

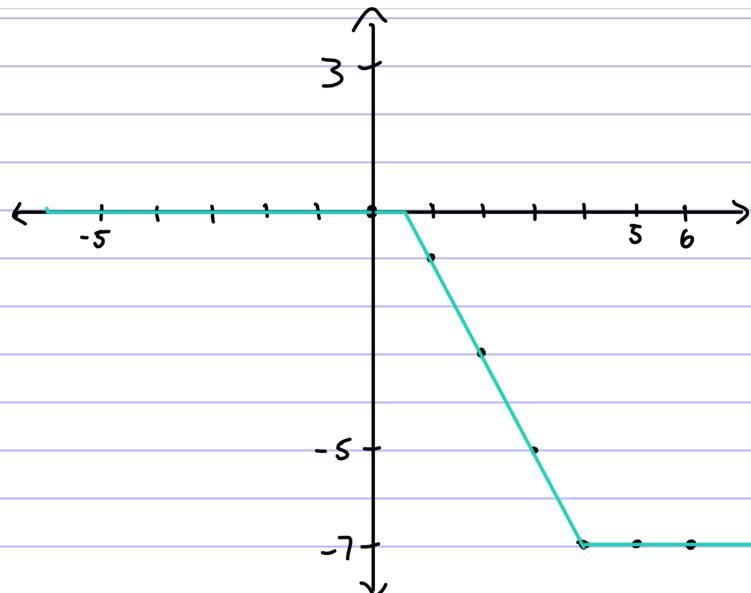
Assignment 10

1. Neural net functions

- a) Sketch the function generated by the following 3-neuron ReLU neural network.

$$f(x) = 2(x - 0.5)_+ - 2(2x - 1)_+ + 4(0.5x - 2)_+$$

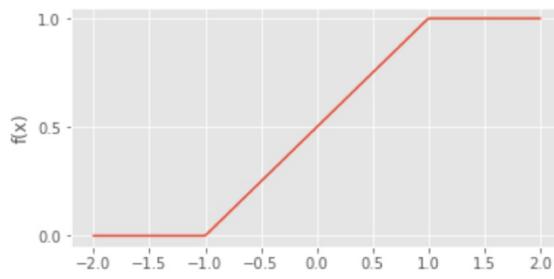
where $x \in \mathbb{R}$ and where $(z)_+ = \max(0, z)$ for any $z \in \mathbb{R}$. Note that this is a single-input, single-output function. Plot $f(x)$ vs x by hand.



- b) Consider the continuous function depicted below. Approximate this function with ReLU neural network with 2 neurons. The function should be in the form

$$f(x) = \sum_{j=1}^2 v_j(w_j x + b_j)_+$$

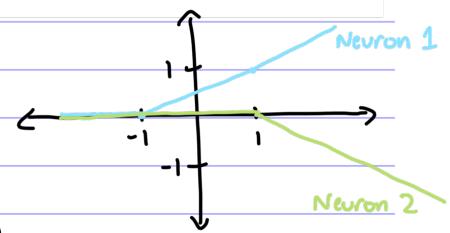
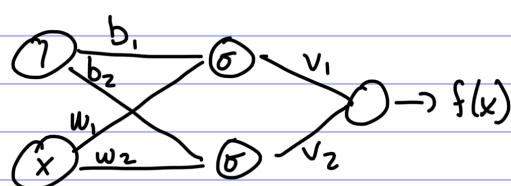
Indicate the weights and biases of each neuron and sketch the neural network function.



