCLA's engineering division is helping out with this emergency, and you have been summoned to contribute to the same.* You quickly take a bird to the secret facility and meet with director in charge of this operation. The director gives you a dataset consisting of images of these creatures along with their species type and instructs you to design a machine learning model to classify the images into species type. The only design constraint that the director has imposed is that the model should not have a very large number of parameters because some of UCLA's compute facilities are overloaded due to the power outages.

You have just learned about Fully connected neural networks (FC net) in class and decide to use it for accomplishing the task. To satisfy the design constraint, you have decided to build a 2-layer FC net and train it using the provided dataset. The trained model will not only



enable you to classify the images into species type but the hidden representations (outputs of intermediate layers) can be used to analyze the various properties of the species. A pictorial representation of the 2-layer FC net is shown above:

In the architecture shown, D represents the number of neurons in input layer, H represents the number of neurons in hidden layer , C represents the number of neurons in the output layer (in our design $C = 7$). The output is then passed through a softmax classifier. Although we learned about the ReLu activation in class, but we decided to use the Swish activation function (introduced by google brain) for the hidden layer. Swish activation function for any scalar input k is defined as,

$$\text{swish}(k) = \frac{k}{1 + e^{-k}} = k\sigma(k),$$

where, $\sigma(k)$, is the sigmoid activation function you have seen in lectures

You will train the 2-layer FC net using gradient descent and for that you will need to compute the gradients. For the gradient computations, you are allowed to keep your final answer in terms of $\frac{\partial L}{\partial z_2}$.

- (a) (3 points) Draw the computational graph for the 2-layer FC net.
 - (b) (5 points) Compute $\nabla_{W_2} L$, $\nabla_{b_2} L$.
 - (c) (7 points) Compute $\nabla_{W_1} L$, $\nabla_{b_1} L$.
4. (30 points) **2-layer neural network.**

Complete the two-layer neural network Jupyter notebook. Print out the entire notebook and relevant code and submit it as a pdf to gradescope. Download the CIFAR-10 dataset, as you did in HW #2.

5. (20 points) **General FC neural network.**

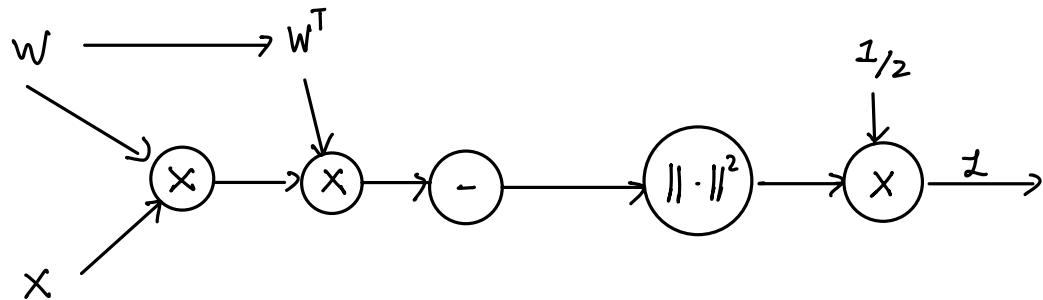
Complete the FC Net Jupyter notebook. Print out the entire notebook and relevant code and submit it as a pdf to gradescope.

1. a) $x \in \mathbb{R}^n$ $w \in \mathbb{R}^{m \times n}$ ($m < n$)

$$L = \frac{1}{2} \| w^T w x - x \|_2^2$$

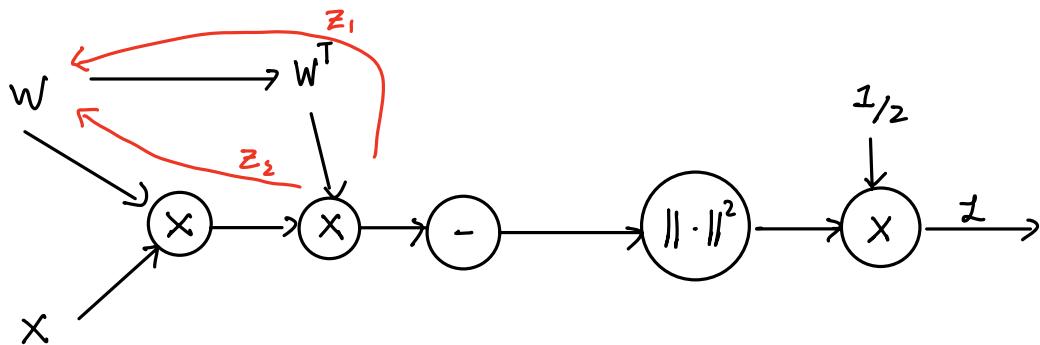
This minimization preserves the information about x by projecting important features into a lower dimensional space constrained by L2-distance. This preserves information about x because w is optimized to retain the most relevant features, minimizing the loss of critical information in the process.

b) Computational graph



$$d = w x ; \quad m = w^T w x - x ; \quad k = \|m\|^2 ; \quad g = \frac{1}{2} k$$

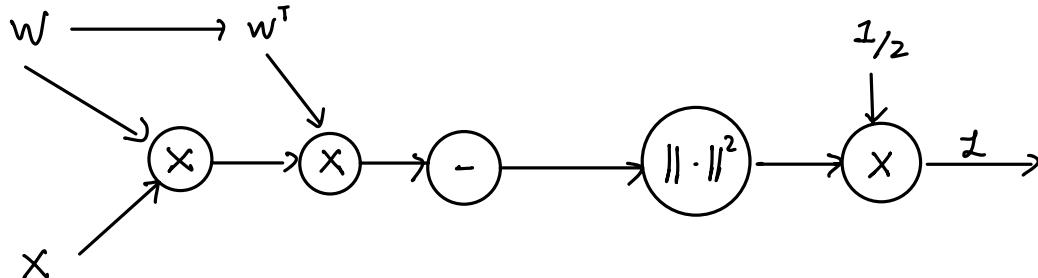
c)



- we have 1 path going through w^T and the other one going through the multiplication gate.
- we can take the gradient for each path individually & add them to find $\nabla_w L$

By chain rule $\nabla_w \mathcal{L} = \frac{\partial z_1}{\partial w} \frac{\partial \mathcal{L}}{\partial z_1} + \frac{\partial z_2}{\partial w} \frac{\partial \mathcal{L}}{\partial z_2}$
 & the law of
 total derivatives $z_1 = w ; z_2 = w, \frac{\partial z_1}{\partial w} = \frac{\partial z_2}{\partial w} = 1$
 we have

d)



$$l = w^T x \quad j \quad m = w^T w x - x \quad ; \quad k = \|m\|^2 \quad ; \quad g = \frac{1}{2} k$$

$w \in \mathbb{R}^m \quad m \in \mathbb{R}^n$

$$\frac{\partial L}{\partial g} = \frac{1}{2} \quad L = \frac{1}{2} m^T m$$

$$\frac{\partial L}{\partial k} = \frac{\partial g}{\partial k} \frac{\partial L}{\partial g} = \frac{1}{2} \cdot 2$$

$$\frac{\partial L}{\partial m} = m = w^T w x - x$$

$$\frac{\partial L}{\partial l} = \frac{\partial l}{\partial m} \frac{\partial L}{\partial m} = w m \quad \frac{\partial L}{\partial w^T} = \frac{\partial m}{\partial w^T} \frac{\partial L}{\partial m} = m (w x)^T$$

$$\frac{\partial L}{\partial w} = \frac{\partial m}{\partial w} \frac{\partial L}{\partial m} = w m x^T$$

$$\therefore \frac{\partial L}{\partial w} = w ((w^T w x - x) x^T) + w x (w^T w x - x)^T$$

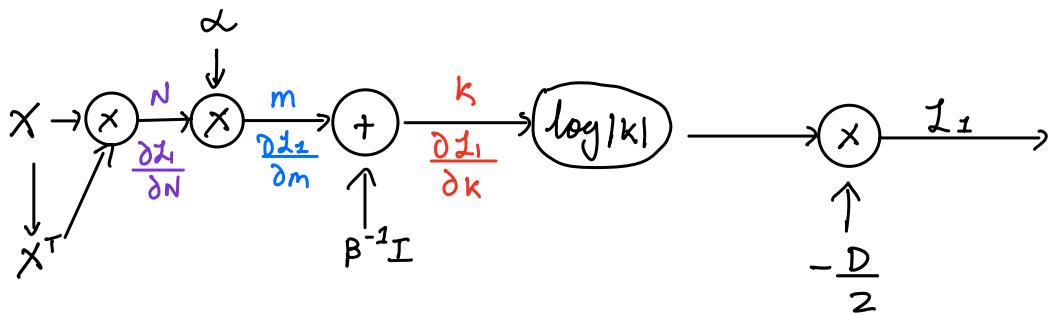
$$2. \quad a) \text{ & } b) \quad \mathcal{L} = -c - \frac{D}{2} \log |K| - \frac{1}{2} \text{tr}(K^{-1} Y Y^T)$$

$$K = \alpha X X^T + \beta^{-1} I$$

$$\mathcal{L}_1 = -\frac{D}{2} \log |\alpha X X^T + \beta^{-1} I|$$

$$\mathcal{L}_2 = -\frac{1}{2} \text{tr}((\alpha X X^T + \beta^{-1} I)^{-1} Y Y^T)$$

$$\text{Note: } \frac{\partial \mathcal{L}_1}{\partial K} = -K^{-T} \frac{\partial \mathcal{L}_1}{\partial K^{-1}} K^{-T}$$



$$\mathcal{L}_1 = -\frac{D}{2} \log |K| = \frac{-D}{2} \log (\det(K))$$

$$\frac{\partial \mathcal{L}_1}{\partial K} = -\frac{D}{2} (K^{-1})^T = -\frac{D}{2} (K^T)^{-1}$$

$$\mathcal{L}_1 = \frac{D}{2} \log(K) \quad \text{④ gradient flow} \quad \frac{\partial \mathcal{L}_1}{\partial m} = -\frac{D}{2} (K^T)^{-1}$$

$$K = mI \quad \text{⑤ Switch gradients} \quad \frac{\partial \mathcal{L}_1}{\partial N} = -\frac{\alpha D}{2} (K^T)^{-1}$$

$$N = X X^T; \quad \frac{\partial \mathcal{L}_1}{\partial X} = \frac{\partial \mathcal{L}_1}{\partial N} (X^T)^T = -\frac{\alpha D}{2} (K^T)^{-1} \cdot X$$

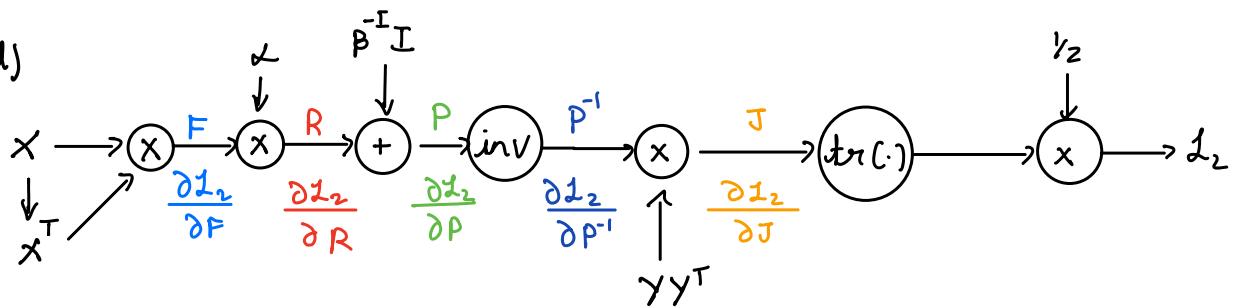
$$\frac{\partial \mathcal{L}_1}{\partial X^T} = X^T \frac{\partial \mathcal{L}_1}{\partial N} = X^T \left[-\frac{\alpha D}{2} (K^T)^{-1} \right] = -\frac{\alpha D}{2} (K^{-1} X)^T$$

