

Fall 2021 CS 4641\7641 A: Machine Learning Homework 4 ¶

Instructor: Dr. Mahdi Roozbahani

Deadline: December 2, Thursday, AOE

- No unapproved extension of the deadline is allowed. Late submission will lead to 0 credit.
- Discussion is encouraged on Ed as part of the Q/A. However, all assignments should be done individually.
- Plagiarism is a **serious offense**. You are responsible for completing your own work. You are not allowed to copy and paste, or paraphrase, or submit materials created or published by others, as if you created the materials. All materials submitted must be your own.
- All incidents of suspected dishonesty, plagiarism, or violations of the Georgia Tech Honor Code will be subject to the institute's Academic Integrity procedures (e.g., reported to and directly handled by the Office of Student Integrity (OSI)). **Consequences can be severe, e.g., academic probation or dismissal, grade penalties, a 0 grade for assignments concerned, and prohibition from withdrawing from the class.**

Instructions for the assignment

- This assignment consists of both programming and theory questions.
- To switch between cell for code and for markdown, see the menu -> Cell -> Cell Type
- You can directly type Latex equations into markdown cells.
- If a question requires a picture, you could use this syntax "`<imgsrc ="" style = " width : 300px; " / >`" to include them within your ipython notebook.
- Your write up must be submitted in PDF form. You may use either Latex, markdown, or any word processing software. We will **NOT** accept handwritten work. Make sure that your work is formatted correctly, for example submit $\sum_{i=0}^n x_i$ instead of `\text{sum}_{i=0} x_i`
- When submitting the non-programming part of your assignment, you must correctly map pages of your PDF to each question/subquestion to reflect where they appear. Improperly mapped questions may not be graded correctly.
- Discussion is encouraged on Edstem as part of the Q/A. You may discuss high-level ideas with other students at the "whiteboard" level (e.g. how cross validation works, using matmul instead of dot) and review any relevant materials online. However, all assignments should be done individually, each student must write up and submit their own answers.
- **Graduate Students:** You are required to complete any sections marked as Bonus for Undergrads

Using the autograder

- You will find three assignments (for grads) on Gradescope that correspond to HW4: "Assignment 4 Programming", "Assignment 4 - Non-programming" and "Assignment 4 Programming - Bonus for all". Undergrads will have an additional assignment called "Assignment 4 Programming - Bonus for Undergrads".
- You will submit your code for the autograder in the Assignment 4 Programming sections. Please refer to the Deliverables and Point Distribution section for what parts are considered required, bonus for undergrads, and bonus for all".
- We provided you different .py files and we added libraries in those files please DO NOT remove those lines and add your code after those lines. Note that these are the only allowed libraries that you can use for the homework
- You are allowed to make as many submissions until the deadline as you like. Additionally, note that the autograder tests each function separately, therefore it can serve as a useful tool to help you debug your code if you are not sure of what part of your implementation might have an issue
- For the "Assignment 4 - Non-programming" part, you will download your Jupyter Notebook as html and submit it as a PDF on Gradescope. To download the notebook as html, click on "File" on the top left corner of this page and select "Download as > html". Then, open the html file and print to PDF. Please refer to the Deliverables and Point Distribution section for an outline of the non-programming questions.
- When submitting to Gradescope, please make sure to mark the page(s) corresponding to each problem/sub-problem.

Deliverables and Points Distribution

Q1: Two Layer NN [70 pts; 55pts + 15pts Undergrad Bonus]

Deliverables: NN.py and Notebook Graphs

- **Implementation** [55pts; 45pts + 10pts Bonus for CS 4641] - *programming*
 - relu [5pts]

- tanh [5pts]
- loss [5pts]
- forward [10pts]
- backward [10pts]
- Gradient Descent [10pts]
- Stochastic Gradient Descent [10pts Bonus for CS 4641]
- **Questions** [10pts: 5pts + 5pts Bonus for CS 4641] - *non-programming*
 - Loss plot and MSE for Gradient Descent [5pts]
 - Loss plot and MSE for Stochastic Gradient Descent [5pts Bonus for CS 4641]

