

## EEE443 Neural Networks Assignment #1

### Question 1)

In this question, the prior probability distribution of the network weights for the following optimization problem is asked:

$$W_{MAP} = \operatorname{argmin}_W \sum_n (y^n - h(x^n, W))^2 + \beta \sum_i \omega_i^2$$

where  $W$  is the vector of weights  $\omega_i$ ,  $\beta$  is a scalar constant, and  $h(\cdot)$  is the output of the neuron.

We are trying to find the weight vector that maximizes the posterior. It can be written as follows:

$$W_{MAP} = \operatorname{arg max}_W P(W|D)$$

Where  $D$  is the dataset comprised of inputs and outputs,  $\{(x^n, y^n)\}$  where  $(n \in [1 N])$ .

To calculate this, first we need to recognize Bayes Rule:

$$P(W|D) = \frac{P(D|W)P(W)}{P(D)}$$

Since  $P(D)$  does not depend on  $W$ , we can take  $P(D)$  as a constant from now on. Therefore, our maximization problem for  $P(W|D)$  is equal to maximizing the numerator, which is  $P(D|W)P(W)$ . Our new problem is in the following form:

$$W_{MAP} = \operatorname{arg max}_W P(D|W)P(W)$$

After that, to make mathematical manipulation, we take the negative natural logarithm of both sides. Taking the logarithm will not change the optimum location because it is  $\ln$  function is a monotonically increasing function. Therefore, there is no problem doing this mathematical manipulation. The resulting equations are given below, step by step:

$$W_{MAP} = \operatorname{arg max}_W [\ln(P(D|W)) + \ln(P(W))]$$

$$W_{MAP} = \operatorname{arg min}_W [-\ln(P(D|W)) - \ln(P(W))]$$

The resulting equation found above is similar to the optimization problem equation that is given to us initially. To make connections between these two equations, I separated the found equation into two parts and analyze the relation between them further.

$$\sum_n (y^n - h(x^n, W))^2 - \ln(P(D|W))$$

As seen above, the equations are related to each because both have the variables  $x$  and  $y$ , input and output variables. Also the remaining terms are related to each other:

$$\beta \sum_i \omega_i^2$$

$$- \ln(P(W))$$

Again, these terms contain the same variable again. So, we derived a relation between the terms of the equations. We can show them as follows:

**Log-Likelihood:**  $-\ln(P(D|W))$  directly proportional with  $\sum_n (y^n - h(x^n, W))^2$

**Log-Prior:**  $-\ln(P(W))$  directly proportional with  $\beta \sum_i \omega_i^2$

Starting with the Log-Prior terms, we equate both sides (There is a constant C due to proportionality):

$$- \ln(P(W)) = \beta \sum_i \omega_i^2 + C$$

Then I exponentiate both sides to leave P(W) alone:

$$P(W) = \exp\left(-\beta \sum_i \omega_i^2 - C\right) = A * \exp\left(-\beta \sum_i \omega_i^2\right), \text{Where } A = e^{-C}$$

The next task is finding the constant A. To find the constant A, we can use the fundamental of sum over all probability distributions is equal to 1.

$$\int_{-\infty}^{\infty} \left[ A * \exp\left(-\beta \sum_i \omega_i^2\right) \right] dW = 1$$

We can factor out the constant A, out of the integration.

$$A * \int_{-\infty}^{\infty} \left[ \exp\left(-\beta \sum_i \omega_i^2\right) \right] dW = 1$$

The summation of exponentials can be written as the product of each individual term.

$$\exp\left(-\beta \sum_i \omega_i^2\right) = \prod_{i=1}^m \exp(-\beta \omega_i^2)$$

Using this conversion, we can evaluate a m-dimensional integral to find the constant A.

$$\int_{-\infty}^{\infty} \prod_{i=1}^m \exp(-\beta \omega_i^2) d\omega_i$$

This is the basic form of a Gaussian Integral. Looking at this table [1], we can derive its answers:

$$\int_{-\infty}^{\infty} \exp(-\beta \omega_i^2) d\omega_i = \sqrt{\frac{\pi}{\beta}}$$

$$\xrightarrow{\text{yields}} \int_{-\infty}^{\infty} \prod_{i=1}^m \exp(-\beta \omega_i^2) d\omega_i = \left(\sqrt{\frac{\pi}{\beta}}\right)^m$$

So, the value of constant can be found using these derivations:

$$A = \frac{1}{\int_{-\infty}^{\infty} [\exp(-\beta \sum_i \omega_i^2)] dW} = \sqrt{\frac{\beta}{\pi}}^m$$

Thus, we can rewrite the P(W), this time knowing the constant A. The Prior distribution is given below.

$$P(W) = \left( \sqrt{\frac{\beta}{\pi}} \right)^m * \exp \left( -\beta \sum_i \omega_i^2 \right)$$

Therefore, the prior probability distribution of the network weights is a joint zero-mean Gaussian (Normal) distribution, where each independent weight has a variance of  $\frac{1}{2\beta}$ .

## Question 2)

In this question, it is asked to implement the following logic equation to a neural network with a single hidden layer with four input neurons (with binary inputs) and a single output neuron:

$$(X_1 \text{ OR } \text{NOT } X_2) \text{ XOR } (\text{NOT } X_3 \text{ OR } \text{NOT } X_4)$$

This logic equation can also be written as follows:

$$f = (X_1 \vee \bar{X}_2) \text{ XOR } (\bar{X}_3 \vee \bar{X}_4)$$

Assume a hidden layer with four hidden units, and a unipolar activation function (i.e., the step function).

**Part a)** For each hidden unit, analytically derive the set of inequalities based on which a set of weights and an activation threshold can be selected.

To do so, we first need to write XOR as a combination of AND, OR logic functions. The XOR identity is given below:

$$X \text{ XOR } Y = (X \wedge \bar{Y}) \vee (\bar{X} \wedge Y)$$

$$X = (X_1 \vee \bar{X}_2), Y = (\bar{X}_3 \vee \bar{X}_4)$$

$$\bar{X} = (\bar{X}_1 \wedge X_2), \bar{Y} = (X_3 \wedge X_4)$$

Using this identity, the logic equation can be written again in terms of AND, OR functions.

$$f = [(X_1 \vee \bar{X}_2) \wedge (X_3 \wedge X_4)] \vee [(\bar{X}_1 \wedge X_2) \wedge (\bar{X}_3 \vee \bar{X}_4)]$$

Distributing the function will yield into the following logic function:

$$f = (X_1 \wedge X_3 \wedge X_4) \vee (\bar{X}_2 \wedge X_3 \wedge X_4) \vee (\bar{X}_1 \wedge X_2 \wedge \bar{X}_3) \vee (\bar{X}_1 \wedge X_2 \wedge \bar{X}_4)$$

This is a basic Perceptron model with a unipolar step using 4 hidden layers, which are depicted below.

$$h_1 = (X_1 \wedge X_3 \wedge X_4)$$

$$h_2 = (\bar{X}_2 \wedge X_3 \wedge X_4)$$

$$h_3 = (\bar{X}_1 \wedge X_2 \wedge \bar{X}_3)$$

$$h_4 = (\bar{X}_1 \wedge X_2 \wedge \bar{X}_4)$$

The weight vector  $W$  can be defined as the matrix shown below.  $\theta$  is the bias term for each neuron, and  $\omega$  represent the weight associated with each neuron.

$$W = \begin{bmatrix} \omega_{11} & \omega_{12} & \omega_{13} & \omega_{14} & \theta_1 \\ \omega_{21} & \omega_{22} & \omega_{23} & \omega_{24} & \theta_2 \\ \omega_{31} & \omega_{32} & \omega_{33} & \omega_{34} & \theta_3 \\ \omega_{41} & \omega_{42} & \omega_{43} & \omega_{44} & \theta_4 \end{bmatrix}$$

$W_h$  represents the weights of the output layer, which is built from the hidden layer results.

$$W_h = [\omega_{51} \quad \omega_{52} \quad \omega_{53} \quad \omega_{54} \quad \theta_5]$$

The output of the function with individual weights are given below:

$$y = \omega_{11}X_1 + \omega_{12}X_2 + \omega_{13}X_3 + \omega_{14}X_4 + \text{bias terms}$$

The output of the neurons will be determined regarding the output of the step function. To determine this one has to look of the output of the step function whether it is 0 or 1.

$$h(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

So, the output will be as follows:

$$OUTPUT = h(\omega_{11}X_1 + \omega_{12}X_2 + \omega_{13}X_3 + \omega_{14}X_4 + \text{bias terms})$$

The truth table for each hidden layer will be examined. Below is the truth table analysis for the first hidden layer. Because  $X_2$  is not present, it's weight is assigned to zero and not shown in the table. This will be the general case for the following hidden layers, the missing input will have pre-assigned 0 weight.

$$h_1 = (X_1 \wedge X_3 \wedge X_4)$$

Table 1: Truth table for first hidden layer

$X_1$	$X_3$	$X_4$	$h_1$	Inequality
0	0	0	0	$-\theta_1 < 0$
0	0	1	0	$\omega_{14} - \theta_1 < 0$
0	1	0	0	$\omega_{13} - \theta_1 < 0$
0	1	1	0	$\omega_{13} + \omega_{14} - \theta_1 < 0$
1	0	0	0	$\omega_{11} - \theta_1 < 0$
1	0	1	0	$\omega_{11} + \omega_{14} - \theta_1 < 0$
1	1	0	0	$\omega_{11} + \omega_{13} - \theta_1 < 0$
1	1	1	1	$\omega_{11} + \omega_{13} + \omega_{14} - \theta_1 > 0$

$$h_2 = (\bar{X}_2 \wedge X_3 \wedge X_4)$$

Table 2: Truth table for second hidden layer

$\bar{X}_2$	$X_3$	$X_4$	$h_2$	Inequality
0	0	0	0	$-\theta_2 < 0$
0	0	1	0	$\omega_{24} - \theta_2 < 0$