Adding the gradient path through $X \otimes X^T$ we get

$$\frac{\partial \mathcal{L}_1}{\partial X} = -\frac{\alpha D}{2} (K^T)^{-1} X - \left[\frac{\alpha D}{2} (K^{-1} X)^T \right]^T$$

K is symmetric \therefore $\frac{\partial \mathcal{L}_1}{\partial X} = -\alpha D K^{-1} X$

c) & d)



$$P = \alpha x x^T + \beta^{-1} I$$

$$L_2 = -\frac{1}{2} \text{tr}(J)$$

$$F = Y Y^T$$

$$\frac{\partial L_2}{\partial J} = -\frac{1}{2} I$$

$$\frac{\partial L_2}{\partial x} = -\alpha P^{-1} \frac{\partial J}{\partial P^{-1}} P^{-1} x$$

$$J = P^{-1} (Y Y^T)$$

$$\frac{\partial L_2}{\partial x} = -\frac{\alpha}{2} P^{-1} Y Y^T P^{-1} x$$

$$\frac{\partial L_2}{\partial P^{-1}} = -\frac{1}{2} (Y Y^T)$$

$$\frac{\partial L_2}{\partial x^T} = -\frac{\alpha}{2} x^T P^{-1} Y Y^T P^{-1}$$

$$\frac{\partial L_2}{\partial P} = -P^{-1} \frac{\partial J}{\partial P^{-1}} P^{-1}$$

$$= \frac{\partial L_2}{\partial x} + \left[\frac{\partial L_2}{\partial x^T} \right]^T$$

$$\frac{\partial L_2}{\partial R} = -P^{-1} \frac{\partial J}{\partial P^{-1}} P^{-1}$$

$$= \frac{\alpha}{2} P^{-1} Y Y^T P^{-1} x + \left[\frac{\alpha}{2} x^T P^{-1} Y Y^T P^{-1} \right]^T$$

$$\frac{\partial L_2}{\partial F} = -\alpha P^{-1} \frac{\partial J}{\partial P^{-1}} P^{-1}$$

$$= \frac{\alpha}{2} P^{-1} Y Y^T P^{-1} x + \frac{\alpha}{2} (P^{-1})^T Y Y^T (P^T)^{-1} x$$

Pin Symmetrie 10:

$$= \frac{\alpha}{2} P^{-1} Y Y^T P^{-1} x + \frac{\alpha}{2} P^{-1} Y Y^T P^{-1} x$$

$$\boxed{\frac{\partial L_2}{\partial x} = \alpha P^{-1} Y Y^T P^{-1} x}$$

e) By the law of total derivation, we have:

$$\begin{aligned} \frac{\partial L}{\partial x} &= \frac{\partial L_1}{\partial x} + \frac{\partial L_2}{\partial x} \\ &= -\alpha D (\alpha x x^T + \beta^{-1} I)^{-1} x + \alpha (\alpha x x^T + \beta^{-1} I)^{-1} Y Y^T (\alpha x x^T + \beta^{-1} I)^{-1} x \\ &= \alpha [-D I + (\alpha x x^T + \beta^{-1} I)^{-1} Y Y^T] (\alpha x x^T + \beta^{-1} I)^{-1} x \end{aligned}$$

3. a) 2-layer FC

$H \rightarrow$ hidden layers

$C \rightarrow \neq$ (neurons in the output layer)

$$\text{Swish}(k) = \frac{k}{1 + e^{-k}} = k \sigma(k), \text{ where } \sigma(k) \text{ is the sigmoid activation function.}$$

$$= k (1 + e^{-k})^{-1}$$

$$z_1 = w_1 x + b_1, \quad z_1 \in \mathbb{R}^H \text{ (affine function)}, \quad x \in \mathbb{R}^D$$

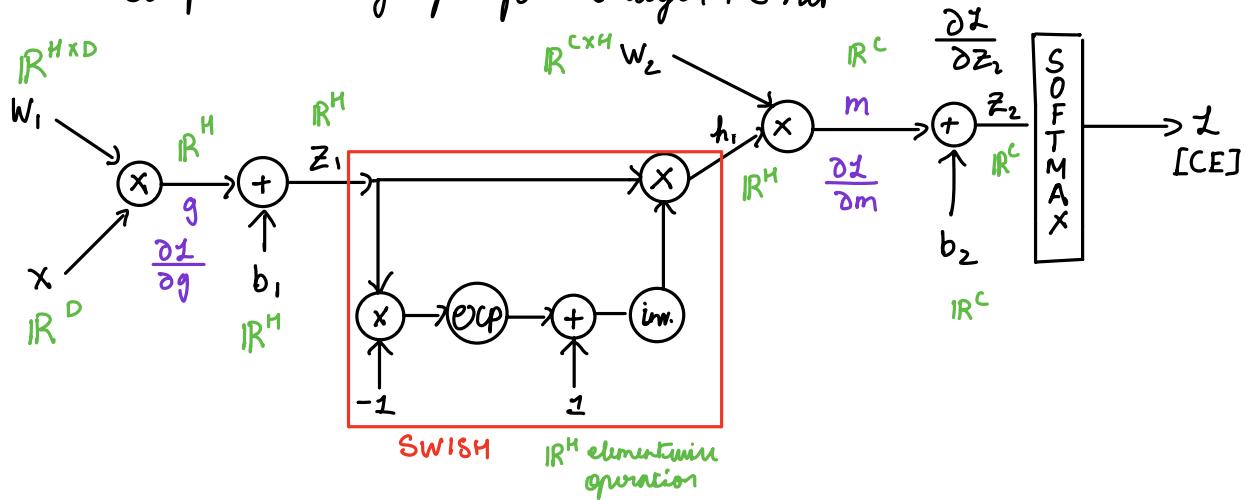
$$h_1 = \text{Swish}(z_1) = k \sigma(z_1) \quad h'_1 = z_1 + \sigma(z_1)(1 - \sigma(z_1))$$

$$z_2 = w_2 h_1 + b_2, \quad z_2 \in \mathbb{R}^C$$

$$\hat{y} = \text{Softmax}(z_2)$$

$$L = \text{Cross entropy} \Rightarrow L = - \sum_i y_i \log(\hat{y}_i)$$

Computational graph for 2-layer FC net:



$$b) \& c) \nabla_{w_2} L, \nabla_{b_2} L, \nabla_{w_1} L, \nabla_{b_1} L$$

$$\frac{\partial L}{\partial m} = \frac{\partial L}{\partial z_2} \quad \& \quad \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2} \quad \therefore z_2 = m + b_2$$

$$m = w_2 h_1, \quad \frac{\partial m}{\partial h_1} = w_2^T \mathbb{R}^{HxL} \quad \frac{\partial m}{\partial w_2} = h_1^T \mathbb{R}^{CxH \times C}$$

$$\frac{\partial L}{\partial h_1} = \frac{\partial L}{\partial m} \frac{\partial m}{\partial h_1} = w_2^T \frac{\partial L}{\partial m}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial m}{\partial w_2} \frac{\partial L}{\partial m} = \frac{\partial L}{\partial m} h_1^T$$

$$Z_1 = \text{Sigmoid}(h_1)$$

$$Z_1 = \begin{bmatrix} Z_{1,1} \\ Z_{1,2} \\ \vdots \\ Z_{1,M} \end{bmatrix} \quad h_1 = \text{Sigmoid}(Z_1) = \begin{bmatrix} \text{Sigmoid}(Z_{1,1}) \\ \text{Sigmoid}(Z_{1,2}) \\ \vdots \\ \text{Sigmoid}(Z_{1,M}) \end{bmatrix} = \begin{bmatrix} h_{1,1} \\ h_{1,2} \\ \vdots \\ h_{1,M} \end{bmatrix}$$

JACOBIAN

$$\frac{\partial h_1}{\partial Z_1} = \begin{bmatrix} \frac{\partial h_1}{\partial Z_{1,1}} \\ \frac{\partial h_1}{\partial Z_{1,2}} \\ \vdots \\ \frac{\partial h_1}{\partial Z_{1,M}} \end{bmatrix} = \begin{bmatrix} \frac{\partial h_{1,1}}{\partial Z_{1,1}} & \frac{\partial h_{1,2}}{\partial Z_{1,1}} & \cdots & \frac{\partial h_{1,M}}{\partial Z_{1,1}} \\ \frac{\partial h_{1,1}}{\partial Z_{1,2}} & \frac{\partial h_{1,2}}{\partial Z_{1,2}} & \cdots & \frac{\partial h_{1,M}}{\partial Z_{1,2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_{1,1}}{\partial Z_{1,M}} & \frac{\partial h_{1,2}}{\partial Z_{1,M}} & \cdots & \frac{\partial h_{1,M}}{\partial Z_{1,M}} \end{bmatrix}$$

Only diagonal elements remain non-zero

$$\frac{\partial h_{1,i}}{\partial Z_{1,j}} = \begin{cases} \frac{\partial}{\partial Z_{1,i}} \left[\frac{Z_{1,i}}{1 + \exp(-Z_{1,i})} \right], & \text{if } i=j \\ 0 & \text{in } i \neq j \end{cases}$$

by Quotient rule we have

$$\begin{aligned} \frac{\partial}{\partial k} &= \left[\frac{k}{1 + \exp(-k)} \right] = \frac{1 \times (1 + \exp(-k)) - [-\exp(-k)]k}{(1 + \exp(-k))^2} \\ &= \frac{1 + (k+1)e^{-k}}{(1 + e^{-k})} \end{aligned}$$

$$\frac{\partial h_1}{\partial Z_1} = \text{diag} \left(\frac{1 + (Z_{1,1} + 1)e^{-Z_{1,1}}}{(1 + e^{-Z_{1,1}})^2}, \frac{1 + (Z_{1,2} + 1)e^{-Z_{1,2}}}{(1 + e^{-Z_{1,2}})^2}, \dots \right)$$

$$\frac{\partial L}{\partial Z_1} = \frac{\partial h_1}{\partial Z_1} \frac{\partial L}{\partial h_1} = S_1 \odot W_2^T \frac{\partial L}{\partial Z_2}$$

where $S_1 =$

$$\begin{bmatrix} \frac{1 + (Z_{1,1} + 1)e^{-Z_{1,1}}}{(1 + e^{-Z_{1,1}})^2} \\ \frac{1 + (Z_{1,2} + 1)e^{-Z_{1,2}}}{(1 + e^{-Z_{1,2}})^2} \\ \vdots \\ \frac{1 + (Z_{1,M} + 1)e^{-Z_{1,M}}}{(1 + e^{-Z_{1,M}})^2} \end{bmatrix}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial z_1}{\partial w_1} \frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial z_1} x^T = \left[S_1 \odot w_2^T \frac{\partial L}{\partial z_2} \right] x^T$$

$$\frac{\partial L}{\partial x} = \frac{\partial x}{\partial w_1} \frac{\partial L}{\partial x} = w_1^T \frac{\partial L}{\partial x} = w_1^T \left[S_1 \odot w_2^T \frac{\partial L}{\partial z_2} \right]$$

$$\nabla_{w_2} L = \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2} h_1^T$$