Q2: CNN [15pts; 12pts Undergrad Bonus + 3pts Bonus for All]

Deliverables: `cnn.py` and **Written Report**

- 2.1.3 [2pts Bonus for CS 4641] - *programming*
- 2.1.4 [8pts Bonus for CS 4641] - *non-programming*
- 2.1.5 [2pts Bonus for CS 4641] - *non-programming*
- 2.2 [3pts Bonus for All] - *non-programming*

Q3: Random Forest [50pts]

Deliverables: `random_forest.py` and **Written Report**

- 3.1 Random Forest Implementation [35pts] - *programming*
- 3.2 Hyperparameter Tuning with a Random Forest [5pts] - *non-programming*
- 3.3 Plotting Feature Importance [5pts] - *non-programming*
- 3.4 Improvement [5pts] - *non-programming*

Q4: SVM [30pts Bonus for all]

Deliverables: `feature.py` and **Written Report**

- 4.1: Fitting an SVM Classifier by hand [20pts] - *non programming*
- 4.2: Feature Mapping [10pts] - *programming*

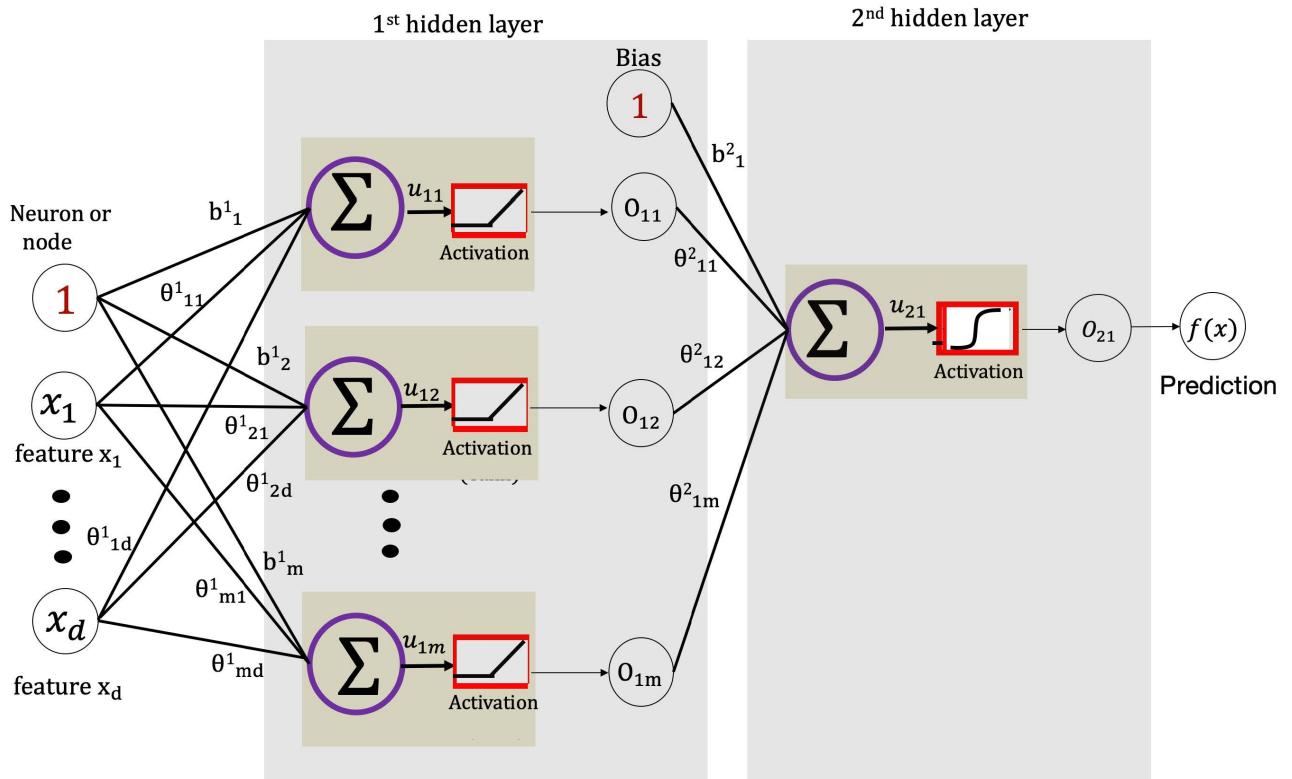
Environment Setup

```
In [3]: 1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.datasets import load_diabetes
4 from sklearn.preprocessing import MinMaxScaler
5
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import classification_report
8 from sklearn.metrics import plot_confusion_matrix
9
10 from collections import Counter
11 from scipy import stats
12 from math import log2, sqrt
13 import pandas as pd
14 import time
15 from sklearn.model_selection import train_test_split
16 from sklearn.preprocessing import LabelEncoder
17 from sklearn.tree import DecisionTreeClassifier
18
19 from sklearn.datasets import make_moons
20 from sklearn.metrics import accuracy_score
21 from sklearn import svm
22 %load_ext autoreload
23
24 %autoreload 2
25
26 %reload_ext autoreload
```