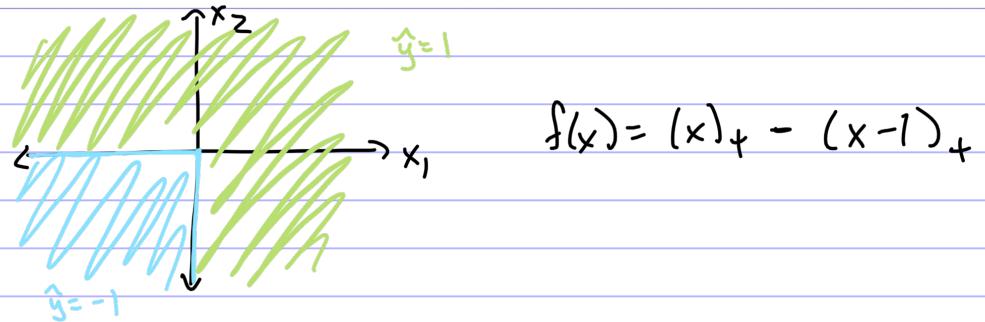
$$f(x) = 0.5(x+1)_+ - 0.5(x-1)_+$$

$$\begin{aligned} v_1 &= 0.5 \\ w_1 &= 1 \\ b_1 &= 1 \end{aligned}$$

$$\begin{aligned} v_2 &= -0.5 \\ w_2 &= 1 \\ b_2 &= -1 \end{aligned}$$



- c) A neural network f_w can be used for binary classification by predicting the label as $\hat{y} = \text{sign}(f_w(\mathbf{x}))$. Consider a setting where $\mathbf{x} \in \mathbb{R}^2$ and the desired classifier is -1 if both elements of \mathbf{x} are less than or equal to zero and $+1$ otherwise. Sketch the desired classification regions in the two-dimensional plane, and provide a formula for a ReLU network with 2-neurons that can produce the desired classification. For simplicity, assume in this question that $\text{sign}(0) = -1$.



2. **Gradients of a neural net.** Consider a 2 layer neural network of the form $f(\mathbf{x}) = \sum_{j=1}^J v_j (\mathbf{w}_j^T \mathbf{x})_+$. Suppose we want to train our network on a dataset of N samples \mathbf{x}_i with corresponding labels y_i , using a least squares loss function $\mathcal{L} = \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2$. Derive the gradient descent update steps for the input weights \mathbf{w}_j and output weights v_j .

1. Initialize w_j^0 and v_j^0 randomly

2. Iterate:

Forward: traverse through the current neural network to get $f(\mathbf{x})$

Backward:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial v_j} &= \frac{\partial f}{\partial x_i} \frac{\partial x_i}{\partial v} \\ &= 2(f(x_i) - y) \cdot (w_j^T x_i)_+\end{aligned}$$

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial w_j} &= \frac{\partial f}{\partial x_i} \frac{\partial x_i}{\partial w} \\ &= 2(f(x_i) - y) \cdot v_j \sigma'(w_j^T x_i) \cdot x_i\end{aligned}$$

$$v^{t+1} = v^t - \alpha_t (2(f(x_i) - y) \cdot (w_j^T x_i)_+)$$

$$w^{t+1} = w^t - \alpha_t (2(f(x_i) - y) \cdot v_j \sigma'(w_j^T x_i) \cdot x_i)$$

3. **Compressing neural nets.** Large neural network models can be approximated by considering low rank approximations to weight matrices. The neural network $f(\mathbf{x}) = \sum_{j=1}^J \mathbf{v}_j (\mathbf{w}_j^T \mathbf{x})_+$ can be written as

$$f(\mathbf{x}) = \mathbf{v}^T (\mathbf{W} \mathbf{x})_+.$$

where \mathbf{v} is a $J \times 1$ vector of the output weights and \mathbf{W} is a $J \times d$ matrix with i th row \mathbf{w}_j^T . Let $\sigma_1, \sigma_2, \dots$ denote the singular values of \mathbf{W} and assume that $\sigma_i \leq \epsilon$ for $i > r$. Let f_r denote the neural network obtained by replacing \mathbf{W} with its best rank r approximation $\hat{\mathbf{W}}_r$. Assuming that \mathbf{x} has unit norm, find an upper bound to the difference $\max_x |f(\mathbf{x}) - f_r(\mathbf{x})|$. (Hint: for any pair of vectors \mathbf{a} and \mathbf{b} , the following inequality holds $\|\mathbf{a}_+ - \mathbf{b}_+\|_2 \leq \|\mathbf{a} - \mathbf{b}\|_2$).

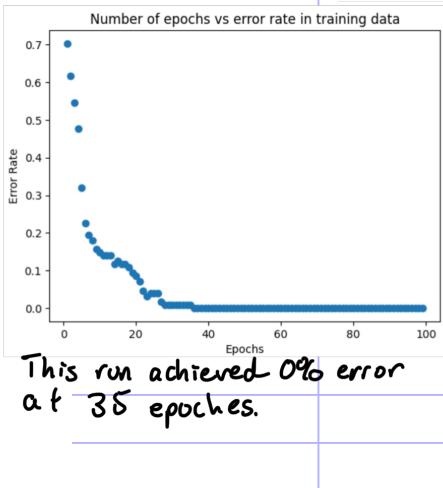
Upper bound would occur when $\hat{\mathbf{W}}_r$ is a rank 1 approximation. Rank 1 is the farthest approximation from full rank (\mathbf{W}), so the weights would be the farthest from the actual weights and thus $f_r(\mathbf{x})$ would have the largest difference from the actual $f(\mathbf{x})$.