$$\nabla_{b_2} L = \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2}$$

$$\nabla_{w_1} L = \frac{\partial L}{\partial w_1} = \frac{\partial z_1}{\partial w_1} \frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial z_1} x^T = \left[S_1 \odot w_2^T \frac{\partial L}{\partial z_2} \right] x^T$$

$$\nabla_{b_1} L = \frac{\partial L}{\partial z_1} = \frac{\partial h_1}{\partial z_1} \frac{\partial L}{\partial h_1} = S_1 \odot w_2^T \frac{\partial L}{\partial z_2}$$

```

import numpy as np
import matplotlib.pyplot as plt

class TwoLayerNet(object):
    """
    A two-layer fully-connected neural network. The net has an input dimension of
    D, a hidden layer dimension of H, and performs classification over C classes.
    We train the network with a softmax loss function and L2 regularization on the
    weight matrices. The network uses a ReLU nonlinearity after the first fully
    connected layer.

    In other words, the network has the following architecture:

    input - fully connected layer - ReLU - fully connected layer - softmax

    The outputs of the second fully-connected layer are the scores for each class.
    """

    def __init__(self, input_size, hidden_size, output_size, std=1e-4):
        """
        Initialize the model. Weights are initialized to small random values and
        biases are initialized to zero. Weights and biases are stored in the
        variable self.params, which is a dictionary with the following keys:

        W1: First layer weights; has shape (H, D)
        b1: First layer biases; has shape (H,)
        W2: Second layer weights; has shape (C, H)
        b2: Second layer biases; has shape (C,)

        Inputs:
        - input_size: The dimension D of the input data.
        - hidden_size: The number of neurons H in the hidden layer.
        - output_size: The number of classes C.
        """
        self.params = {}
        self.params['W1'] = std * np.random.randn(hidden_size, input_size)
        self.params['b1'] = np.zeros(hidden_size)
        self.params['W2'] = std * np.random.randn(output_size, hidden_size)
        self.params['b2'] = np.zeros(output_size)

    def loss(self, X, y=None, reg=0.0):
        """
        Compute the loss and gradients for a two layer fully connected neural
        network.

        Inputs:
        - X: Input data of shape (N, D). Each X[i] is a training sample.
        - y: Vector of training labels. y[i] is the label for X[i], and each y[i] is
          an integer in the range 0 <= y[i] < C. This parameter is optional; if it
          is not passed then we only return scores, and if it is passed then we
          instead return the loss and gradients.
        - reg: Regularization strength.

        Returns:
        If y is None, return a matrix scores of shape (N, C) where scores[i, c] is
        the score for class c on input X[i].
        """
        if y is not None:
            # Compute the forward pass
            scores = None

            # YOUR CODE HERE:
            # =====#
            # Calculate the output scores of the neural network. The result
            # should be (N, C). As stated in the description for this class,
            # there should not be a ReLU layer after the second FC layer.
            # The output of the second FC layer is the output scores. Do not
            # use a for loop in your implementation.
            # =====#
            Hin = X@W1.T + b1

```

```

Hout = Hin.copy()
Hout[Hout < 0] = 0
Z = Hout@W2.T + b2
scores = Z

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

# If the targets are not given then jump out, we're done
if y is None:
    return scores

# Compute the loss
loss = None

# ===== #
# YOUR CODE HERE:
#   Calculate the loss of the neural network. This includes the
#   softmax loss and the L2 regularization for W1 and W2. Store the
#   total loss in teh variable loss. Multiply the regularization
#   loss by 0.5 (in addition to the factor reg).
# ===== #
exp_x = np.exp(scores - np.max(scores, axis=1).reshape(-1, 1))
sum_x = np.sum(exp_x, axis=1).reshape(-1, 1)
prob = exp_x / sum_x
loss = 0.5 * reg * (np.power(np.linalg.norm(W1, "fro"), 2) + np.power(np.linalg.norm(W2, "fro"), 2)) # L2 regularization
loss += -np.sum(np.log(prob[np.arange(N), y])) / N # softmax loss
# ===== #
# END YOUR CODE HERE
# ===== #

grads = {}

# ===== #
# YOUR CODE HERE:
#   Implement the backward pass. Compute the derivatives of the
#   weights and the biases. Store the results in the grads
#   dictionary. e.g., grads['W1'] should store the gradient for
#   W1, and be of the same size as W1.
# ===== #
C, H = W2.shape

indicator = np.zeros((N, C))
indicator[np.arange(N), y] = 1

softmax = (prob.T - indicator.T) / N

grads['W2'] = softmax@Hout + reg * W2
grads['b2'] = np.sum(softmax, axis=1)

relu_indicator = Hout/Hin
H_grad = np.multiply(relu_indicator, (W2.T@softmax).T).T

grads['W1'] = H_grad@X + reg * W1
grads['b1'] = np.sum(H_grad, axis=1)

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

return loss, grads

def train(self, X, y, X_val, y_val,
          learning_rate=1e-3, learning_rate_decay=0.95,
          reg=1e-5, num_iters=100,
          batch_size=200, verbose=False):
    """
    Train this neural network using stochastic gradient descent.

    Inputs:
    - X: A numpy array of shape (N, D) giving training data.
    - y: A numpy array f shape (N,) giving training labels; y[i] = c means that
      X[i] has label c, where 0 <= c < C.
    - X_val: A numpy array of shape (N_val, D) giving validation data.
    - y_val: A numpy array of shape (N_val,) giving validation labels.
    - learning_rate: Scalar giving learning rate for optimization.
    - learning_rate_decay: Scalar giving factor used to decay the learning rate
      after each epoch.
    - reg: Scalar giving regularization strength.
    """

```

```

- num_iters: Number of steps to take when optimizing.
- batch_size: Number of training examples to use per step.
- verbose: boolean; if true print progress during optimization.
"""
num_train = X.shape[0]
iterations_per_epoch = max(num_train / batch_size, 1)

# Use SGD to optimize the parameters in self.model
loss_history = []
train_acc_history = []
val_acc_history = []

for it in np.arange(num_iters):

    X_batch = None
    y_batch = None

    # ===== #
    # YOUR CODE HERE:
    #   Create a minibatch by sampling batch_size samples randomly.
    # ===== #
    batch_indices = np.random.choice(num_train, batch_size, replace=True)
    X_batch = X[batch_indices]
    y_batch = y[batch_indices]

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

    # Compute loss and gradients using the current minibatch
    loss, grads = self.loss(X_batch, y=y_batch, reg=reg)
    loss_history.append(loss)

    # ===== #
    # YOUR CODE HERE:
    #   Perform a gradient descent step using the minibatch to update
    #   all parameters (i.e., W1, W2, b1, and b2).
    # ===== #
    self.params['W1'] -= learning_rate * grads['W1']
    self.params['b1'] -= learning_rate * grads['b1']
    self.params['W2'] -= learning_rate * grads['W2']
    self.params['b2'] -= learning_rate * grads['b2']

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

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

    # Every epoch, check train and val accuracy and decay learning rate.
    if it % iterations_per_epoch == 0:
        # Check accuracy
        train_acc = (self.predict(X_batch) == y_batch).mean()
        val_acc = (self.predict(X_val) == y_val).mean()
        train_acc_history.append(train_acc)
        val_acc_history.append(val_acc)

        # Decay learning rate
        learning_rate *= learning_rate_decay

return {
    'loss_history': loss_history,
    'train_acc_history': train_acc_history,
    'val_acc_history': val_acc_history,
}
}

def predict(self, X):
    """
    Use the trained weights of this two-layer network to predict labels for
    data points. For each data point we predict scores for each of the C
    classes, and assign each data point to the class with the highest score.

    Inputs:
    - X: A numpy array of shape (N, D) giving N D-dimensional data points to
    classify.

    Returns:
    - y_pred: A numpy array of shape (N,) giving predicted labels for each of
    the elements of X. For all i, y_pred[i] = c means that X[i] is predicted
    to have class c, where 0 <= c < C.
    """
y_pred = None

```

```
# ===== #
# YOUR CODE HERE:
#   Predict the class given the input data.
# ===== #

Hin = X @ self.params['W1'].T + self.params['b1']
Hout = Hin.copy()
Hout[Hout < 0] = 0
Z = Hout @ self.params['W2'].T + self.params['b2']
y_pred = np.argmax(Z, axis=1)

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

return y_pred
```

This is the 2-layer neural network notebook for ECE C147/C247 Homework #3

Please follow the notebook linearly to implement a two layer neural network.

Please print out the notebook entirely when completed.

The goal of this notebook is to give you experience with training a two layer neural network.

In [24]:

```
1 import random
2 import numpy as np
3 from utils.data_utils import load_CIFAR10
4 import matplotlib.pyplot as plt
5
6 %matplotlib inline
7 %load_ext autoreload
8 %autoreload 2
9
10 def rel_error(x, y):
11     """ returns relative error """
12     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
```

Toy example

Before loading CIFAR-10, there will be a toy example to test your implementation of the forward and backward pass. Make sure to read the description of TwoLayerNet class in neural_net.py file , understand the architecture and initializations

In [25]:

```
1 from nnndl.neural_net import TwoLayerNet
```

In [26]:

```
1 # Create a small net and some toy data to check your implementations.
2 # Note that we set the random seed for repeatable experiments.
3
4 input_size = 4
5 hidden_size = 10
6 num_classes = 3
7 num_inputs = 5
8
9 def init_toy_model():
10     np.random.seed(0)
11     return TwoLayerNet(input_size, hidden_size, num_classes, std=1e-1)
12
13 def init_toy_data():
14     np.random.seed(1)
15     X = 10 * np.random.randn(num_inputs, input_size)
16     y = np.array([0, 1, 2, 2, 1])
17     return X, y
18
19 net = init_toy_model()
20 X, y = init_toy_data()
```

Compute forward pass scores

```
In [27]: 1 ## Implement the forward pass of the neural network.  
2 ## See the loss() method in TwoLayerNet class for the same  
3  
4 # Note, there is a statement if y is None: return scores, which is why  
5 # the following call will calculate the scores.  
6 scores = net.loss(X)  
7 print('Your scores:')8 print(scores)  
9 print()  
10 print('correct scores:')11 correct_scores = np.asarray([  
12     [-1.07260209, 0.05083871, -0.87253915],  
13     [-2.02778743, -0.10832494, -1.52641362],  
14     [-0.74225908, 0.15259725, -0.39578548],  
15     [-0.38172726, 0.10835902, -0.17328274],  
16     [-0.64417314, -0.18886813, -0.41106892]])  
17 print(correct_scores)  
18 print()  
19  
20 # The difference should be very small. We get < 1e-7  
21 print('Difference between your scores and correct scores:')22 print(np.sum(np.abs(scores - correct_scores)))
```

Your scores:
[[-1.07260209 0.05083871 -0.87253915]
 [-2.02778743 -0.10832494 -1.52641362]
 [-0.74225908 0.15259725 -0.39578548]
 [-0.38172726 0.10835902 -0.17328274]
 [-0.64417314 -0.18886813 -0.41106892]]

correct scores:
[[-1.07260209 0.05083871 -0.87253915]
 [-2.02778743 -0.10832494 -1.52641362]
 [-0.74225908 0.15259725 -0.39578548]
 [-0.38172726 0.10835902 -0.17328274]
 [-0.64417314 -0.18886813 -0.41106892]]

Difference between your scores and correct scores:
3.381231233889892e-08

Forward pass loss

```
In [28]: 1 loss, _ = net.loss(X, y, reg=0.05)  
2 correct_loss = 1.071696123862817  
3  
4 # should be very small, we get < 1e-12  
5 print("Loss:", loss)  
6 print('Difference between your loss and correct loss:')7 print(np.sum(np.abs(loss - correct_loss)))
```