0	1	0	0	$\omega_{23} - \theta_2 < 0$
0	1	1	1	$\omega_{23} + \omega_{24} - \theta_2 > 0$
1	0	0	0	$\omega_{22} - \theta_2 < 0$
1	0	1	0	$\omega_{22} + \omega_{24} - \theta_2 < 0$
1	1	0	0	$\omega_{22} + \omega_{23} - \theta_2 < 0$
1	1	1	0	$\omega_{22} + \omega_{23} + \omega_{24} - \theta_2 < 0$

$$h_3 = (\bar{X}_1 \wedge X_2 \wedge \bar{X}_3)$$

Table 3: Truth table for third hidden layer

$\bar{X}_1$	$X_2$	$\bar{X}_3$	$h_3$	Inequality
0	0	0	0	$-\theta_3 < 0$
0	0	1	0	$\omega_{33} - \theta_3 < 0$
0	1	0	1	$\omega_{32} - \theta_3 > 0$
0	1	1	0	$\omega_{32} + \omega_{33} - \theta_3 < 0$
1	0	0	0	$\omega_{31} - \theta_3 < 0$
1	0	1	0	$\omega_{31} + \omega_{33} - \theta_3 < 0$
1	1	0	0	$\omega_{31} + \omega_{32} - \theta_3 < 0$
1	1	1	0	$\omega_{31} + \omega_{32} + \omega_{33} - \theta_3 < 0$

$$h_4 = (\bar{X}_1 \wedge X_2 \wedge \bar{X}_4)$$

Table 4: Truth table for fourth hidden layer

$\bar{X}_1$	$X_2$	$\bar{X}_4$	$h_4$	Inequality
0	0	0	0	$-\theta_4 < 0$
0	0	1	0	$\omega_{44} - \theta_4 < 0$
0	1	0	1	$\omega_{42} - \theta_4 > 0$
0	1	1	0	$\omega_{42} + \omega_{44} - \theta_4 < 0$
1	0	0	0	$\omega_{41} - \theta_4 < 0$
1	0	1	0	$\omega_{41} + \omega_{44} - \theta_4 < 0$
1	1	0	0	$\omega_{41} + \omega_{42} - \theta_4 < 0$
1	1	1	0	$\omega_{41} + \omega_{42} + \omega_{44} - \theta_4 < 0$

$$h_5 = (h_1 \vee h_2 \vee h_3 \vee h_4)$$

Table 5: Truth table for output hidden layer

$h_1$	$h_2$	$h_3$	$h_4$	$h_5$	Inequality
0	0	0	0	0	$-\theta_5 < 0$
0	0	0	1	1	$\omega_{54} - \theta_5 > 0$
0	0	1	0	1	$\omega_{53} - \theta_5 > 0$
0	0	1	1	1	$\omega_{53} + \omega_{54} - \theta_5 > 0$
0	1	0	0	1	$\omega_{52} - \theta_5 > 0$
0	1	0	1	1	$\omega_{52} + \omega_{54} - \theta_5 > 0$
0	1	1	0	1	$\omega_{52} + \omega_{53} - \theta_5 > 0$
0	1	1	1	1	$\omega_{52} + \omega_{53} + \omega_{54} - \theta_5 > 0$
1	0	0	0	1	$\omega_{51} - \theta_5 > 0$
1	0	0	1	1	$\omega_{51} + \omega_{54} - \theta_5 > 0$

1	0	1	0	1	$\omega_{51} + \omega_{53} - \theta_5 > 0$
1	0	1	1	1	$\omega_{51} + \omega_{53} + \omega_{54} - \theta_5 > 0$
1	1	0	0	1	$\omega_{51} + \omega_{52} - \theta_5 > 0$
1	1	0	1	1	$\omega_{51} + \omega_{52} + \omega_{54} - \theta_5 > 0$
1	1	1	0	1	$\omega_{51} + \omega_{52} + \omega_{53} - \theta_5 > 0$
1	1	1	1	1	$\omega_{51} + \omega_{52} + \omega_{53} + \omega_{54} - \theta_5 > 0$

**Part b)** To find a particular weight vector including the bias term, I used Python to find them. All of the Python code was appended at the end of this report. The comments inside the code explain the work done thoroughly, due to the page limit, explicit explanations are not done here.

Basically I defined an interval from -3 to 3 to pick valid numbers for weights and bias terms. Then I wrote the individual inequalities. The resulting weight vector including the bias terms is given below. After I wrote the logic functions again and analyzed whether the found weights and bias term values reach a 100% performance.

$$W = \begin{bmatrix} 1 & 0 & 1 & 1 & 2.5 \\ 0 & -3 & 1 & 1 & 1.5 \\ -3 & 1 & -3 & 0 & 0.5 \\ -3 & 1 & 0 & 3 & 0.5 \end{bmatrix}$$

$$W_h = [1 \quad 1 \quad 1 \quad 1 \quad 0.5]$$

```

0 0 0 0 | 0 | 0 | True
0 0 0 1 | 0 | 0 | True
0 0 1 0 | 0 | 0 | True
0 0 1 1 | 1 | 1 | True
0 1 0 0 | 1 | 1 | True
0 1 0 1 | 1 | 1 | True
0 1 1 0 | 1 | 1 | True
0 1 1 1 | 0 | 0 | True
1 0 0 0 | 0 | 0 | True
1 0 0 1 | 0 | 0 | True
1 0 1 0 | 0 | 0 | True
1 0 1 1 | 1 | 1 | True
1 1 0 0 | 0 | 0 | True
1 1 0 1 | 0 | 0 | True
1 1 1 0 | 0 | 0 | True
1 1 1 1 | 1 | 1 | True
-----
Total Network Accuracy: 100.0%

```

Fig. 1: 100% performance has achieved with the found weights and bias terms

**Part c)** In this part, it is asked to find another weight matrix when the input terms ( $x_i$ 's) are subject to a small noise. To get the most robust function to work with best performance under noisy conditions, I need to keep on one thing, which is the amplitude of the individual weights'. When the amplitude gets bigger, unexpectedly the noise variance will be automatically amplified. Therefore, to achieve the most robust function, the weights must have been kept as small as possible. Thus, rather than the weights that are calculated in part b), I will use only  $\pm 1$  values for the weights. So the weights can take values  $\{-1, 0, 1\}$ . +1 if the input is +1 and -1 if the input is 0. After that, I need to tackle bias terms in addition. In order to stay in one side of the inequalities of the hidden layers, bias terms must be placed right at the center of these inequalities, so that logic gate results will not flip (step response function). Below is given an example to explain the procedure:

$$h_1 = (X_1 \wedge X_3 \wedge X_4)$$

$\omega_{11} + \omega_{13} - \theta_1 < 0$
$\omega_{11} + \omega_{13} + \omega_{14} - \theta_1 > 0$

Using these two inequalities, and assigning +1 weights to present weights, we can achieve an interval for the bias term  $\theta_1$ .

$$\omega_{11} = 1, \omega_{13} = 1, \omega_{14} = 1$$

$2 < \theta_1$
$3 > \theta_1$

So:

$$3 > \theta_1 > 2$$

$$\theta_1 = 2.5$$

To find the best value, we choose the center point, which is 2.5 for the first bias term. The same procedure is done to other hidden layers and each bias term was found.

$$h_2 = (\bar{X}_2 \wedge X_3 \wedge X_4)$$

$\omega_{23} - \theta_2 < 0$
$\omega_{23} + \omega_{24} - \theta_2 > 0$

$$\omega_{22} = -1, \omega_{23} = 1, \omega_{24} = 1$$

$1 < \theta_2$
$2 > \theta_2$

$$\theta_2 = 1.5$$

$$h_3 = (\bar{X}_1 \wedge X_2 \wedge \bar{X}_3)$$

$\omega_{32} - \theta_3 > 0$
$\omega_{32} + \omega_{33} - \theta_3 < 0$

$$\omega_{31} = -1, \omega_{32} = 1, \omega_{33} = -1$$

$0 < \theta_3$
$1 > \theta_3$

$$\theta_3 = 0.5$$

$$h_4 = (\bar{X}_1 \wedge X_2 \wedge \bar{X}_4)$$

$\omega_{42} - \theta_4 > 0$
$\omega_{42} + \omega_{44} - \theta_4 < 0$

$$\omega_{41} = -1, \omega_{42} = 1, \omega_{44} = -1$$

$0 < \theta_4$
$1 > \theta_4$

$$\theta_4 = 0.5$$

$$h_5 = (h_1 \vee h_2 \vee h_3 \vee h_4)$$

$$\begin{array}{|c|} \hline -\theta_5 < 0 \\ \hline \omega_{54} - \theta_5 > 0 \\ \hline \omega_{51} = 1, \omega_{52} = 1, \omega_{53} = 1, \omega_{54} = 1 \\ \hline \end{array}$$

$$\begin{array}{|c|} \hline 0 < \theta_5 \\ \hline 1 > \theta_5 \\ \hline \theta_5 = 0.5 \\ \hline \end{array}$$

Hidden layer weights and output neuron weights are given in below two matrices. Basically, bias terms are the negative values of the  $\theta$ 's.

$$W = \begin{bmatrix} 1 & 0 & 1 & 1 & 2.5 \\ 0 & -1 & 1 & 1 & 1.5 \\ -1 & 1 & -1 & 0 & 0.5 \\ -1 & 1 & 0 & 1 & 0.5 \end{bmatrix}$$

$$W_h = [1 \quad 1 \quad 1 \quad 1 \quad 0.5]$$

**Part d)** The code for this part can be found at the end of the report. The comments written explains the work done in detail. For this part, to test the networks, 25 copies of each of the 16 possible 4-bit input vectors were generated, creating a validation dataset of 400 samples. For the part a, calculated hidden layer and output neuron weight matrix, the performance with zero mean and 0.2 standard deviation Gaussian noise is calculated as 80.75%. For the matrix calculated and found in part c, the performance is calculated as 92.25%.

```
Part A Performance: 80.75%
Part C Performance: 92.25%
```

Fig. 2: Comparison of the noisy input performances of Part a and Part c

**Question 3)** The visual results of this part can be seen at the end of the code.

**Part a)** First, the 4 arrays inside the h5 file is found, and their shapes are printed. Also, the number of distinct classes is found by checking the letters using indexing.

```
testims (1300, 28, 28)
testlbls (1300,)
trainims (5200, 28, 28)
trainlbls (5200,)
26
```

Fig. 3: Array shapes

The samples in trainims visualized using matplotlib.