4. **Face Emotion Classification with a three layer neural network.** In this problem we return to the face emotion data studied previously. You may find it very helpful to use code from an activity (or libraries such as Keras and Tensorflow).

- a) Build a classifier using a full connected three layer neural network with logistic activation functions. Your network should
- take a vector $\mathbf{x} \in \mathbb{R}^{10}$ as input (nine features plus a constant offset),
 - have a single, fully connected hidden layer with 32 neurons
 - output a scalar \hat{y} .

Note that since the logistic activation function is always positive, your decision should be as follows: $\hat{y} > 0.5$ corresponds to a ‘happy’ face, while $\hat{y} \leq 0.5$ is not happy.

- b) Train your classifier using stochastic gradient descent (start with a step size of $\alpha = 0.05$) and create a plot with the number of epochs on the horizontal axis, and training accuracy on the vertical axis. Does your classifier achieve 0% training error? If so, how many epoch does it take for your classifier to achieve perfect classification on the training set?

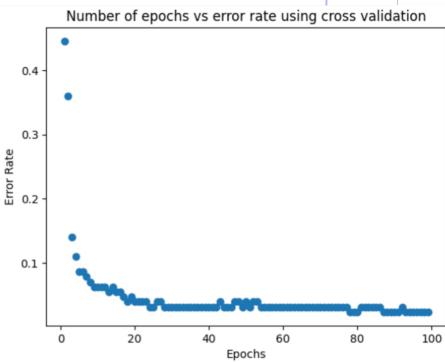


The classifier does achieve 0% training error. Since the initial weights are initialized randomly and then the same initial weights are used for each epoch amount, the number of epochs needed to reach 0% error varies between runs.

Based on my runs, I've seen 0% error achieved anywhere between 30 and 80 epochs.

- c) Find a more realistic estimate of the accuracy of your classifier by using 8-fold cross validation. Can you achieve perfect test accuracy?

We do not achieve 0% test error. However we can get very close.
With the same random initial weights but different epoch amounts, the test error converges after around 20 epochs.



```

import numpy as np
import matplotlib.pyplot as plt
from scipy.io import loadmat

dataset = loadmat('face_emotion_data.mat')

X, y = dataset['X'], dataset['y']
n, p = np.shape(X)

y[y== -1] = 0 # use 0/1 for labels instead of -1/+1
X = np.hstack((np.ones((n,1)), X)) # append a column of ones

```

Neural network definition

```

def logsig(_x):
    return 1/(1+np.exp(-_x))

def neural_net(L, X, y, V, W):
    alpha = 0.05 #step size

    for epoch in range(L):
        ind = np.random.permutation(np.shape(X)[0])
        for i in ind:
            # Forward-propagate
            H = logsig(np.hstack((np.ones((1,1)), X[[i],:]@W)))
            Yhat = logsig(H@V)
            # Backpropagate
            delta = (Yhat-y[[i],:])*Yhat*(1-Yhat)
            Vnew = V-alpha*H.T@delta
            gamma = delta@V[1:,:].T*H[:,1:]* (1-H[:,1:])
            Wnew = W - alpha*X[[i],:].T@gamma
            V = Vnew
            W = Wnew
            #print('epoch: ', epoch)

    return W, V

```

Varied epoch numbers

```

error_rates = []

q = np.shape(y)[1] #number of classification problems
M = 32 #number of hidden nodes

## initial weights
V_init = np.random.randn(M+1, q);
W_init = np.random.randn(p+1, M);

for e in range(1, 100):
    W, V = neural_net(e, X, y, V_init, W_init)

    ## Final predicted labels (on training data)
    H = logsig(np.hstack((np.ones((n,1)), X@W)))
    Yhat = logsig(H@V)
    Yhat = np.array([[1 if i > 0.5 else 0 for i in Yhat]]).transpose()

    # error calculation
    errors = [0 if i[0]==i[1] else 1 for i in np.hstack((y, Yhat))]
    error_rates.append(sum(errors)/np.shape(y)[0])

```

```

print(list(error_rates).index(0.0))