Loss: 1.071696123862817
Difference between your loss and correct loss:
0.0

Backward pass

Implements the backwards pass of the neural network. Check your gradients with the gradient check utilities provided.

```
In [29]: 1 from utils.gradient_check import eval_numerical_gradient
2
3 # Use numeric gradient checking to check your implementation of the backward pass.
4 # If your implementation is correct, the difference between the numeric and
5 # analytic gradients should be less than 1e-8 for each of W1, W2, b1, and b2.
6
7 loss, grads = net.loss(X, y, reg=0.05)
8
9 # these should all be less than 1e-8 or so
10 for param_name in grads:
11     f = lambda W: net.loss(X, y, reg=0.05)[0]
12     param_grad_num = eval_numerical_gradient(f, net.params[param_name], verbose=False)
13     print('{} max relative error: {}'.format(param_name, rel_error(param_grad_num, grads[param_name])))

W2 max relative error: 2.9632227682005116e-10
b2 max relative error: 1.248270530283678e-09
W1 max relative error: 1.2832823337649917e-09
b1 max relative error: 3.172680092703762e-09
```

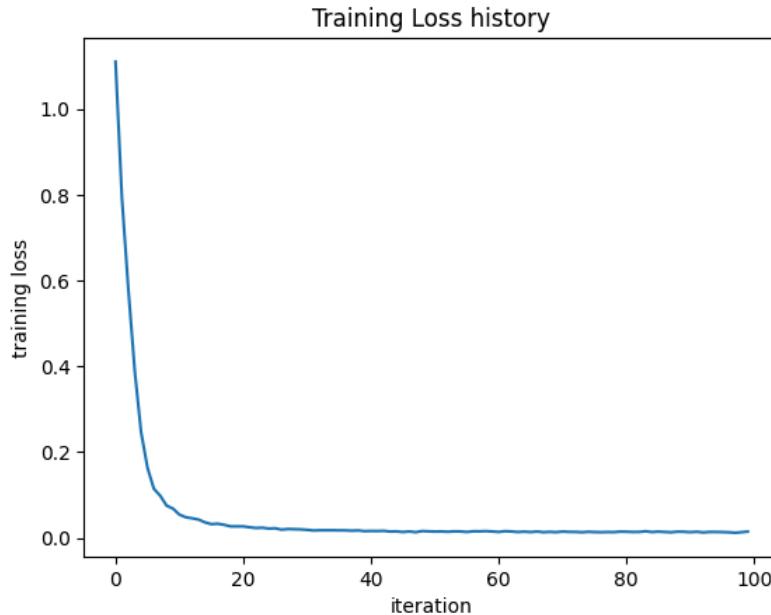
Training the network

Implement `neural_net.train()` to train the network via stochastic gradient descent, much like the softmax and SVM.

```
In [30]: 1 net = init_toy_model()
2 stats = net.train(X, y, X, y,
3                   learning_rate=1e-1, reg=5e-6,
4                   num_iters=100, verbose=False)
5
6 print('Final training loss: ', stats['loss_history'][-1])
7
8 # plot the loss history
9 plt.plot(stats['loss_history'])
10 plt.xlabel('iteration')
11 plt.ylabel('training loss')
12 plt.title('Training Loss history')
13 plt.show()

Final training loss:  0.014497864587765886

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "orientation" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "dpi" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "facecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "edgecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "bbox_inches_restore" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
```



Classify CIFAR-10

Do classification on the CIFAR-10 dataset.

```
In [32]: 1 from utils.data_utils import load_CIFAR10
2
3 def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000):
4     """
5         Load the CIFAR-10 dataset from disk and perform preprocessing to prepare
6         it for the two-layer neural net classifier.
7     """
8     # Load the raw CIFAR-10 data
9     cifar10_dir = '/Users/sujitsilas/Desktop/UCLA/Winter 2025/EE ENGR 247/Homeworks/HW2/student_copy/cif
10    X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)
11
12    # Subsample the data
13    mask = list(range(num_training, num_training + num_validation))
14    X_val = X_train[mask]
15    y_val = y_train[mask]
16    mask = list(range(num_training))
17    X_train = X_train[mask]
18    y_train = y_train[mask]
19    mask = list(range(num_test))
20    X_test = X_test[mask]
21    y_test = y_test[mask]
22
23    # Normalize the data: subtract the mean image
24    mean_image = np.mean(X_train, axis=0)
25    X_train -= mean_image
26    X_val -= mean_image
27    X_test -= mean_image
28
29    # Reshape data to rows
30    X_train = X_train.reshape(num_training, -1)
31    X_val = X_val.reshape(num_validation, -1)
32    X_test = X_test.reshape(num_test, -1)
33
34    return X_train, y_train, X_val, y_val, X_test, y_test
35
36
37 # Invoke the above function to get our data.
38 X_train, y_train, X_val, y_val, X_test, y_test = get_CIFAR10_data()
39 print('Train data shape: ', X_train.shape)
40 print('Train labels shape: ', y_train.shape)
41 print('Validation data shape: ', X_val.shape)
42 print('Validation labels shape: ', y_val.shape)
43 print('Test data shape: ', X_test.shape)
44 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,)
```

Running SGD

If your implementation is correct, you should see a validation accuracy of around 28-29%.

```
In [ ]: 1 input_size = 32 * 32 * 3
2 hidden_size = 50
3 num_classes = 10
4 net = TwoLayerNet(input_size, hidden_size, num_classes)
5
6 # Train the network
7 stats = net.train(X_train, y_train, X_val, y_val,
8                     num_iters=1000, batch_size=200,
9                     learning_rate=1e-4, learning_rate_decay=0.95,
10                    reg=0.25, verbose=True)
11
12 # Predict on the validation set
13 val_acc = (net.predict(X_val) == y_val).mean()
14 print('Validation accuracy: ', val_acc)
15
16 # Save this net as the variable subopt_net for later comparison.
17 subopt_net = net
```

```
iteration 0 / 1000: loss 2.302757518613176
iteration 100 / 1000: loss 2.302120159207236
iteration 200 / 1000: loss 2.2956136007408703
iteration 300 / 1000: loss 2.2518259043164135
iteration 400 / 1000: loss 2.188995235046776
iteration 500 / 1000: loss 2.1162527791897747
iteration 600 / 1000: loss 2.064670827698217
iteration 700 / 1000: loss 1.990168862308394
iteration 800 / 1000: loss 2.002827640124685
iteration 900 / 1000: loss 1.94651768178565
Validation accuracy: 0.283
```

Questions:

The training accuracy isn't great.

(1) What are some of the reasons why this is the case? Take the following cell to do some analyses and then report your answers in the cell following the one below.

(2) How should you fix the problems you identified in (1)?

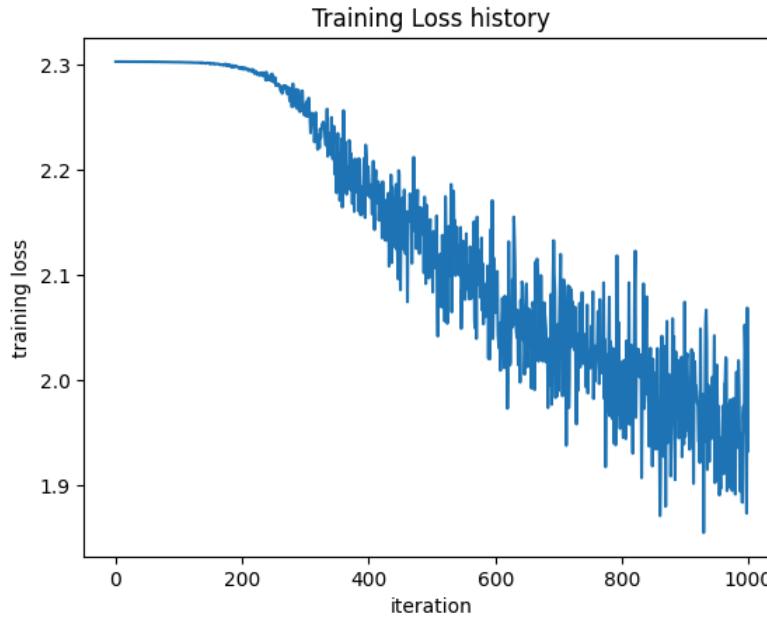
```
In [34]: 1 stats['train_acc_history']
```

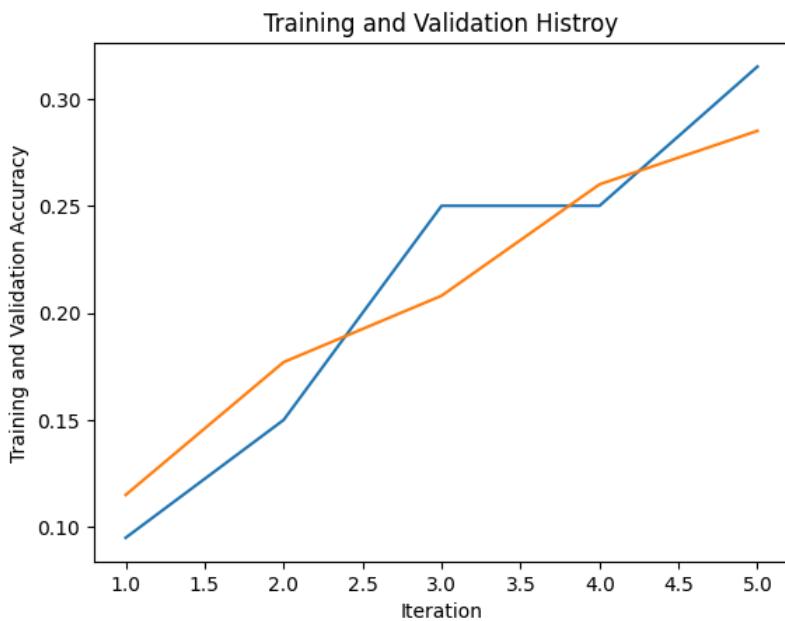
```
Out[34]: [0.095, 0.15, 0.25, 0.25, 0.315]
```

In []:

```
1 # ===== #
2 # YOUR CODE HERE:
3 # Do some debugging to gain some insight into why the optimization
4 # isn't great.
5 # ===== #
6
7 # Plot the loss function and train / validation accuracies
8
9 plt.plot(stats['loss_history'], label='loss')
10 plt.xlabel('iteration')
11 plt.ylabel('training loss')
12 plt.title('Training Loss history')
13 plt.show()
14
15 plt.plot([1,2,3,4,5], stats['train_acc_history'], label='Train Accuracy')
16 plt.plot([1,2,3,4,5], stats['val_acc_history'], label='Validation Accuracy')
17 plt.xlabel('Iteration')
18 plt.ylabel('Training and Validation Accuracy')
19 plt.title('Training and Validation Histroy')
20 plt.show()
21
22 # ===== #
23 # END YOUR CODE HERE
24 # ===== #
```

```
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "orientation" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "dpi" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "facecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "edgecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "bbox_inches_restore" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
```





Answers:

- (1) The training performance isn't great because the loss is highly volatile (zigzagging) instead of smoothly converging. This instability can be caused by improper weight initialization with large values, insufficient training time, overly strong regularization, or learning rate issues—either too high, causing overshooting, or too low, leading to slow convergence.
- (2) The problem can be fixed by extending training time, improving weight initialization with methods such as Xavier initialization to prevent gradients from exploding. We can also experiment with different regularization strengths and learning rates.