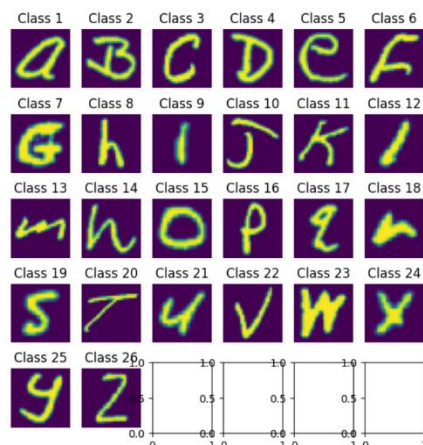


Fig. 4: Visualized letters

After that, using the built in correlation function in Python, the correlation matrix is created.

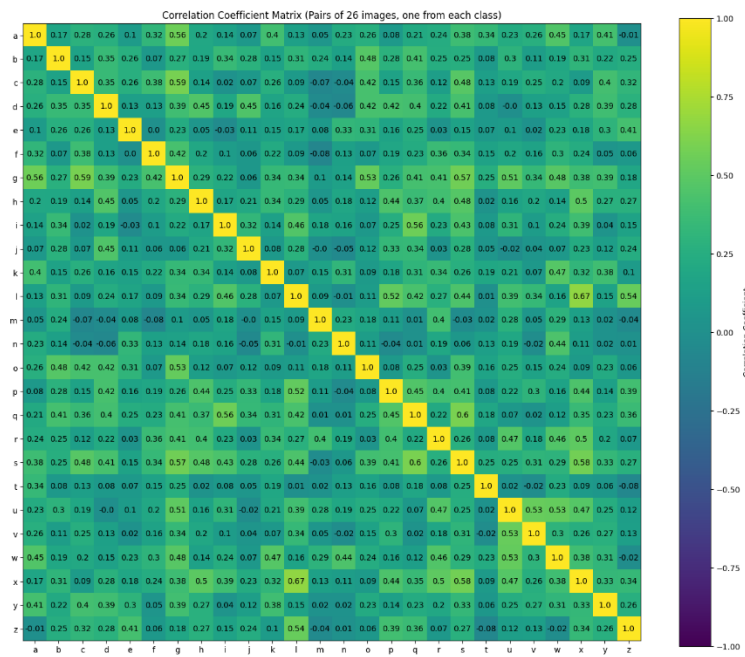


Fig. 5: Correlation matrix

As expected for the same letter pair, the correlation coefficient is equal to 1, which is the within-class variability. Also, for different letters, the coefficient is much smaller, which is the across-class variability. Because the dataset demonstrates high across-class variability and low within-class variability, the feature vectors for the different letters are highly separable. This makes the dataset an excellent candidate for classification using a single-layer perceptron or neural network.

**Part b)** For this part a single layer Perceptron with an output neuron for each digit is created. The optimum value for learning rate ( $\eta$ ) is chosen by looking at the results of MSE. After that, network weights for each letter is depicted. As seen from Fig. 1, the optimum value for learning rate is chosen as "0.12". Fig. 2 shows the final network weights for each digit.

eta=0.001	final MSE=0.034964
eta=0.01	final MSE=0.014692
eta=0.1	final MSE=0.010058
eta=0.12	final MSE=0.009721
eta=0.13	final MSE=0.009729

Fig. 6: Learning rate vs MSE

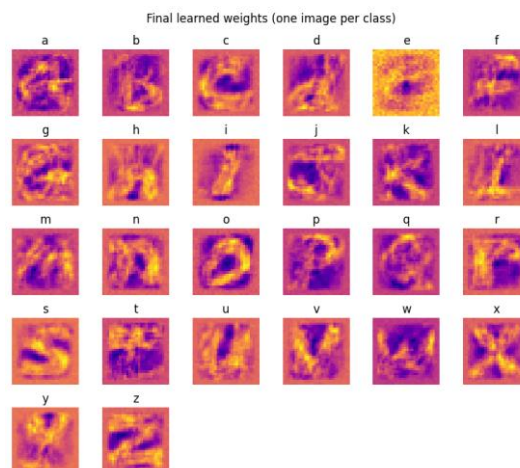
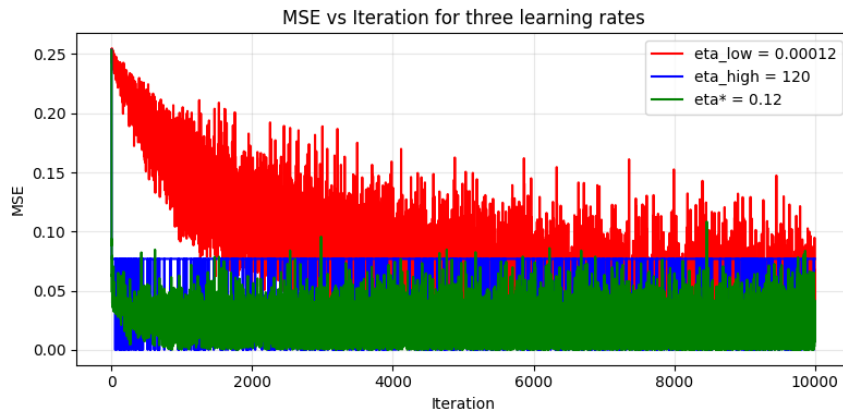


Fig. 7: Final network weights

**Part c)** As stated, the  $\eta$  value is chosen as 0.12 previously. For the low and high right, I respectively divide and multiply with 1000, which resulted in 0.00012 and 120 as the new rates. The MSE curves are given below.



*Fig. 8: MSE curves for different learning rates*

For the low  $\eta$  value (red curve), which is 0.00012, at the beginning, MSE curve drops down smoothly but the network learns slowly regarding other  $\eta$  values, also the step sizes (thickness of the red curve) are too small, which makes it under trained. For the high  $\eta$  value (blue curve), which is 120, MSE curve drops rapidly, however, the oscillation interval is too large for the MSE values, which makes the algorithm fluctuate and fail to settle, additionally, the step sizes are too big (thickness of the blue curve). For the normal  $\eta$  value (green curve), which is found to be 0.12, again the curve drops rapidly and stabilizes close to zero. The step sizes are great, no overshooting or under training in this case.

**Part d)** In this last part, the test samples are fed into the algorithm with above found three distinct learning rates, and the individual test performances are observed. As seen from Fig. 9, the optimal learning rate, 0.12 achieved the highest accuracy among all three networks, which is 55.77%. It successfully converged after 10000 iterations as seen from Fig. 8. Also, because this is a single layer neural network problem, 55.77% accuracy is a solid performance. For the low  $\eta$  value (0.00012), the test accuracy is 20.15%. The reason for low accuracy is slow learning and under training. Lastly for the high  $\eta$  value (120), this network performed the worst among all three networks, resulting in a 3.85% accuracy, the reason is again mentioned above, overshooting and being too aggressive. This results in a poor prediction ability. Therefore, as expected, when we go far away from the found learning rate, the accuracy decreases significantly.

```

Test performance (eta_low) : 20.15%
Test performance (eta*)   : 55.77%
Test performance (eta_high): 3.85%

```

*Fig. 9: Network accuracy results for different learning rates*

#### Question 4)

Because I reached the 10-page limit, the comments and answers to the inline questions can be found at the end of the report.

## References

- [1] J. Alison, "Gaussian Integrals," Univ. of Pennsylvania, Mar. 8, 2008. [Online]. Available: <https://www.hep.upenn.edu/~johnda/Papers/GausInt.pdf>.
- [2] Google, Gemini (Version 3.1 Pro). Google, 2026. [Online]. Available: <https://gemini.google.com>
- I get partial help from Gemini 3.1 Pro to write the code of Question 2 Part d, printing out the performance metric; and Question 3 Part b, I asked how to write the code for getting the MSE.
- [3] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA, USA: MIT Press, 2016. [Online]. Available: <http://www.deeplearningbook.org>
- [4] Prateek Kumar, "How to read and plot HDF5 file in python ?," *YouTube*, May 2, 2020. [Online]. Available: <https://www.youtube.com/watch?v=Jz9srOETL1M>
- [2] D. Chuan-En Lin, "8 Simple Techniques to Prevent Overfitting," TDS Archive (Medium), Jun. 7, 2020. [Online]. Available: <https://medium.com/data-science/8-simple-techniques-to-prevent-overfitting-4d443da2ef7d>.
- [6] A. Olaoye, "4 Effective Ways to Prevent Overfitting and Why They Work," TDS Archive (Towards Data Science), Medium, Oct. 27, 2022. [Online]. Available: <https://towardsdatascience.com/4-effective-ways-to-prevent-overfitting-and-why-they-work-f6e3b98aefda/>.