```

35

```

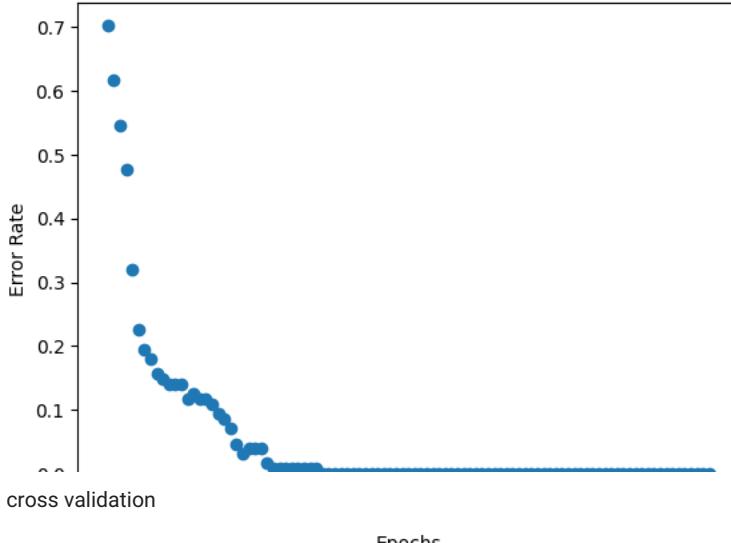
plt.scatter(range(1,100), error_rates)

plt.title("Number of epochs vs error rate in training data")
plt.xlabel("Epochs")
plt.ylabel("Error Rate")

```

```
Text(0, 0.5, 'Error Rate')
```

Number of epochs vs error rate in training data



```
8-fold cross validation
```

```
error_rate = []  
  
q = np.shape(y)[1] #number of classification problems  
M = 32 #number of hidden nodes  
  
## initial weights  
V_init = np.random.randn(M+1, q);  
W_init = np.random.randn(p+1, M);  
  
# for each epoch size  
for e in range(1, 100):  
    eight_sets = np.split(X, 8)  
    eight_y = np.split(y, 8)  
    total_error = 0  
  
    # cross validation  
    for i in range(8):  
        # remove one set from the eight  
        training = np.delete(eight_sets, i, axis=0).reshape(-1, 10)  
        training_y = np.delete(eight_y, i, axis=0).reshape(112, -1)  
  
        W, V = neural_net(e, training, training_y, V_init, W_init)  
  
        ## Final predicted labels (on training data)  
        H = logsig(np.hstack((np.ones((16,1)), eight_sets[i]@W)))  
        Yhat = logsig(H@V)  
  
        Yhat = np.array([[1 if i > 0.5 else 0 for i in Yhat]]).transpose()  
  
        errors = [0 if i[0]==i[1] else 1 for i in np.hstack((eight_y[i], Yhat))]  
        total_error += (sum(errors)/np.shape(eight_y[i])[0])  
  
    error_rate.append(total_error/8) # average error of the eight holdout sets  
  
plt.scatter(range(1,100), error_rate)
```

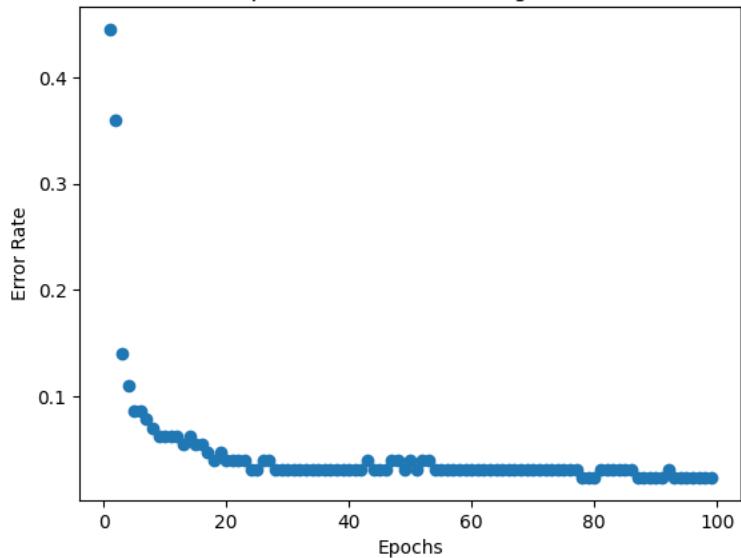
```
plt.title("Number of epochs vs error rate using cross validation")
```

```
plt.xlabel("Epochs")
```

```
plt.ylabel("Error Rate")
```

Text(0, 0.5, 'Error Rate')

Number of epochs vs error rate using cross validation



Colab paid products - Cancel contracts here

✓ 0s completed at 6:29 PM