Optimize the neural network

Use the following part of the Jupyter notebook to optimize your hyperparameters on the validation set. Store your nets as `best_net`.

In [48]:

```
1 # ===== #
2 # YOUR CODE HERE:
3 #   Optimize over your hyperparameters to arrive at the best neural
4 #   network. You should be able to get over 50% validation accuracy.
5 #   For this part of the notebook, we will give credit based on the
6 #   accuracy you get. Your score on this question will be multiplied by:
7 #       min(floor((X - 28%) / %22, 1)
8 #   where if you get 50% or higher validation accuracy, you get full
9 #   points.
10 #
11 #
12 #   Note, you need to use the same network structure (keep hidden_size = 50) !
13 # ===== #
14 input_size = 32 * 32 * 3
15 hidden_size = 50
16 num_classes = 10
17 net = TwoLayerNet(input_size, hidden_size, num_classes)
18
19 # Hyperparameter tuning with Xavier initialization
20 best_val_acc = 0
21 best_net = None
22
23 # Grid search over hyperparameters
24 learning_rates=[0.001, 0.005, 0.0001, 0.00005]
25 regularization_strengths = [0.1, 0.25]
26 learning_rate_decays = [0.55, 0.60, 0.80]
27
28 for lr in learning_rates:
29     for reg in regularization_strengths:
30         for lr_decay in learning_rate_decays:
31             print(f"Training with learning_rate={lr}, reg={reg}, lr_decay={lr_decay}")
32
33     # Initialize the network with Xavier initialization
34     net = TwoLayerNet(input_size, hidden_size, num_classes)
35
36     # Train the network
37     stats = net.train(X_train, y_train, X_val, y_val,
38                       num_iters=1000,
39                       batch_size=200,
40                       learning_rate=lr,
41                       learning_rate_decay=lr_decay,
42                       reg=reg,
43                       verbose=False)
44
45     # Predict on the validation set
46     val_acc = (net.predict(X_val) == y_val).mean()
47     print(f"Validation accuracy: {val_acc}")
48
49     # Save the best model
50     if val_acc > best_val_acc:
51         best_val_acc = val_acc
52         best_net = net
53
54 print('Best validation accuracy: ', best_val_acc)
55
56 # ===== #
57 # END YOUR CODE HERE
58 # ===== #
```

```
Training with learning_rate=0.001, reg=0.1, lr_decay=0.55
Validation accuracy: 0.422
Training with learning_rate=0.001, reg=0.1, lr_decay=0.6
Validation accuracy: 0.424
Training with learning_rate=0.001, reg=0.1, lr_decay=0.8
Validation accuracy: 0.467
Training with learning_rate=0.001, reg=0.25, lr_decay=0.55
Validation accuracy: 0.413
Training with learning_rate=0.001, reg=0.25, lr_decay=0.6
Validation accuracy: 0.435
Training with learning_rate=0.001, reg=0.25, lr_decay=0.8
Validation accuracy: 0.472
Training with learning_rate=0.005, reg=0.1, lr_decay=0.55
Validation accuracy: 0.489
Training with learning_rate=0.005, reg=0.1, lr_decay=0.6
Validation accuracy: 0.477
Training with learning_rate=0.005, reg=0.1, lr_decay=0.8
```

```
/Users/sujitsilas/Desktop/UCLA/Winter 2025/EE ENGR 247/Homeworks/HW3/code_student_version/nndl/neural_net.p
y:113: RuntimeWarning: divide by zero encountered in log
    loss += -np.sum(np.log(prob[np.arange(N), y])) / N # softmax loss

Validation accuracy: 0.284
Training with learning_rate=0.005, reg=0.25, lr_decay=0.55
Validation accuracy: 0.496
Training with learning_rate=0.005, reg=0.25, lr_decay=0.6
Validation accuracy: 0.513
Training with learning_rate=0.005, reg=0.25, lr_decay=0.8
Validation accuracy: 0.227
Training with learning_rate=0.0001, reg=0.1, lr_decay=0.55
Validation accuracy: 0.182
Training with learning_rate=0.0001, reg=0.1, lr_decay=0.6
Validation accuracy: 0.183
Training with learning_rate=0.0001, reg=0.1, lr_decay=0.8
Validation accuracy: 0.246
Training with learning_rate=0.0001, reg=0.25, lr_decay=0.55
Validation accuracy: 0.161
Training with learning_rate=0.0001, reg=0.25, lr_decay=0.6
Validation accuracy: 0.194
Training with learning_rate=0.0001, reg=0.25, lr_decay=0.8
Validation accuracy: 0.244
Training with learning_rate=0.0005, reg=0.1, lr_decay=0.55
Validation accuracy: 0.337
Training with learning_rate=0.0005, reg=0.1, lr_decay=0.6
Validation accuracy: 0.365
Training with learning_rate=0.0005, reg=0.1, lr_decay=0.8
Validation accuracy: 0.427
Training with learning_rate=0.0005, reg=0.25, lr_decay=0.55
Validation accuracy: 0.334
Training with learning_rate=0.0005, reg=0.25, lr_decay=0.6
Validation accuracy: 0.363
Training with learning_rate=0.0005, reg=0.25, lr_decay=0.8
Validation accuracy: 0.421
Best validation accuracy: 0.513
```

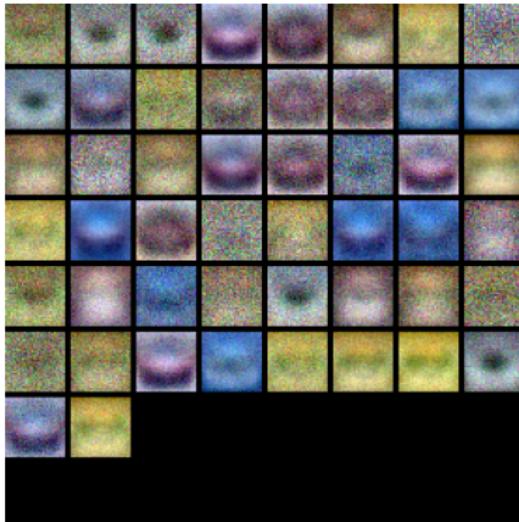
Best Validation Accuracy: 0.513

- Training with learning_rate=0.005, reg=0.25, lr_decay=0.6

In [49]:

```
1 from utils.vis_utils import visualize_grid
2
3 # Visualize the weights of the network
4
5 def show_net_weights(net):
6     W1 = net.params['W1']
7     W1 = W1.T.reshape(32, 32, 3, -1).transpose(3, 0, 1, 2)
8     plt.imshow(visualize_grid(W1, padding=3).astype('uint8'))
9     plt.gca().axis('off')
10    plt.show()
11
12 show_net_weights(subopt_net)
13 show_net_weights(best_net)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "orientation" which is no longer supported as of 3.3 and will become an error two minor releases later
fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "dpi" which is no longer supported as of 3.3 and will become an error two minor releases later
fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "facecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "edgecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "bbox_inches_restore" which is no longer supported as of 3.3 and will become an error two minor releases later
fig.canvas.print_figure(bytes_io, **kw)



Question:

(1) What differences do you see in the weights between the suboptimal net and the best net you arrived at?

Answer:

(1) The best net has more evenly distributed and well-scaled weights compared to the suboptimal net, which may have erratic or extreme values due to small learning rates. The features are more distinctly visible on the best net leading to higher validation accuracy. The weights of the best net look more like a template that the model can use to assign images to certain classes. The suboptimal net's weights might exhibit a wider spread, leading to vanishing or exploding gradients, while the best net maintains controlled weight distributions for stable optimization.

Evaluate on test set

```
In [50]: 1 test_acc = (best_net.predict(X_test) == y_test).mean()
          2 print('Test accuracy: ', test_acc)
```

Test accuracy: 0.49

```

import numpy as np
import pdb


def affine_forward(x, w, b):
    """
    Computes the forward pass for an affine (fully-connected) layer.

    The input x has shape (N, d_1, ..., d_k) and contains a minibatch of N
    examples, where each example x[i] has shape (d_1, ..., d_k). We will
    reshape each input into a vector of dimension D = d_1 * ... * d_k, and
    then transform it to an output vector of dimension M.

    Inputs:
    - x: A numpy array containing input data, of shape (N, d_1, ..., d_k)
    - w: A numpy array of weights, of shape (D, M)
    - b: A numpy array of biases, of shape (M,)

    Returns a tuple of:
    - out: output, of shape (N, M)
    - cache: (x, w, b)
    """
    # ===== #
    # YOUR CODE HERE:
    #   Calculate the output of the forward pass. Notice the dimensions
    #   of w are D x M, which is the transpose of what we did in earlier
    #   assignments.
    # ===== #

    N = x.shape[0]
    D = np.prod(x.shape[1:])  # Calculate flattened input dimension
    x_reshaped = x.reshape(N, D)  # Reshape x for matrix multiplication
    out = x_reshaped @ w + b  # Vectorized affine transformation (@ is shorthand for np.matmul)
    cache = (x, w, b)
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    return out, cache


def affine_backward(dout, cache):
    """
    Computes the backward pass for an affine layer.

    Inputs:
    - dout: Upstream derivative, of shape (N, M)
    - cache: Tuple of:
        - x: Input data, of shape (N, d_1, ..., d_k)
        - w: Weights, of shape (D, M)

    Returns a tuple of:
    - dx: Gradient with respect to x, of shape (N, d1, ..., d_k)
    - dw: Gradient with respect to w, of shape (D, M)
    - db: Gradient with respect to b, of shape (M,)
    """
    x, w, b = cache
    dx, dw, db = None, None, None

    # ===== #
    # YOUR CODE HERE:
    #   Calculate the gradients for the backward pass.
    # ===== #

    # dout is N x M
    # dx should be N x d1 x ... x dk; it relates to dout through multiplication with w, which is D x M

```

```

# dw should be D x M; it relates to dout through multiplication with x, which is N x D after reshaping
# db should be M; it is just the sum over dout examples

x, w, b = cache
dx = dout @ w.T
dx = dx.reshape(x.shape) # Reshape to match input shape

dw = x.reshape(x.shape[0], -1).T @ dout # Compute weight gradient
db = np.sum(dout, axis=0) # Compute bias gradient

return dx, dw, db

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

return dx, dw, db

def relu_forward(x):
    """
    Computes the forward pass for a layer of rectified linear units (ReLUs).

    Input:
    - x: Inputs, of any shape

    Returns a tuple of:
    - out: Output, of the same shape as x
    - cache: x
    """
    # ===== #
    # YOUR CODE HERE:
    # Implement the ReLU forward pass.
    # ===== #

    out = np.maximum(0, x)

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

    cache = x
    return out, cache

def relu_backward(dout, cache):
    """
    Computes the backward pass for a layer of rectified linear units (ReLUs).

    Input:
    - dout: Upstream derivatives, of any shape
    - cache: Input x, of same shape as dout

    Returns:
    - dx: Gradient with respect to x
    """
    x = cache

    # ===== #
    # YOUR CODE HERE:
    # Implement the ReLU backward pass
    # ===== #

    # ReLU directs linearly to those > 0
    dx = dout * (x > 0)

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

```

```
return dx

def softmax_loss(x, y):
    """
    Computes the loss and gradient for softmax classification.