## Question 2 Part a

```

In [ ]: import itertools
import numpy as np

# I created a search interval between -3 to 3. This can be changed, trying diffe
# give 100% efficiency all the times because the code couldn't find any valid va
# of the weights.
space = [x * 0.5 for x in range(-6, 7)]

print("Searching for valid weights and thresholds...\n")

#h1
for w11, w13, w14, t1 in itertools.product(space, repeat=4):
    if (-t1 < 0 and
        w14 - t1 < 0 and
        w13 - t1 < 0 and
        w13 + w14 - t1 < 0 and
        w11 - t1 < 0 and
        w11 + w14 - t1 < 0 and
        w11 + w13 - t1 < 0 and
        w11 + w13 + w14 - t1 > 0):
        print(f"h1 parameters: w11={w11}, w12=0, w13={w13}, w14={w14}, theta1={t1}")
        break

#h2
for w22, w23, w24, t2 in itertools.product(space, repeat=4):
    if (-t2 < 0 and
        w24 - t2 < 0 and
        w23 - t2 < 0 and
        w23 + w24 - t2 > 0 and
        w22 - t2 < 0 and
        w22 + w24 - t2 < 0 and
        w22 + w23 - t2 < 0 and
        w22 + w23 + w24 - t2 < 0):
        print(f"h2 parameters: w21=0, w22={w22}, w23={w23}, w24={w24}, theta2={t2}")
        break

#h3
for w31, w32, w33, t3 in itertools.product(space, repeat=4):
    if (-t3 < 0 and
        w33 - t3 < 0 and
        w32 - t3 > 0 and
        w32 + w33 - t3 < 0 and
        w31 - t3 < 0 and
        w31 + w33 - t3 < 0 and
        w31 + w32 - t3 < 0 and
        w31 + w32 + w33 - t3 < 0):
        print(f"h3 parameters: w31={w31}, w32={w32}, w33={w33}, w34=0, theta3={t3}")
        break

#h4
for w41, w42, w44, t4 in itertools.product(space, repeat=4):
    if (-t4 < 0 and
        w44 - t4 < 0 and
        w42 - t4 > 0 and
        w42 + w44 - t4 < 0 and
        w41 - t4 < 0 and
    
```

```

w41 + w44 - t4 < 0 and
w41 + w42 - t4 < 0 and
w41 + w42 + w44 - t4 < 0):
print(f"h4 parameters: w41={w41}, w42={w42}, w43=0, w44={w44}, theta4={t4}")
break

#h5
for w51, w52, w53, w54, t5 in itertools.product(space, repeat=5):
    # Because h5 is OR gate of all other hidden layers, the inequalities can be
    if (-t5 < 0 and
        w54 - t5 > 0 and
        w53 - t5 > 0 and
        w52 - t5 > 0 and
        w51 - t5 > 0 and
        w51 + w52 + w53 + w54 - t5 > 0):
        print(f"h5 parameters: w51={w51}, w52={w52}, w53={w53}, w54={w54}, theta5={t5}")
        break

#The activation function definition, which is the step function in this problem.
def step_function(x):
    return np.where(x >= 0, 1, 0)

#Assigned initial 0 weights regarding the h1,h2,h3, and h4 hidden layer combination
w12 = 0
w21 = 0
w34 = 0
w43 = 0

#W matrix is created with the above found weights values.
W_hidden = np.array([
    [ w11, w12, w13, w14], # h1 weights
    [ w21, w22, w23, w24], # h2 weights
    [ w31, w32, w33, w34], # h3 weights
    [ w41, w42, w43, w44] # h4 weights
])

#Bias matrix is created, there is a minus sign in front of each, because bias terms
b_hidden = np.array([-t1, -t2, -t3, -t4])

#Output hidden layer matrix
W_out = np.array([w51, w52, w53, w54])

#To check the performance, we should look at the results of h5, which consists of
# h1, h2, h3, and h4, so we need 2^4 = 16 conditions in total. So from [0,0,0,0]
inputs = np.array(list(itertools.product([0, 1], repeat=4)))
correct_predictions = 0

print("X1 X2 X3 X4 | Target | NN Output | Match?")
print("-" * 45)

for X in inputs:
    #Ground truth weight
    term1 = X[0] or (not X[1]) #(X1 OR NOT X2)
    term2 = (not X[2]) or (not X[3]) #(NOT X3 OR NOT X4)

    # F = (X1 OR NOT X2) XOR (NOT X3 OR NOT X4), != is equal to XOR gate
    target = int(term1 != term2)

    #y = W*X + bias term

```

```

#output = h(y)
hidden_net = np.dot(W_hidden, X) + b_hidden
hidden_out = step_function(hidden_net)

#FinalOutput = Wh*X + bias term for output hidden Layer (theta5)
output_net = np.dot(W_out, hidden_out) - t5
nn_output = int(step_function(output_net))

#Check the equality, therefore the performance
match = (target == nn_output)
if match:
    correct_predictions += 1

print(f" {X[0]} {X[1]} {X[2]} {X[3]} | {target} | {nn_output}")

accuracy = (correct_predictions / 16.0) * 100
print("-" * 45)
print(f"Total Network Accuracy: {accuracy}%")

```

Searching for valid weights and thresholds...

h1 parameters: w11=1.0, w12=0, w13=1.0, w14=1.0, theta1=2.5

h2 parameters: w21=0, w22=-3.0, w23=1.0, w24=1.0, theta2=1.5

h3 parameters: w31=-3.0, w32=1.0, w33=-3.0, w34=0, theta3=0.5

h4 parameters: w41=-3.0, w42=1.0, w43=0, w44=-3.0, theta4=0.5

h5 parameters: w51=1.0, w52=1.0, w53=1.0, w54=1.0, theta5=0.5

X1 X2 X3 X4 | Target | NN Output | Match?

```

-----
0 0 0 0 | 0 | 0 | True
0 0 0 1 | 0 | 0 | True
0 0 1 0 | 0 | 0 | True
0 0 1 1 | 1 | 1 | True
0 1 0 0 | 1 | 1 | True
0 1 0 1 | 1 | 1 | True
0 1 1 0 | 1 | 1 | True
0 1 1 1 | 0 | 0 | True
1 0 0 0 | 0 | 0 | True
1 0 0 1 | 0 | 0 | True
1 0 1 0 | 0 | 0 | True
1 0 1 1 | 1 | 1 | True
1 1 0 0 | 0 | 0 | True
1 1 0 1 | 0 | 0 | True
1 1 1 0 | 0 | 0 | True
1 1 1 1 | 1 | 1 | True
-----

```

Total Network Accuracy: 100.0%

Question 2 Part c

In [143...

```

#Calculated new weight matrix for part c
W_partc = np.array([
    [ 1.0,  0.0,  1.0,  1.0],
    [ 0.0, -1.0,  1.0,  1.0],
    [-1.0,  1.0, -1.0,  0.0],
    [-1.0,  1.0,  0.0, -1.0]
])

#"inputs" was already created 16 combination of 4 inputs
#Then we generate 100 samples from these 16 items that is created earlier. repla

```

```

#after choosing one of the 16 samples, take it back to the pool so that it can be
#otherwise after choosing 16 samples, there weren't any sample to pick from.
base_inputs = np.array(list(itertools.product([0, 1], repeat=4)))
clean_samples = np.repeat(base_inputs, 25, axis=0)

#zero mean, std 0.2 gaussian noise generation
noise = np.random.normal(0, 0.2, size=(400, 4))
noisy_samples = clean_samples + noise

correct_a = 0
correct_c = 0

#part a and part c function calculation
for i in range(len(noisy_samples)):
    clean_X = clean_samples[i]
    noisy_X = noisy_samples[i]

    #F = (X1 OR NOT X2) XOR (NOT X3 OR NOT X4), != means XOR, first do it with n
    # noise will be added to it after
    target = int((clean_X[0] or not clean_X[1]) != (not clean_X[2] or not clean_X[3]))

    #part a
    h_a = step_function(np.dot(W_hidden, noisy_X) + b_hidden)
    out_a = int(step_function(np.dot(W_out, h_a) - t5))
    if out_a == target: correct_a += 1

    #part c
    h_c = step_function(np.dot(W_partc, noisy_X) + b_hidden)
    out_c = int(step_function(np.dot(W_out, h_c) - t5))
    if out_c == target: correct_c += 1
np.random.seed(45) #to get fixated results
#performance calculation
print(f"Part A Performance: {correct_a / 400.0 * 100}%")
print(f"Part C Performance: {correct_c / 400.0 * 100}%")

```

Part A Performance: 80.75%

Part C Performance: 92.25%

## Question 3 Part a Sample Visualisation

```
In [4]: import h5py
import numpy as np
import matplotlib.pyplot as plt
h5_filename = "assign1_data1.h5"
f = h5py.File(h5_filename, 'r')
for key in f.keys():
    shape = f[key].shape
    print(key, shape)

#take the trainims to plot and visualize them
img = f["trainims"][:]
img1 = np.squeeze(img)

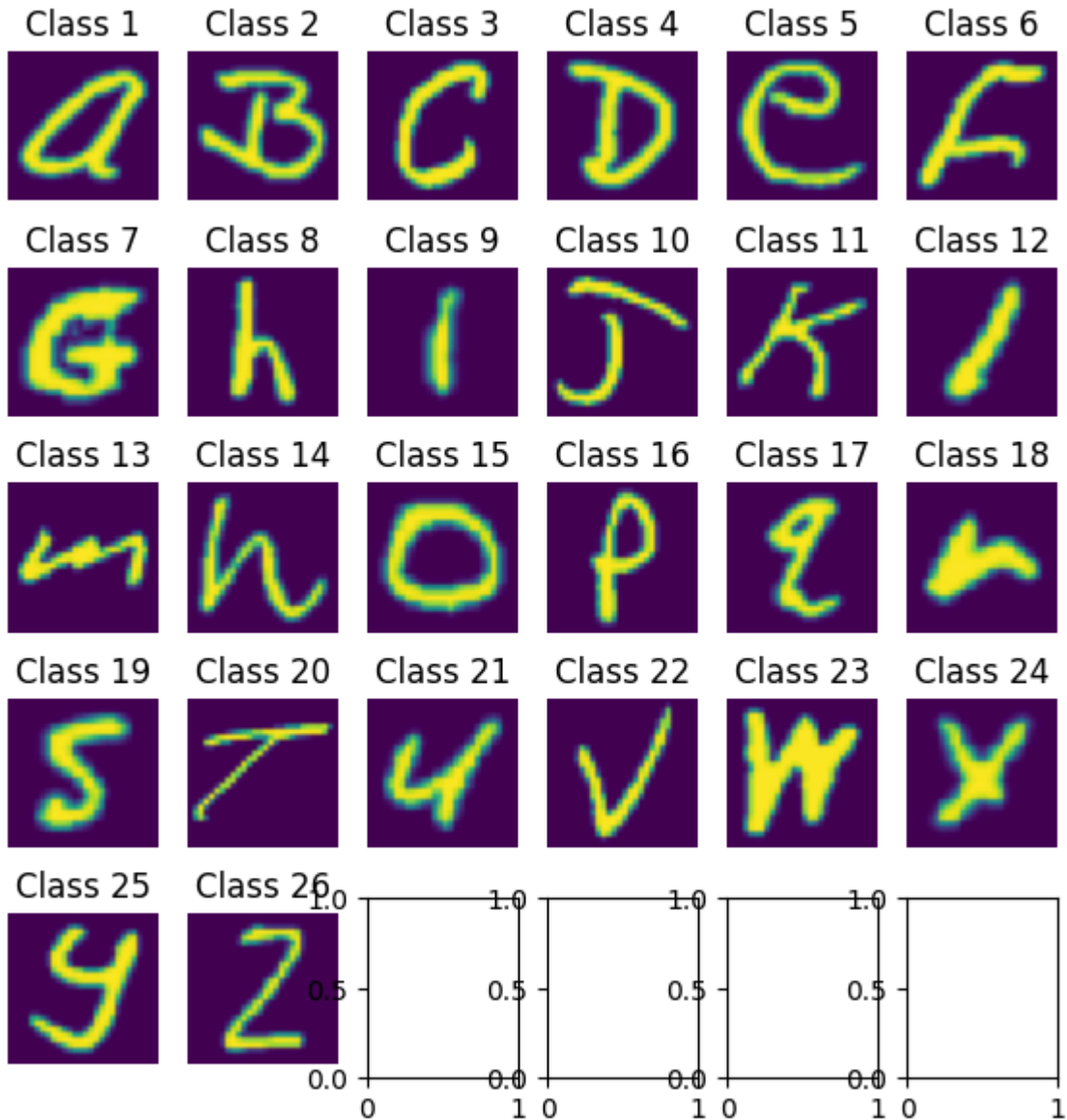
#see how many classes in there, so we get the total number
#distinct letters are there
lbls = np.squeeze(f["trainlbls"][:])
classes = np.unique(lbls) #we use unique to find the number of DISTINCT letters
print(len(classes))

#create subplots to order the visuals
fig, axes = plt.subplots(5, 6, figsize=(7, 7))
axes = axes.flatten()

#the iteration will be done for the number of distinct classes
for i in range(len(classes)):
    #to find the index of the beginning of each distinct letter, we need to find
    #for example the index where the samples of letter "A" changes to samples "B"
    correct_index = np.where(lbls == classes[i])[0][0]
    axes[i].imshow(img1[correct_index].T)
    axes[i].set_title(f"Class {int(classes[i])}")
    axes[i].axis('off')

plt.show()
```

```
testims (1300, 28, 28)
testlbls (1300,)
trainims (5200, 28, 28)
trainlbls (5200,)
26
```



Question 3 Part a Correlation Matrix

```
In [5]: import string
#Rather than indexing the change from letter "A" to letter "B", I know that there
#so, every 200 samples the letter will be changing.
selected_images = img1[0:5200:200]

#Make the selected_images 1D so I can calculate the correlation matrix
flattened_images = selected_images.reshape(26, -1)

#Correlation matrix
corr_matrix = np.corrcoef(flattened_images)

#Matrix plotting
fig, ax = plt.subplots(figsize=(14, 14))
im = ax.imshow(corr_matrix, vmin=-1, vmax=1)

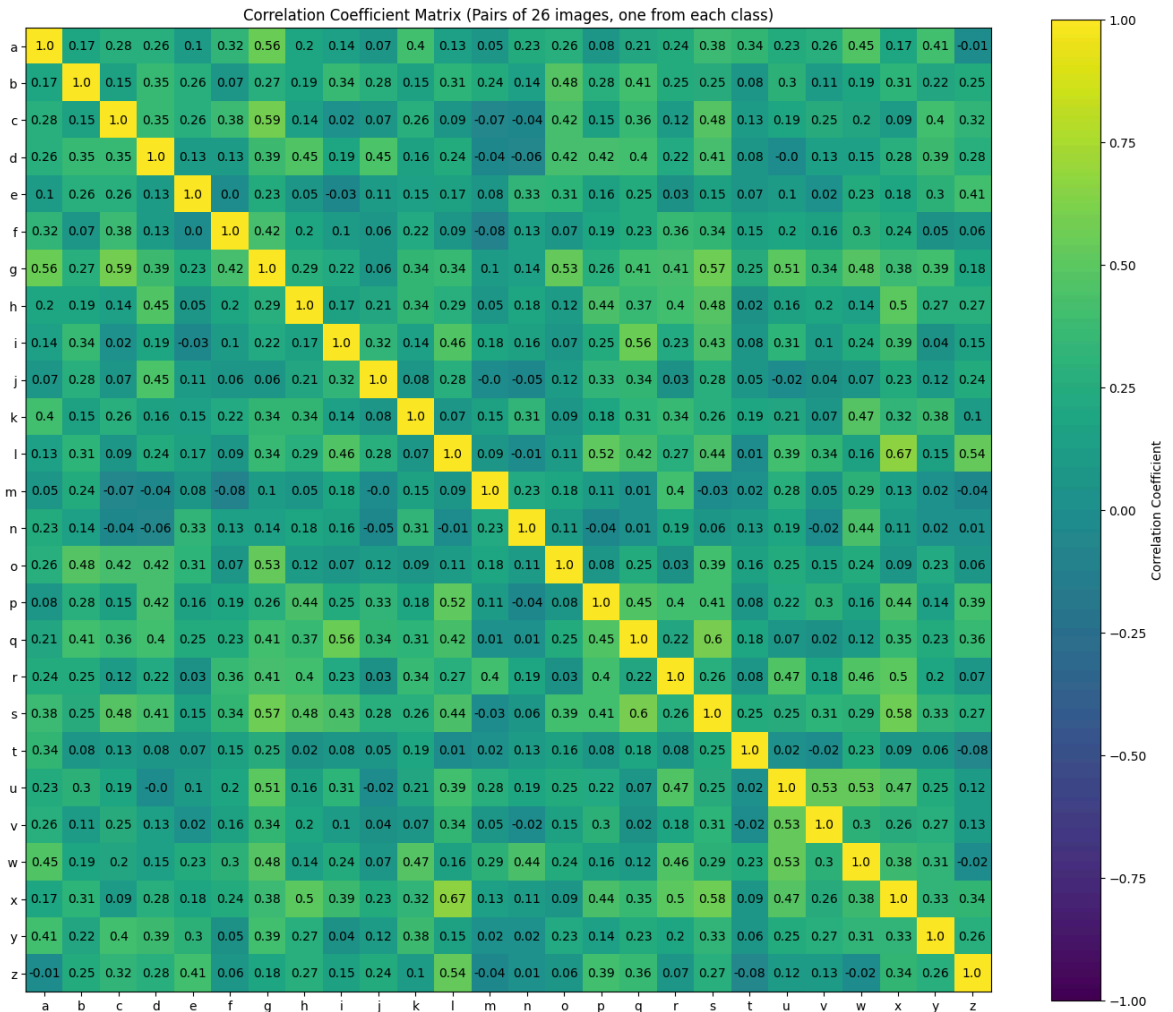
#showing the correlation in each box
for i in range(len(classes)):
    for j in range(len(classes)):
        val = float(np.round(corr_matrix[i, j], 2))
        ax.text(j, i, val, ha="center", va="center", color = "k")
```

```

#X and Y Labels
letters = list(string.ascii_lowercase)
ax.set_xticks(np.arange(len(classes)))
ax.set_yticks(np.arange(len(classes)))
ax.set_xticklabels(letters)
ax.set_yticklabels(letters)
ax.set_title("Correlation Coefficient Matrix (Pairs of 26 images, one from each

fig.colorbar(im, ax=ax, shrink=0.8, label='Correlation Coefficient')
plt.tight_layout()
plt.show()

```



Question 3 Part b

```

In [6]: import numpy as np
import matplotlib.pyplot as plt
import string

trainims = img1
trainbls_raw = lbls.astype(int)

#find the number of distinc classes and make them sort from 0 to 25
classes = np.unique(trainbls_raw)
C = len(classes)
y_idx = np.searchsorted(classes, trainbls_raw) # -> 0..C-1

#shape (5200, 28, 28)
N, H, W = trainims.shape

```

```

#make it 1D array
X_train = trainims.reshape(N, -1).astype(np.float32)

#reshape
if X_train.max() > 1.5:
    X_train /= 255.0 # Normalize 8-bit pixel intensities (0-255) to a continuous

D = X_train.shape[1]
Y_train = np.eye(C, dtype=np.float32)[y_idx]

def sigmoid(z):
    return 1.0 / (1.0 + np.exp(-z))

#Trains a single-layer perceptron using online Stochastic Gradient Descent (SGD)
#Uses a compound matrix to update all 26 class connections simultaneously.

def train_one_layer(eta, iterations=10000, seed=42):
    rng = np.random.default_rng(seed)

    #Initialize compound weight matrix W and bias vector by using a Gaussian dis
    W = rng.normal(0.0, 0.01, size=(D, C)).astype(np.float32)
    b = rng.normal(0.0, 0.01, size=(1, C)).astype(np.float32)
    mse_hist = np.zeros(iterations, dtype=np.float32)

    for it in range(iterations):
        # Pick exactly ONE random training sample (Stochastic/Online update)
        idx = rng.integers(0, N)
        x = X_train[idx:idx+1]
        y = Y_train[idx:idx+1]

        #Forward Pass
        yhat = sigmoid(x @ W + b)
        #error calculation
        err = y - yhat
        mse_hist[it] = np.mean(err**2)

        #Backward Pass
        delta = err * (yhat * (1.0 - yhat))

        #update the weights in each iteration
        W += eta * (x.T @ delta)
        b += eta * delta

    return W, b, mse_hist

#Finding the best eta value
eta_candidates = [0.001, 0.01, 0.1, 0.12, 0.13]
best_eta, best_final = None, np.inf

for eta in eta_candidates:
    W_tmp, b_tmp, mse_tmp = train_one_layer(eta, iterations=10000, seed=42)
    final_mse = float(mse_tmp[-1])
    print(f"eta={eta:<6} final MSE={final_mse:.6f}")
    if final_mse < best_final:
        best_final = final_mse
        best_eta = eta
    W_best, b_best, mse_best = W_tmp, b_tmp, mse_tmp

print("\nChosen eta* =", best_eta)