    Inputs:
    - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
          for the ith input.
    - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
          0 <= y[i] < C

    Returns a tuple of:
    - loss: Scalar giving the loss
    - dx: Gradient of the loss with respect to x
    """
    probs = np.exp(x - np.max(x, axis=1, keepdims=True))
    probs /= np.sum(probs, axis=1, keepdims=True)
    N = x.shape[0]
    loss = -np.sum(np.log(probs[np.arange(N), y])) / N
    dx = probs.copy()
    dx[np.arange(N), y] -= 1
    dx /= N
    return loss, dx
```

```

import numpy as np

from .layers import *
from .layer_utils import *

class TwoLayerNet(object):
    """
    A two-layer fully-connected neural network with ReLU nonlinearity and
    softmax loss that uses a modular layer design. We assume an input dimension
    of D, a hidden dimension of H, and perform classification over C classes.

    The architecture should be affine - relu - affine - softmax.

    Note that this class does not implement gradient descent; instead, it
    will interact with a separate Solver object that is responsible for running
    optimization.

    The learnable parameters of the model are stored in the dictionary
    self.params that maps parameter names to numpy arrays.

    """
    def __init__(self, input_dim=3*32*32, hidden_dims=100, num_classes=10,
                 dropout=0, weight_scale=1e-3, reg=0.0):
        """
        Initialize a new network.

        Inputs:
        - input_dim: An integer giving the size of the input
        - hidden_dims: An integer giving the size of the hidden layer
        - num_classes: An integer giving the number of classes to classify
        - dropout: Scalar between 0 and 1 giving dropout strength.
        - weight_scale: Scalar giving the standard deviation for random
            initialization of the weights.
        - reg: Scalar giving L2 regularization strength.
        """
        self.params = {}
        self.reg = reg

        # ===== #
        # YOUR CODE HERE:
        #   Initialize W1, W2, b1, and b2. Store these as self.params['W1'],
        #   self.params['W2'], self.params['b1'] and self.params['b2']. The
        #   biases are initialized to zero and the weights are initialized
        #   so that each parameter has mean 0 and standard deviation weight_scale.
        #   The dimensions of W1 should be (input_dim, hidden_dim) and the
        #   dimensions of W2 should be (hidden_dims, num_classes)
        # ===== #
        self.params["W1"] = weight_scale * np.random.randn(input_dim, hidden_dims)
        self.params["b1"] = np.zeros(hidden_dims)
        self.params["W2"] = weight_scale * np.random.randn(hidden_dims, num_classes)
        self.params["b2"] = np.zeros(num_classes)

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

    def loss(self, X, y=None):
        """
        Compute loss and gradient for a minibatch of data.

        Inputs:
        - X: Array of input data of shape (N, d_1, ..., d_k)
        - y: Array of labels, of shape (N,). y[i] gives the label for X[i].
        """
        Compute loss and gradient for a minibatch of data.

        Returns:
        If y is None, then run a test-time forward pass of the model and return:
        - scores: Array of shape (N, C) giving classification scores, where
            scores[i, c] is the classification score for X[i] and class c.

        If y is not None, then run a training-time forward and backward pass and
        return a tuple of:
        - loss: Scalar value giving the loss
        - grads: Dictionary with the same keys as self.params, mapping parameter
            names to gradients of the loss with respect to those parameters.
        """
        scores = None

        # ===== #

```

```

# YOUR CODE HERE:
#   Implement the forward pass of the two-layer neural network. Store
#   the class scores as the variable 'scores'. Be sure to use the layers
#   you prior implemented.
# ===== #

# Forward pass: affine - relu - affine - softmax
out1, cache1 = affine_forward(X, self.params['W1'], self.params['b1'])
relu_out, relu_cache = relu_forward(out1)
scores, cache2 = affine_forward(relu_out, self.params['W2'], self.params['b2'])

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

# If y is None then we are in test mode so just return scores
if y is None:
    return scores

loss, grads = 0, {}
# ===== #
# YOUR CODE HERE:
#   Implement the backward pass of the two-layer neural net. Store
#   the loss as the variable 'loss' and store the gradients in the
#   'grads' dictionary. For the grads dictionary, grads['W1'] holds
#   the gradient for W1, grads['b1'] holds the gradient for b1, etc.
#   i.e., grads[k] holds the gradient for self.params[k].
#
#   Add L2 regularization, where there is an added cost 0.5*self.reg*W^2
#   for each W. Be sure to include the 0.5 multiplying factor to
#   match our implementation.
#
#   And be sure to use the layers you prior implemented.
# ===== #

loss, dscores = softmax_loss(scores, y)
loss += 0.5 * self.reg * (np.sum(self.params['W1'] ** 2) + np.sum(self.params['W2'] ** 2))

# Backward pass
dx2, dW2, db2 = affine_backward(dscores, cache2)
dW2 += self.reg * self.params['W2']

drelu = relu_backward(dx2, relu_cache)
dx1, dW1, db1 = affine_backward(drelu, cache1)
dW1 += self.reg * self.params['W1']

grads['W1'], grads['b1'] = dW1, db1
grads['W2'], grads['b2'] = dW2, db2

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

return loss, grads

class FullyConnectedNet(object):
    """
    A fully-connected neural network with an arbitrary number of hidden layers,
    ReLU nonlinearities, and a softmax loss function. This will also implement
    dropout and batch normalization as options. For a network with L layers,
    the architecture will be

    {affine - [batch norm] - relu - [dropout]} x (L - 1) - affine - softmax

    where batch normalization and dropout are optional, and the {...} block is
    repeated L - 1 times.

    Similar to the TwoLayerNet above, learnable parameters are stored in the
    self.params dictionary and will be learned using the Solver class.
    """
    def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
                 dropout=0, use_batchnorm=False, reg=0.0,
                 weight_scale=1e-2, dtype=np.float32, seed=None):
        """
        Initialize a new FullyConnectedNet.
        """

```

```

Inputs:
- hidden_dims: A list of integers giving the size of each hidden layer.
- input_dim: An integer giving the size of the input.
- num_classes: An integer giving the number of classes to classify.
- dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0 then
the network should not use dropout at all.
- use_batchnorm: Whether or not the network should use batch normalization.
- reg: Scalar giving L2 regularization strength.
- weight_scale: Scalar giving the standard deviation for random
initialization of the weights.
- dtype: A numpy datatype object; all computations will be performed using
this datatype. float32 is faster but less accurate, so you should use
float64 for numeric gradient checking.
- seed: If not None, then pass this random seed to the dropout layers. This
will make the dropout layers deterministic so we can gradient check the
model.

"""
self.use_batchnorm = use_batchnorm
self.use_dropout = dropout > 0
self.reg = reg
self.num_layers = 1 + len(hidden_dims)
self.dtype = dtype
self.params = {}

# ===== #
# YOUR CODE HERE:
#   Initialize all parameters of the network in the self.params dictionary.
#   The weights and biases of layer 1 are W1 and b1; and in general the
#   weights and biases of layer i are Wi and bi. The
#   biases are initialized to zero and the weights are initialized
#   so that each parameter has mean 0 and standard deviation weight_scale.
# ===== #
self.params = {}
input_dim_current = input_dim

for i, hidden_dim in enumerate(hidden_dims):
    self.params[f'W{i+1}'] = weight_scale * np.random.randn(input_dim_current, hidden_dim)
    self.params[f'b{i+1}'] = np.zeros(hidden_dim)
    input_dim_current = hidden_dim

    if self.use_batchnorm:
        self.params[f'gamma{i+1}'] = np.ones(hidden_dim)
        self.params[f'beta{i+1}'] = np.zeros(hidden_dim)

self.params[f'W{len(hidden_dims)+1}'] = weight_scale * np.random.randn(input_dim_current, num_classes)
self.params[f'b{len(hidden_dims)+1}'] = np.zeros(num_classes)

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

# When using dropout we need to pass a dropout_param dictionary to each
# dropout layer so that the layer knows the dropout probability and the mode
# (train / test). You can pass the same dropout_param to each dropout layer.
self.dropout_param = {}
if self.use_dropout:
    self.dropout_param = {'mode': 'train', 'p': dropout}
    if seed is not None:
        self.dropout_param['seed'] = seed

# With batch normalization we need to keep track of running means and
# variances, so we need to pass a special bn_param object to each batch
# normalization layer. You should pass self.bn_params[0] to the forward pass
# of the first batch normalization layer, self.bn_params[1] to the forward
# pass of the second batch normalization layer, etc.
self.bn_params = []
if self.use_batchnorm:
    self.bn_params = [{ 'mode': 'train'} for i in np.arange(self.num_layers - 1)]

# Cast all parameters to the correct datatype
for k, v in self.params.items():
    self.params[k] = v.astype(dtype)

def loss(self, X, y=None):
    """
    Compute loss and gradient for the fully-connected net.


```

```

Input / output: Same as TwoLayerNet above.
"""

X = X.astype(self.dtype)
mode = 'test' if y is None else 'train'

# Set train/test mode for batchnorm params and dropout param since they
# behave differently during training and testing.
if self.dropout_param is not None:
    self.dropout_param['mode'] = mode
if self.use_batchnorm:
    for bn_param in self.bn_params:
        bn_param[mode] = mode

scores = None
caches = {}
layer_input = X

# ===== #
# YOUR CODE HERE:
#   Implement the forward pass of the FC net and store the output
#   scores as the variable "scores".
# ===== #

for i in range(1, self.num_layers):
    Wi, bi = self.params[f'W{i}'], self.params[f'b{i}']
    layer_input, caches[i] = affine_relu_forward(layer_input, Wi, bi)

# Last layer (affine only)
scores, caches[self.num_layers] = affine_forward(layer_input,
                                                self.params[f'W{self.num_layers}'],
                                                self.params[f'b{self.num_layers}'])

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

# If test mode return early
if mode == 'test':
    return scores

loss, grads = 0.0, {}
# ===== #
# YOUR CODE HERE:
#   Implement the backwards pass of the FC net and store the gradients
#   in the grads dict, so that grads[k] is the gradient of self.params[k]
#   Be sure your L2 regularization includes a 0.5 factor.
# ===== #

loss, dscores = softmax_loss(scores, y)
loss += 0.5 * self.reg * sum(np.sum(self.params[f'W{i}']**2) for i in range(1, self.num_layers + 1))

dout, grads[f'W{self.num_layers}'], grads[f'b{self.num_layers}'] = affine_backward(dscores, caches[self.num_layers])
grads[f'W{self.num_layers}'] += self.reg * self.params[f'W{self.num_layers}']

for i in reversed(range(1, self.num_layers)):
    dout, grads[f'W{i}'], grads[f'b{i}'] = affine_relu_backward(dout, caches[i])
    grads[f'W{i}'] += self.reg * self.params[f'W{i}']

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

return loss, grads

```

Fully connected networks

In the previous notebook, you implemented a simple two-layer neural network class. However, this class is not modular. If you wanted to change the number of layers, you would need to write a new loss and gradient function. If you wanted to optimize the network with different optimizers, you'd need to write new training functions. If you wanted to incorporate regularizations, you'd have to modify the loss and gradient function.

Instead of having to modify functions each time, for the rest of the class, we'll work in a more modular framework where we define forward and backward layers that calculate losses and gradients respectively. Since the forward and backward layers share intermediate values that are useful for calculating both the loss and the gradient, we'll also have these functions return "caches" which store useful intermediate values.

The goal is that through this modular design, we can build different sized neural networks for various applications.

In this HW #3, we'll define the basic architecture, and in HW #4, we'll build on this framework to implement different optimizers and regularizations (like BatchNorm and Dropout).

Modular layers

This notebook will build modular layers in the following manner. First, there will be a forward pass for a given layer with inputs (`x`) and return the output of that layer (`out`) as well as cached variables (`cache`) that will be used to calculate the gradient in the backward pass.

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

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

    return out, cache
```

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

```
def layer_backward(dout, cache):
    """
    Receive derivative of loss with respect to outputs and cache,
    and compute derivative with respect to inputs.
    """

    # Unpack cache values
    x, w, z, out = cache

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

    return dx, dw
```

```
In [1]: 1 ## Import and setups
2
3 import time
4 import numpy as np
5 import matplotlib.pyplot as plt
6 from nnndl.fc_net import *
7 from utils.data_utils import get_CIFAR10_data
8 from utils.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
9 from utils.solver import Solver
10
11 %matplotlib inline
12 plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
13 plt.rcParams['image.interpolation'] = 'nearest'
14 plt.rcParams['image.cmap'] = 'gray'
15
16 # for auto-reloading external modules
17 # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
18 %load_ext autoreload
19 %autoreload 2
20
21 def rel_error(x, y):
22     """ returns relative error """
23     return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

```
In [4]: 1 # Load the (preprocessed) CIFAR10 data.
2 data = get_CIFAR10_data()
3 for k in data.keys():
4     print('{0}: {1}'.format(k, data[k].shape))
```

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

Linear layers

In this section, we'll implement the forward and backward pass for the linear layers.

The linear layer forward pass is the function `affine_forward` in `nnndl/layers.py` and the backward pass is `affine_backward`.

After you have implemented these, test your implementation by running the cell below.

Affine layer forward pass

Implement `affine_forward` and then test your code by running the following cell.

```
In [17]: 1 # Test the affine_forward function
2
3 num_inputs = 2
4 input_shape = (4, 5, 6)
5 output_dim = 3
6
7 input_size = num_inputs * np.prod(input_shape)
8 weight_size = output_dim * np.prod(input_shape)
9
10 x = np.linspace(-0.1, 0.5, num=input_size).reshape(num_inputs, *input_shape)
11 w = np.linspace(-0.2, 0.3, num=weight_size).reshape(np.prod(input_shape), output_dim)
12 b = np.linspace(-0.3, 0.1, num=output_dim)
13
14 out, _ = affine_forward(x, w, b)
15 correct_out = np.array([[ 1.49834967,  1.70660132,  1.91485297],
16                         [ 3.25553199,  3.5141327,   3.77273342]])
17
18 # Compare your output with ours. The error should be around 1e-9.
19 print('Testing affine_forward function:')
20 print('difference: {}'.format(rel_error(out, correct_out)))
```

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

Affine layer backward pass

Implement `affine_backward` and then test your code by running the following cell.

```
In [18]: 1 # Test the affine_backward function
2
3 x = np.random.randn(10, 2, 3)
4 w = np.random.randn(6, 5)
5 b = np.random.randn(5)
6 dout = np.random.randn(10, 5)
7
8 dx_num = eval_numerical_gradient_array(lambda x: affine_forward(x, w, b)[0], x, dout)
9 dw_num = eval_numerical_gradient_array(lambda w: affine_forward(x, w, b)[0], w, dout)
10 db_num = eval_numerical_gradient_array(lambda b: affine_forward(x, w, b)[0], b, dout)
11
12 _, cache = affine_forward(x, w, b)
13 dx, dw, db = affine_backward(dout, cache)
14
15 # The error should be around 1e-10
16 print('Testing affine_backward function:')
17 print('dx error: {}'.format(rel_error(dx_num, dx)))
18 print('dw error: {}'.format(rel_error(dw_num, dw)))
19 print('db error: {}'.format(rel_error(db_num, db)))
```

Testing affine_backward function:
dx error: 1.92096764519679e-09
dw error: 1.096709875458812e-10
db error: 5.0588919335003475e-11

Activation layers

In this section you'll implement the ReLU activation.

ReLU forward pass

Implement the `relu_forward` function in `nndl/layers.py` and then test your code by running the following cell.

```
In [19]: 1 # Test the relu_forward function
2
3 x = np.linspace(-0.5, 0.5, num=12).reshape(3, 4)
4
5 out, _ = relu_forward(x)
6 correct_out = np.array([[ 0.,           0.,           0.,           0.,           ],
7                         [ 0.,           0.,           0.04545455,  0.13636364,],
8                         [ 0.22727273,  0.31818182,  0.40909091,  0.5,        ]])
9
10 # Compare your output with ours. The error should be around 1e-8
11 print('Testing relu_forward function:')
12 print('difference: {}'.format(rel_error(out, correct_out)))
```

Testing relu_forward function:
difference: 4.99999798022158e-08

ReLU backward pass

Implement the `relu_backward` function in `nndl/layers.py` and then test your code by running the following cell.

```
In [20]: 1 x = np.random.randn(10, 10)
2 dout = np.random.randn(*x.shape)
3
4 dx_num = eval_numerical_gradient_array(lambda x: relu_forward(x)[0], x, dout)
5
6 _, cache = relu_forward(x)
7 dx = relu_backward(dout, cache)
8
9 # The error should be around 1e-12
10 print('Testing relu_backward function:')
11 print('dx error: {}'.format(rel_error(dx_num, dx)))
```

Testing relu_backward function:
dx error: 3.2756195039972447e-12

Combining the affine and ReLU layers

Often times, an affine layer will be followed by a ReLU layer. So let's make one that puts them together. Layers that are combined are stored in `nndl/layer_utils.py`.

Affine-ReLU layers

We've implemented `affine_relu_forward()` and `affine_relu_backward` in `nndl/layer_utils.py`. Take a look at them to make sure you understand what's going on. Then run the following cell to ensure its implemented correctly.

```
In [21]: 1 from nndl.layer_utils import affine_relu_forward, affine_relu_backward
2
3 x = np.random.randn(2, 3, 4)
4 w = np.random.randn(12, 10)
5 b = np.random.randn(10)
6 dout = np.random.randn(2, 10)
7
8 out, cache = affine_relu_forward(x, w, b)
9 dx, dw, db = affine_relu_backward(dout, cache)
10
11 dx_num = eval_numerical_gradient_array(lambda x: affine_relu_forward(x, w, b)[0], x, dout)
12 dw_num = eval_numerical_gradient_array(lambda w: affine_relu_forward(x, w, b)[0], w, dout)
13 db_num = eval_numerical_gradient_array(lambda b: affine_relu_forward(x, w, b)[0], b, dout)
14
15 print('Testing affine_relu_forward and affine_relu_backward:')
16 print('dx error: {}'.format(rel_error(dx_num, dx)))
17 print('dw error: {}'.format(rel_error(dw_num, dw)))
18 print('db error: {}'.format(rel_error(db_num, db)))
```

Testing affine_relu_forward and affine_relu_backward:
dx error: 2.1153280924652617e-10
dw error: 2.0913929148545147e-09
db error: 7.826653464481358e-12

Softmax loss

You've already implemented it, so we have written it in `layers.py`. The following code will ensure they are working correctly.

```
In [22]: 1 num_classes, num_inputs = 10, 50
2 x = 0.001 * np.random.randn(num_inputs, num_classes)
3 y = np.random.randint(num_classes, size=num_inputs)
4
5
6
7 dx_num = eval_numerical_gradient(lambda x: softmax_loss(x, y)[0], x, verbose=False)
8 loss, dx = softmax_loss(x, y)
9
10 # Test softmax_loss function. Loss should be 2.3 and dx error should be 1e-8
11 print('\nTesting softmax_loss:')
12 print('loss: {}'.format(loss))
13 print('dx error: {}'.format(rel_error(dx_num, dx)))
```

Testing softmax_loss:
loss: 2.30302268998151
dx error: 9.476234199045239e-09

Implementation of a two-layer NN

In `nndl/fc_net.py`, implement the class `TwoLayerNet` which uses the layers you made here. When you have finished, the following cell will test your implementation.

```

In [24]: 1 N, D, H, C = 3, 5, 50, 7
2 X = np.random.randn(N, D)
3 y = np.random.randint(C, size=N)
4
5 std = 1e-2
6 model = TwoLayerNet(input_dim=D, hidden_dims=H, num_classes=C, weight_scale=std)
7
8 print('Testing initialization ... ')
9 W1_std = abs(model.params['W1'].std() - std)
10 b1 = model.params['b1']
11 W2_std = abs(model.params['W2'].std() - std)
12 b2 = model.params['b2']
13 assert W1_std < std / 10, 'First layer weights do not seem right'
14 assert np.all(b1 == 0), 'First layer biases do not seem right'
15 assert W2_std < std / 10, 'Second layer weights do not seem right'
16 assert np.all(b2 == 0), 'Second layer biases do not seem right'
17
18 print('Testing test-time forward pass ... ')
19 model.params['W1'] = np.linspace(-0.7, 0.3, num=D*H).reshape(D, H)
20 model.params['b1'] = np.linspace(-0.1, 0.9, num=H)
21 model.params['W2'] = np.linspace(-0.3, 0.4, num=H*C).reshape(H, C)
22 model.params['b2'] = np.linspace(-0.9, 0.1, num=C)
23 X = np.linspace(-5.5, 4.5, num=N*D).reshape(D, N).T
24 scores = model.loss(X)
25 correct_scores = np.asarray(
26     [[11.53165108, 12.2917344, 13.05181771, 13.81190102, 14.57198434, 15.33206765, 16.09215096],
27      [12.05769098, 12.74614105, 13.43459113, 14.1230412, 14.81149128, 15.49994135, 16.18839143],
28      [12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.05099822, 15.66781506, 16.2846319]])
29 scores_diff = np.abs(scores - correct_scores).sum()
30 assert scores_diff < 1e-6, 'Problem with test-time forward pass'
31
32 print('Testing training loss (no regularization)')
33 y = np.asarray([0, 5, 1])
34 loss, grads = model.loss(X, y)
35 correct_loss = 3.4702243556
36 assert abs(loss - correct_loss) < 1e-10, 'Problem with training-time loss'
37
38 model.reg = 1.0
39 loss, grads = model.loss(X, y)
40 correct_loss = 26.5948426952
41 assert abs(loss - correct_loss) < 1e-10, 'Problem with regularization loss'
42
43 for reg in [0.0, 0.7]:
44     print('Running numeric gradient check with reg = {}'.format(reg))
45     model.reg = reg
46     loss, grads = model.loss(X, y)
47
48     for name in sorted(grads):
49         f = lambda _: model.loss(X, y)[0]
50         grad_num = eval_numerical_gradient(f, model.params[name], verbose=False)
51         print('{} relative error: {}'.format(name, rel_error(grad_num, grads[name])))

```