```

```

#MSE plot
plt.figure(figsize=(7,4))
plt.plot(mse_best)
plt.xlabel("Iteration")
plt.ylabel("MSE")
plt.title(f"Training MSE (eta*={best_eta})")
plt.grid(True, alpha=0.3)
plt.tight_layout()
plt.show()

#Showing the final network weights
letters = list(string.ascii_lowercase)[:C]

fig, axes = plt.subplots(5, 6, figsize=(8, 7))
axes = axes.flatten()

for i in range(C):
    wimg = W_best[:, i].reshape(H, W)
    axes[i].imshow(wimg.T, cmap="plasma")
    axes[i].set_title(letters[i])
    axes[i].axis("off")

for k in range(C, len(axes)):
    axes[k].axis("off")

plt.suptitle("Final learned weights (one image per class)")
plt.tight_layout()
plt.show()

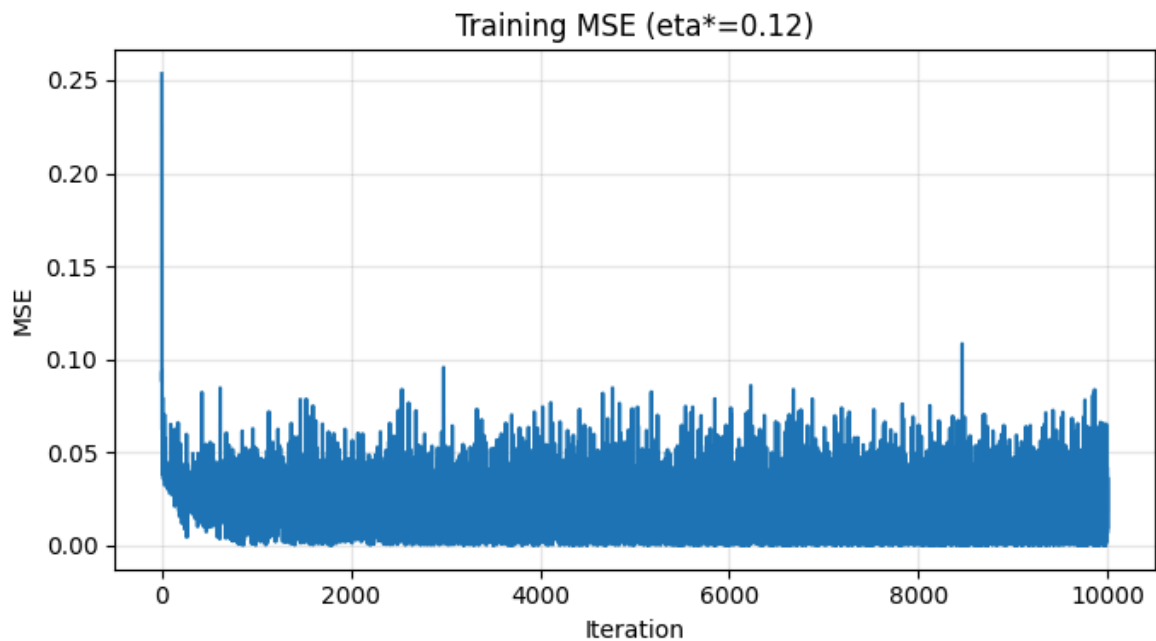
```

```

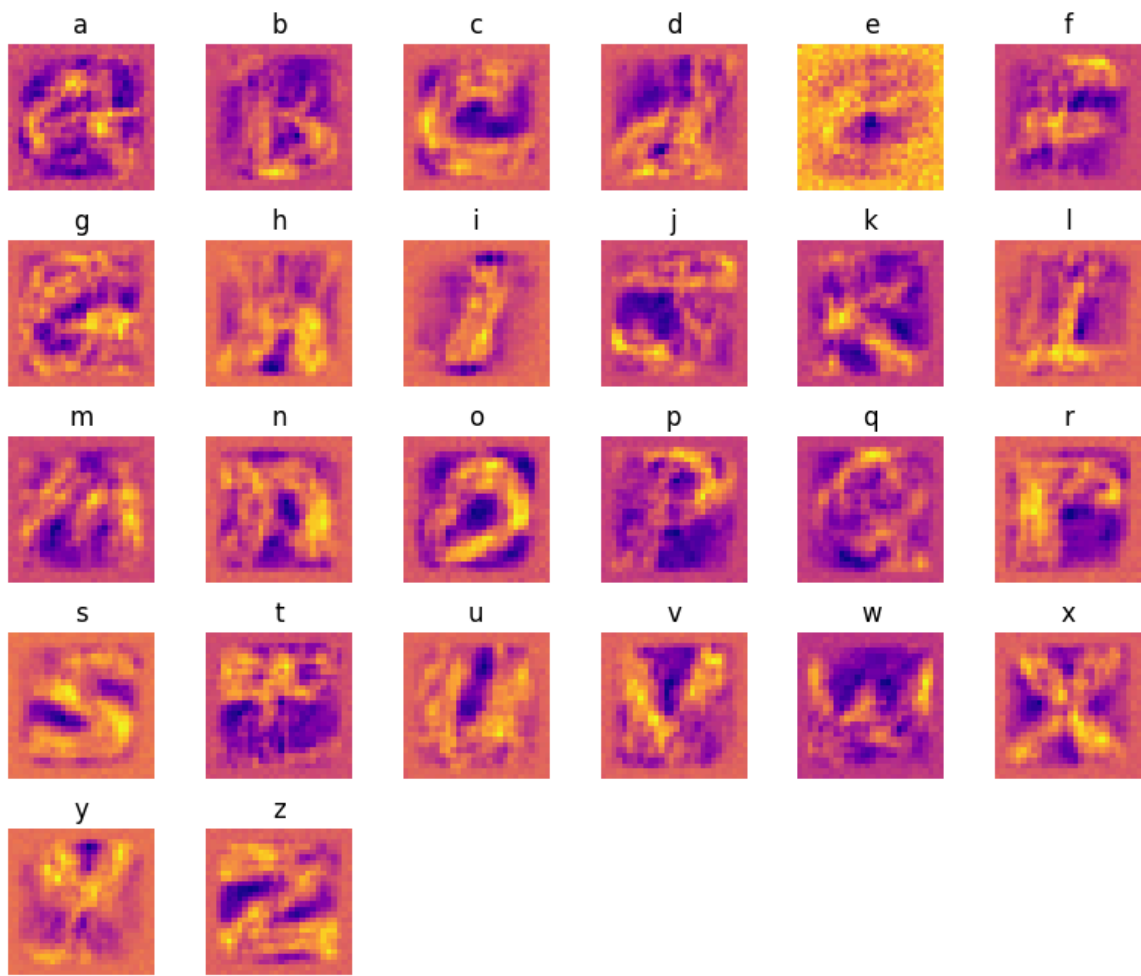
eta=0.001  final MSE=0.034964
eta=0.01   final MSE=0.014692
eta=0.1    final MSE=0.010058
eta=0.12   final MSE=0.009721
eta=0.13   final MSE=0.009729

```

Chosen eta\* = 0.12



Final learned weights (one image per class)



## Question 3 Part c

```
In [9]: #eta high, eta Low, and normal eta values
eta_star = best_eta
eta_low = eta_star / 1000
eta_high = eta_star * 1000

print("eta_low =", eta_low, "eta* =", eta_star, "eta_high =", eta_high)

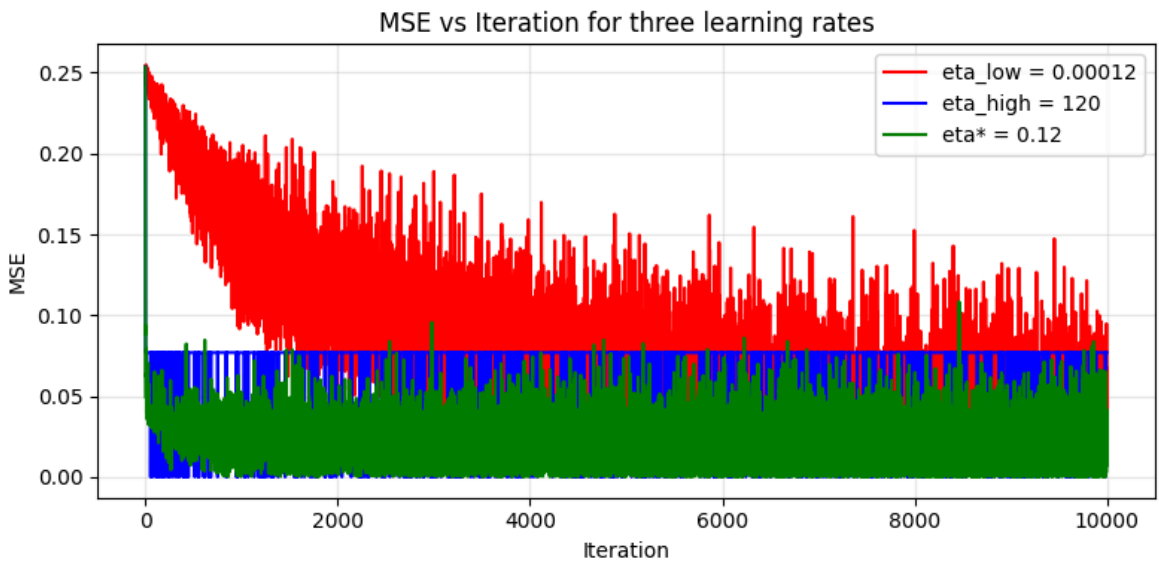
#do the same things for new eta values
W_low, b_low, mse_low = train_one_layer(eta_low, iterations=10000, seed=42)
W_high, b_high, mse_high = train_one_layer(eta_high, iterations=10000, seed=42)

#plot the MSE curve to analyze
plt.figure(figsize=(8,4))
plt.plot(mse_low, label=f"eta_low = {eta_low:g}", color = "r")
plt.plot(mse_high, label=f"eta_high = {eta_high:g}", color = "b")
plt.plot(mse_best, label=f"eta* = {eta_star:g}", color = "g")

plt.xlabel("Iteration")
plt.ylabel("MSE")
plt.title("MSE vs Iteration for three learning rates")
plt.grid(True, alpha=0.3)
plt.legend()
plt.tight_layout()
plt.show()
```

```
eta_low = 0.00011999999999999999 eta* = 0.12 eta_high = 120.0
```

```
C:\Users\DELL\AppData\Local\Temp\ipykernel_2860\554917531.py:21: RuntimeWarning:
overflow encountered in exp
return 1.0 / (1.0 + np.exp(-z))
```