```

Testing initialization ...
Testing test-time forward pass ...
Testing training loss (no regularization)
Running numeric gradient check with reg = 0.0
W1 relative error: 1.5215703686475096e-08
W2 relative error: 3.2068321167375225e-10
b1 relative error: 8.368195737354163e-09
b2 relative error: 4.3291360264321544e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 2.527915175868136e-07
W2 relative error: 2.8508510893102143e-08
b1 relative error: 1.5646801536371197e-08
b2 relative error: 7.759095355706557e-10

```

Solver

We will now use the utils Solver class to train these networks. Familiarize yourself with the API in `utils/solver.py`. After you have done so, declare an instance of a TwoLayerNet with 200 units and then train it with the Solver. Choose parameters so that your validation accuracy is at least 50%.

In [32]:

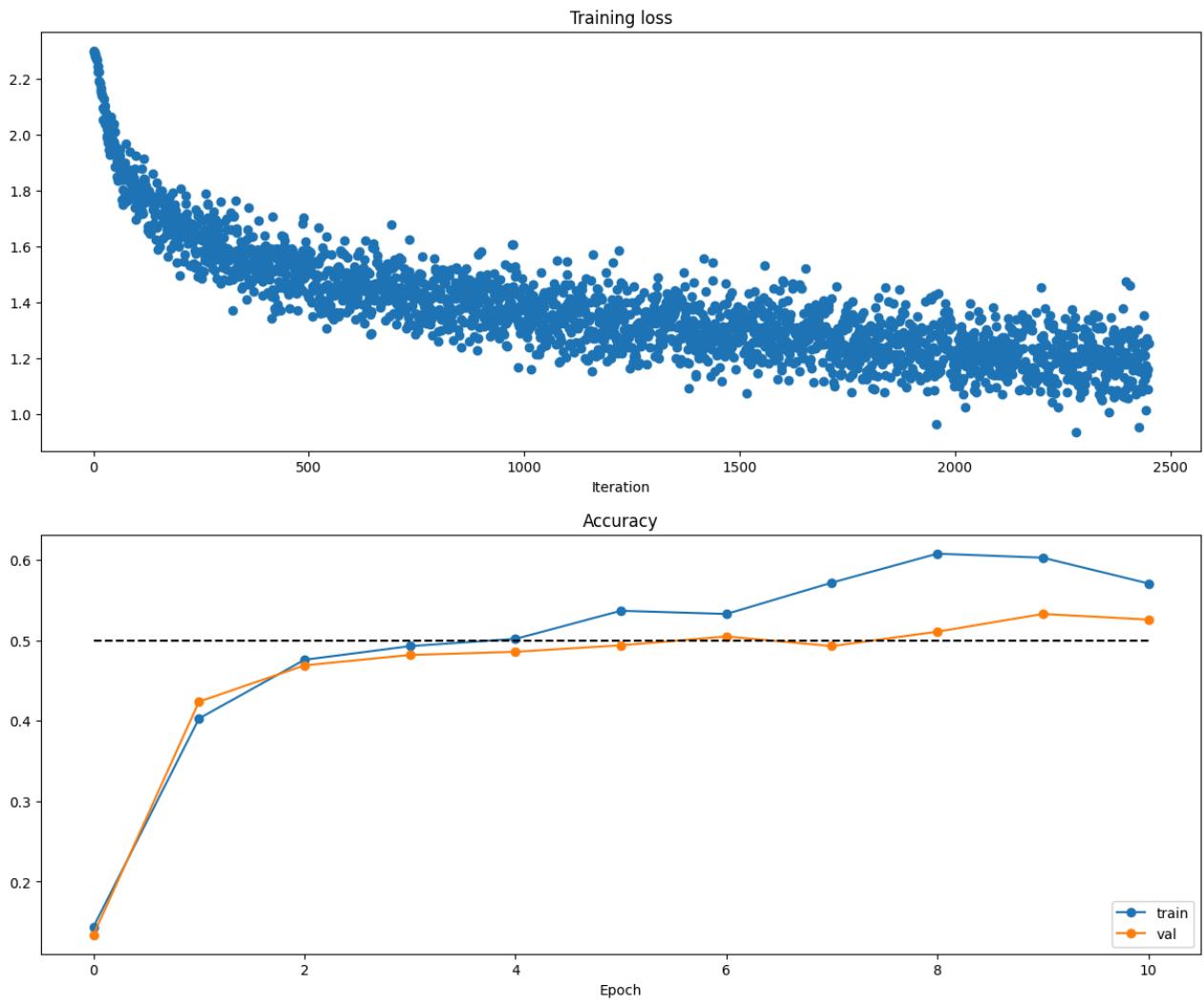
```
1 model = TwoLayerNet()
2 solver = None
3
4 # ===== #
5 # YOUR CODE HERE:
6 #   Declare an instance of a TwoLayerNet and then train
7 #   it with the Solver. Choose hyperparameters so that your validation
8 #   accuracy is at least 50%. We won't have you optimize this further
9 #   since you did it in the previous notebook.
10 #
11 # ===== #
12
13 solver = Solver(model, data,
14                  update_rule='sgd', optim_config={'learning_rate': 1e-3},
15                  lr_decay=0.95,
16                  num_epochs=10, batch_size=200,
17                  print_every=100)
18
19 solver.train()
20
21 # ===== #
22 # END YOUR CODE HERE
23 # ===== #
```

(Iteration 1 / 2450) loss: 2.302087
(Epoch 0 / 10) train acc: 0.144000; val_acc: 0.134000
(Iteration 101 / 2450) loss: 1.757412
(Iteration 201 / 2450) loss: 1.637018
(Epoch 1 / 10) train acc: 0.403000; val_acc: 0.424000
(Iteration 301 / 2450) loss: 1.526684
(Iteration 401 / 2450) loss: 1.537290
(Epoch 2 / 10) train acc: 0.476000; val_acc: 0.469000
(Iteration 501 / 2450) loss: 1.541296
(Iteration 601 / 2450) loss: 1.400074
(Iteration 701 / 2450) loss: 1.424359
(Epoch 3 / 10) train acc: 0.493000; val_acc: 0.482000
(Iteration 801 / 2450) loss: 1.295249
(Iteration 901 / 2450) loss: 1.474847
(Epoch 4 / 10) train acc: 0.502000; val_acc: 0.486000
(Iteration 1001 / 2450) loss: 1.293859
(Iteration 1101 / 2450) loss: 1.499670
(Iteration 1201 / 2450) loss: 1.323325
(Epoch 5 / 10) train acc: 0.537000; val_acc: 0.494000
(Iteration 1301 / 2450) loss: 1.384211
(Iteration 1401 / 2450) loss: 1.290005
(Epoch 6 / 10) train acc: 0.533000; val_acc: 0.505000
(Iteration 1501 / 2450) loss: 1.324360
(Iteration 1601 / 2450) loss: 1.399059
(Iteration 1701 / 2450) loss: 1.288299
(Epoch 7 / 10) train acc: 0.572000; val_acc: 0.493000
(Iteration 1801 / 2450) loss: 1.161164
(Iteration 1901 / 2450) loss: 1.163402
(Epoch 8 / 10) train acc: 0.608000; val_acc: 0.511000
(Iteration 2001 / 2450) loss: 1.148272
(Iteration 2101 / 2450) loss: 1.139933
(Iteration 2201 / 2450) loss: 1.451878
(Epoch 9 / 10) train acc: 0.603000; val_acc: 0.533000
(Iteration 2301 / 2450) loss: 1.079785
(Iteration 2401 / 2450) loss: 1.288225
(Epoch 10 / 10) train acc: 0.571000; val_acc: 0.526000

In [33]:

```
1 # Run this cell to visualize training loss and train / val accuracy
2
3 plt.subplot(2, 1, 1)
4 plt.title('Training loss')
5 plt.plot(solver.loss_history, 'o')
6 plt.xlabel('Iteration')
7
8 plt.subplot(2, 1, 2)
9 plt.title('Accuracy')
10 plt.plot(solver.train_acc_history, '-o', label='train')
11 plt.plot(solver.val_acc_history, '-o', label='val')
12 plt.plot([0.5] * len(solver.val_acc_history), 'k--')
13 plt.xlabel('Epoch')
14 plt.legend(loc='lower right')
15 plt.gcf().set_size_inches(15, 12)
16 plt.show()
```

```
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "orientation" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "dpi" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "facecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "edgecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "bbox_inches_restore" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
```



Multilayer Neural Network

Now, we implement a multi-layer neural network.

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

Implement the initialization, the forward pass, and the backward pass. There will be lines for batchnorm and dropout layers and caches; ignore these all for now. That'll be in HW #4.

```
In [25]: 1 N, D, H1, H2, C = 2, 15, 20, 30, 10
2 X = np.random.randn(N, D)
3 y = np.random.randint(C, size=(N,))
4
5 for reg in [0, 3.14]:
6     print('Running check with reg = {}'.format(reg))
7     model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
8                               reg=reg, weight_scale=5e-2, dtype=np.float64)
9
10    loss, grads = model.loss(X, y)
11    print('Initial loss: {}'.format(loss))
12
13    for name in sorted(grads):
14        f = lambda _: model.loss(X, y)[0]
15        grad_num = eval_numerical_gradient(f, model.params[name], verbose=False, h=1e-5)
16        print('{} relative error: {}'.format(name, rel_error(grad_num, grads[name])))
```

Running check with reg = 0
Initial loss: 2.2985850094266307
W1 relative error: 1.5638062961849955e-07
W2 relative error: 2.1676748168737265e-06
W3 relative error: 1.0089230525404633e-07
b1 relative error: 1.0382484882526912e-08
b2 relative error: 8.753554622271737e-10
b3 relative error: 1.602629914844274e-10
Running check with reg = 3.14
Initial loss: 7.22328144849468
W1 relative error: 4.403277559810739e-08
W2 relative error: 3.934021782261803e-08
W3 relative error: 2.2109530606580983e-08
b1 relative error: 1.8765746368643743e-08
b2 relative error: 1.2119045814856181e-08
b3 relative error: 2.410833343106932e-10

In [29]:

```
1 # Use the three layer neural network to overfit a small dataset.
2
3 num_train = 50
4 small_data = {
5     'X_train': data['X_train'][:num_train],
6     'y_train': data['y_train'][:num_train],
7     'X_val': data['X_val'],
8     'y_val': data['y_val'],
9 }
10
11 #####
12 ##### !!!!!!
13 # Play around with the weight_scale and learning_rate so that you can overfit a small dataset.
14 # Your training accuracy should be 1.0 to receive full credit on this part.
15 weight_scale = 2*1e-2
16 learning_rate = 1e-2
17
18 model = FullyConnectedNet([100, 100],
19                         weight_scale=weight_scale, dtype=np.float64)
20 solver = Solver(model, small_data,
21                  print_every=10, num_epochs=20, batch_size=25,
22                  update_rule='sgd',
23                  optim_config={
24                      'learning_rate': learning_rate,
25                  })
26
27 solver.train()
28
29 plt.plot(solver.loss_history, 'o')
30 plt.title('Training loss history')
31 plt.xlabel('Iteration')
32 plt.ylabel('Training loss')
33 plt.show()
```

```
(Iteration 1 / 40) loss: 3.182387
(Epoch 0 / 20) train acc: 0.360000; val_acc: 0.123000
(Epoch 1 / 20) train acc: 0.560000; val_acc: 0.121000
(Epoch 2 / 20) train acc: 0.600000; val_acc: 0.126000
(Epoch 3 / 20) train acc: 0.700000; val_acc: 0.173000
(Epoch 4 / 20) train acc: 0.760000; val_acc: 0.141000
(Epoch 5 / 20) train acc: 0.880000; val_acc: 0.146000
(Iteration 11 / 40) loss: 0.870777
(Epoch 6 / 20) train acc: 0.900000; val_acc: 0.132000
(Epoch 7 / 20) train acc: 0.980000; val_acc: 0.169000
(Epoch 8 / 20) train acc: 0.980000; val_acc: 0.158000
(Epoch 9 / 20) train acc: 1.000000; val_acc: 0.171000
(Epoch 10 / 20) train acc: 1.000000; val_acc: 0.179000
(Iteration 21 / 40) loss: 0.030894
(Epoch 11 / 20) train acc: 1.000000; val_acc: 0.175000
(Epoch 12 / 20) train acc: 1.000000; val_acc: 0.174000
(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.183000
(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.181000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.182000
(Iteration 31 / 40) loss: 0.013713
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.180000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.176000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.179000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.178000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.176000

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "orientation" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "dpi" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "facecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "edgecolor" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "bbox_inches_restore" which is no longer supported as of 3.3 and will become an error two minor releases later
    fig.canvas.print_figure(bytes_io, **kw)
```

Training loss history