### Question 3 Part d

```
In [8]: import h5py
import numpy as np
#getting the testims and testlbls from the h5 file to test their performance
with h5py.File(h5_filename, "r") as h5f:
    testims = np.squeeze(h5f["testims"][:]) #has a size (130
    testlbls_raw = np.squeeze(h5f["testlbls"][:]).astype(int) #

#in train array, we have 5200 samples, in test array we have 1300 samples!
Nt = testims.shape[0] #getting number of samples in test array, which is 1300

#reshape to meet dimensions
X_test = testims.reshape(Nt, -1).astype(np.float32)
if X_test.max() > 1.5:
    X_test /= 255.0 # Normalize 8-bit pixel intensities (0-255) to a continuous

test_y_idx = np.searchsorted(classes, testlbls_raw)

#sigmoid function
def sigmoid(z):
    return 1.0 / (1.0 + np.exp(-z))

def predict_idx(X, W, b): #^y = W*x + b
    yhat = sigmoid(X @ W + b) #shape of ^y = (1300, 26)
    return np.argmax(yhat, axis=1)

def accuracy_percent(y_true, y_pred):
    return 100.0 * np.mean(y_true == y_pred)

pred_low = predict_idx(X_test, W_low, b_low)
pred_star = predict_idx(X_test, W_best, b_best)
pred_high = predict_idx(X_test, W_high, b_high)

acc_low = accuracy_percent(test_y_idx, pred_low)
acc_star = accuracy_percent(test_y_idx, pred_star)
```

```
acc_high = accuracy_percent(test_y_idx, pred_high)

print(f"Test performance (eta_low) : {acc_low:.2f}%")
print(f"Test performance (eta*)   : {acc_star:.2f}%")
print(f"Test performance (eta_high): {acc_high:.2f}%")
```

```
Test performance (eta_low) : 20.15%
Test performance (eta*)   : 55.77%
Test performance (eta_high): 3.85%
```

```
C:\Users\DELL\AppData\Local\Temp\ipykernel_2860\554917531.py:21: RuntimeWarning:
overflow encountered in exp
  return 1.0 / (1.0 + np.exp(-z))
```

# Implementing a Neural Network

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

```
In [1]: # A bit of setup
from __future__ import print_function
import numpy as np
import matplotlib.pyplot as plt

from cs231n.classifiers.neural_net import TwoLayerNet

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

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

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

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

```
In [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_inputs = 5

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

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

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

## Forward pass: compute scores

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

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

```
In [3]: scores = net.loss(X)
print('Your scores:')
print(scores)
print()
print('correct scores:')
correct_scores = np.asarray([
    [-0.81233741, -1.27654624, -0.70335995],
    [-0.17129677, -1.18803311, -0.47310444],
    [-0.51590475, -1.01354314, -0.8504215 ],
    [-0.15419291, -0.48629638, -0.52901952],
    [-0.00618733, -0.12435261, -0.15226949]])
print(correct_scores)
print()

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

Your scores:

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

correct scores:

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

Difference between your scores and correct scores:

```
3.6802720745909845e-08
```

## Forward pass: compute loss

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

```
In [4]: loss, _ = net.loss(X, y, reg=0.05)
correct_loss = 1.30378789133

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

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

## Backward pass

Implement the rest of the function. This will compute the gradient of the loss with respect to the variables `W1`, `b1`, `W2`, and `b2`. Now that you (hopefully!) have a correctly implemented forward pass, you can debug your backward pass using a numeric gradient check:

```
In [5]: from cs231n.gradient_check import eval_numerical_gradient

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

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

# 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)[0]
    param_grad_num = eval_numerical_gradient(f, net.params[param_name], verbose=False)
    print('%s max relative error: %e' % (param_name, rel_error(param_grad_num, g
```

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

## Train the network

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

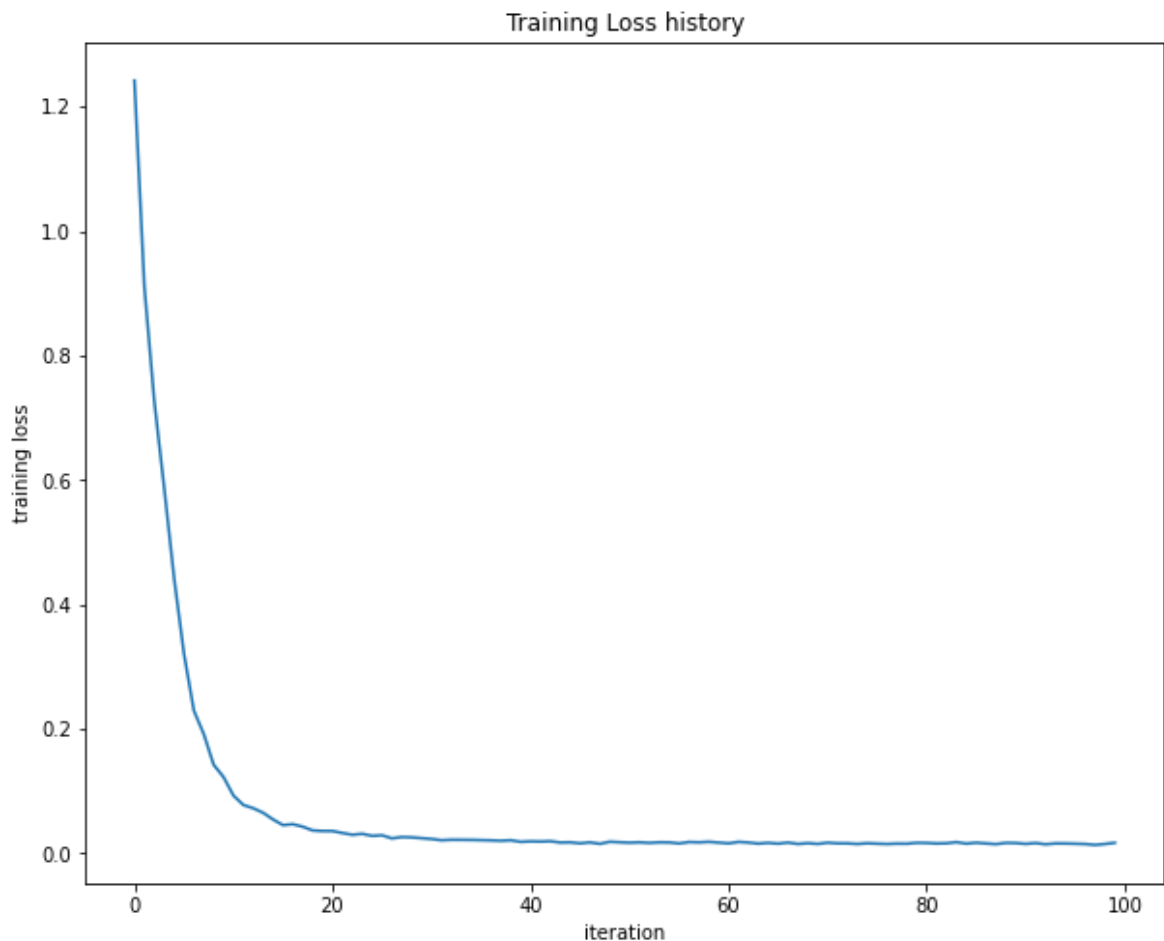
Once you have implemented the method, run the code below to train a two-layer network on toy data. You should achieve a training loss less than 0.2.

```
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: 0.017149607938732093



## 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 cs231n.data_utils import load_CIFAR10

def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000):
    """
    Load the CIFAR-10 dataset from disk and perform preprocessing to prepare
    it for the two-layer neural net classifier. These are the same steps as
    we used for the SVM, but condensed to a single function.
    """
    # Load the raw CIFAR-10 data
    cifar10_dir = 'cs231n/datasets/cifar-10-batches-py'
```

```

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

# Subsample the data
mask = 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

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

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

```

```

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

```

## Train a network

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

```
In [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 2.302954
iteration 100 / 1000: loss 2.302550
iteration 200 / 1000: loss 2.297648
iteration 300 / 1000: loss 2.259602
iteration 400 / 1000: loss 2.204170
iteration 500 / 1000: loss 2.118565
iteration 600 / 1000: loss 2.051535
iteration 700 / 1000: loss 1.988466
iteration 800 / 1000: loss 2.006591
iteration 900 / 1000: loss 1.951473
Validation accuracy: 0.287
```

## Debug the training

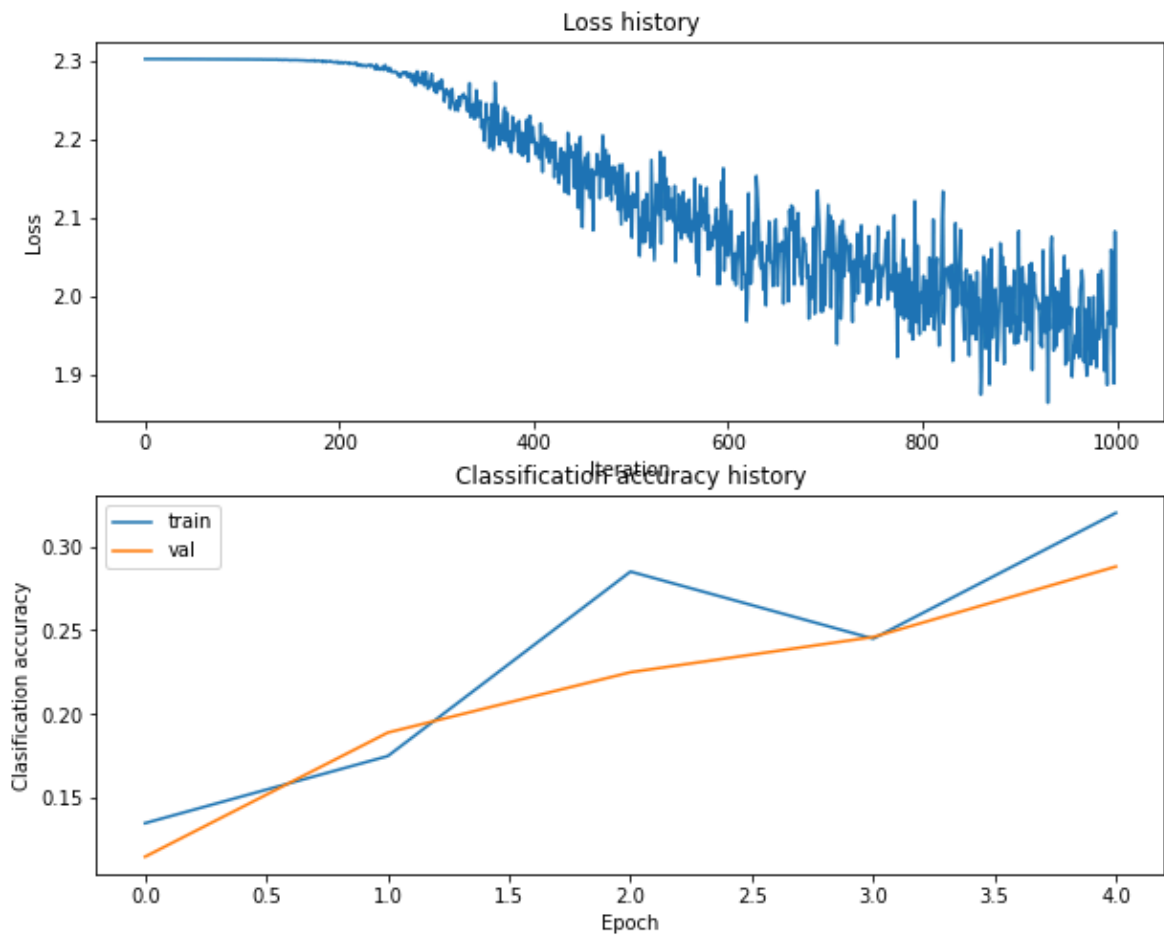
With the default parameters we provided above, you should get a validation accuracy of about 0.29 on the validation set. This isn't very good.

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.legend()
plt.show()
```

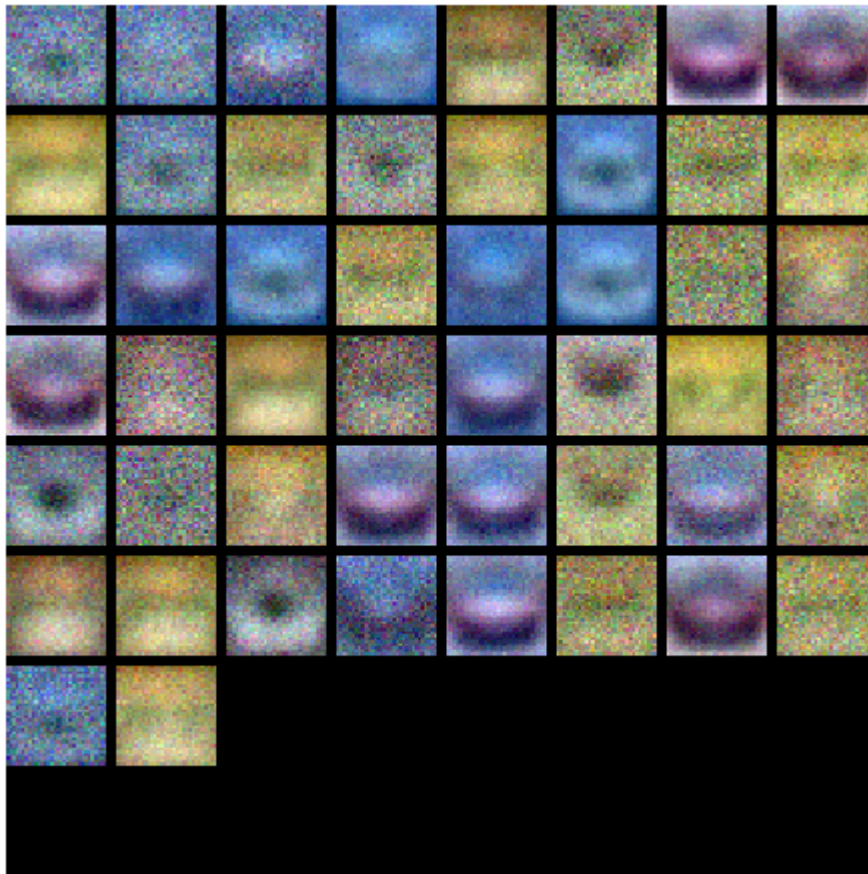


```
In [10]: from cs231n.vis_utils import visualize_grid

# Visualize the weights of the network

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

show_net_weights(net)
```



## Tune your hyperparameters

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

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

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

**Experiment:** Your goal in this exercise is to get as good of a result on CIFAR-10 as you can, with a fully-connected Neural Network. Feel free to implement your own techniques (e.g. PCA to reduce dimensionality, or adding dropout, or adding features to the solver, etc.).

```

In [11]: best_net = None # store the best model into this

#####
# TODO: Tune hyperparameters using the validation set. Store your best trained
# model in best_net.
#
# To help debug your network, it may help to use visualizations similar to the
# ones we used above; these visualizations will have significant qualitative
# differences from the ones we saw above for the poorly tuned network.
#
# Tweaking hyperparameters by hand can be fun, but you might find it useful to
# write code to sweep through possible combinations of hyperparameters
# automatically like we did on the previous exercises.
#####

input_size = X_train.shape[1]
hidden_size = 100
output_size = 10

# Learning_rates = [1, 1e-1, 1e-2, 1e-3]
# regularization_strengths = [1e-6, 1e-5, 1e-4, 1e-3]

# Magic lrs and regs?
# Just copy from: https://github.com/Lightaime/cs231n/blob/master/assignment1/two\_layer\_net.py
# :(
learning_rates = np.array([0.7, 0.8, 0.9, 1, 1.1])*1e-3
regularization_strengths = [0.75, 1, 1.25]

best_val = -1

for lr in learning_rates:
    for reg in regularization_strengths:
        net = TwoLayerNet(input_size, hidden_size, output_size)
        net.train(X_train, y_train, X_val, y_val, learning_rate=lr, reg=reg,
                  num_iters=1500)

        y_val_pred = net.predict(X_val)
        val_acc = np.mean(y_val_pred == y_val)

        print('lr: %f, reg: %f, val_acc: %f' % (lr, reg, val_acc))

        if val_acc > best_val:
            best_val = val_acc
            best_net = net

print('Best validation accuracy: %f' % best_val)
#####
#                                     END OF YOUR CODE
#####

```

```
lr: 0.000700, reg: 0.750000, val_acc: 0.477000
lr: 0.000700, reg: 1.000000, val_acc: 0.477000
lr: 0.000700, reg: 1.250000, val_acc: 0.478000
lr: 0.000800, reg: 0.750000, val_acc: 0.471000
lr: 0.000800, reg: 1.000000, val_acc: 0.464000
lr: 0.000800, reg: 1.250000, val_acc: 0.470000
lr: 0.000900, reg: 0.750000, val_acc: 0.489000
lr: 0.000900, reg: 1.000000, val_acc: 0.474000
lr: 0.000900, reg: 1.250000, val_acc: 0.473000
lr: 0.001000, reg: 0.750000, val_acc: 0.479000
lr: 0.001000, reg: 1.000000, val_acc: 0.477000
lr: 0.001000, reg: 1.250000, val_acc: 0.467000
lr: 0.001100, reg: 0.750000, val_acc: 0.474000
lr: 0.001100, reg: 1.000000, val_acc: 0.491000
lr: 0.001100, reg: 1.250000, val_acc: 0.477000
Best validation accuracy: 0.491000
```

```
In [12]: # visualize the weights of the best network
show_net_weights(best_net)
```



## Run on the test set

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

```
In [16]: test_acc = (best_net.predict(X_test) == y_test).mean()
print('Test accuracy: ', test_acc)
```

```
Test accuracy: 0.482
```

### Inline Question

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

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

Your answer: **1 and 3**

Your explanation: This behavior of the large gap between train and test accuracy is called as "overfitting". What happens is basically the algorithm focuses more on noise and specific features of the train data rather than necessary generalized patterns, therefore it MEMORIZES the train set. Thus, when the test data fed into the algorithm, it performs bad due to memorization. To overcome this problem, i suggest using larger datasets on training and increasing the regularization strength.

### WHY NOT 2?

**2.** Before explaining the first and third options, I want to explain how will adding more hidden units affect the gap between train and test accuracy. Restricting the network will hinder the ability to memorize train dataset of the algorithm. If I would have used more hidden layers, our model will memorize the train dataset more, which will increase the gap between test and train accuracy. However, if I add controllable hidden layers, it may result in better performance and can decrease the gap between train and test accuracies, if the network was UNDERFITTED [5].

**WHY 1 & 3? 1.** Increasing the number of training examples is one of the most reliable ways to reduce overfitting. When a model is trained on a limited dataset, it tends to memorize the specific noise, quirks, and exact pixel arrangements of those few images. Because this limited set does not properly represent the real-world variations of the data, the model fits too closely to it and becomes ungeneralizable to new, unseen datasets. By providing a much larger and more diverse dataset, it becomes statistically impossible for the network to simply memorize every single sample. Instead, the model is forced to generalize and learn the true, underlying structural patterns of the classes. As the amount of data increases, the model generalizes better because it fits a smaller fraction of the specific training samples and focuses on broader trends[5][6].

**3.** Increasing the regularization strength is a fundamental method used to constrain the complexity of a machine learning model. It works by acting as a mathematical penalty during the learning process, specifically targeting and penalizing features or weights that have little to no true predictive contribution. When a network overfits and memorizes noise, it often does so by assigning massive weight values to highly specific artifacts in the training data. Regularization limits this model complexity by forcing the network to keep its weights small and distributed [5][6].

### Overall comments on the code

The custom implementation of the two-layer neural network was successfully verified using a numeric gradient check. The relative errors for the forward pass, loss calculation, and analytical gradients were all well below the required thresholds.

The initial default parameters yielded a poor validation accuracy of 28.7%, which has to be around 29%. To improve this, a hyperparameter grid search was conducted. The model was trained across a range of learning rates ( $0.71e-3$  to  $1.11e-3$ ) and regularization strengths (0.75 to 1.25). The optimal hyperparameter combination was found to be a learning rate of 0.0011 and a regularization strength of 0.75, which boosted the validation accuracy to 49.1%.

The optimized two-layer network was evaluated on the unseen test set, achieving a final test accuracy of 48.2%. This exceeds the assignment's baseline requirement of 48%. The small margin between the validation accuracy (49.1%) and the test accuracy (48.2%) indicates that the chosen hyperparameters are robust and that the model did not artificially overfit to the validation data.