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

1: Two Layer Neural Network [70 pts; 55pts + 15pts Undergrad Bonus] [P][W]

Perceptron



Notation clarification – superscript represents the layer number, subscripts represent the specific units in two adjacent layers being connected by theta

θ^1_{21} - theta of the 1st layer connecting the 2nd hidden unit of the 1st layer and the 1st input unit

θ^2_{12} - theta of the 2nd layer connecting the 1st hidden unit of the 2nd layer and the 2nd hidden unit of the 1st layer

b^1_1 - bias of the 1st hidden unit of the 1st layer

A single layer perceptron can be thought of as a linear hyperplane as in logistic regression followed by a non-linear activation function.

$$u_i = \sum_{j=1}^d \theta_{ij} x_j + b_i$$

$$o_i = \phi \left(\sum_{j=1}^d \theta_{ij} x_j + b_i \right) = \phi(\theta_i^T x + b_i)$$

where x is a d-dimensional vector i.e. $x \in R^d$. It is one datapoint with d features. $\theta_i \in R^d$ is the weight vector for the i^{th} hidden unit, $b_i \in R$ is the bias element for the i^{th} hidden unit and $\phi(\cdot)$ is a non-linear activation function that has been described below. u_i is a linear combination of the features in x_j weighted by θ_i whereas o_i is the i^{th} output unit from the activation layer.

Fully connected Layer

Typically, a modern neural network contains millions of perceptrons as the one shown in the previous image. Perceptrons interact in different configurations such as cascaded or parallel. In this part, we describe a fully connected layer configuration in a neural network which comprises multiple parallel perceptrons forming one layer.

We extend the previous notation to describe a fully connected layer. Each layer in a fully connected network has a number of input/hidden/output units cascaded in parallel. Let us define a single layer of the neural net as follows:

m denotes the number of hidden units in a single layer l whereas n denotes the number of units in the previous layer $l - 1$.

$$u^{[l]} = \theta^{[l]} o^{[l-1]} + b^{[l]}$$

where $u^{[l]} \in R^m$ is a m -dimensional vector pertaining to the hidden units of the l^{th} layer of the neural network after applying linear operations. Similarly, $o^{[l-1]}$ is the n -dimensional output vector corresponding to the hidden units of the $(l - 1)^{th}$ activation layer.

$\theta^{[l]} \in R^{m \times n}$ is the weight matrix of the l^{th} layer where each row of $\theta^{[l]}$ is analogous to θ_i described in the previous section i.e. each row

corresponds to one hidden unit of the i^{th} layer. $b^{[l]} \in \mathbb{R}^m$ is the bias vector of the layer where each element of b pertains to one hidden unit of the i^{th} layer. This is followed by element wise non-linear activation function $o^{[l]} = \phi(u^{[l]})$. The whole operation can be summarized as,

$$o^{[l]} = \phi(\theta^{[l]} o^{[l-1]} + b^{[l]})$$

where $o^{[l-1]}$ is the output of the previous layer.

Activation Function

There are many activation functions in the literature but for this question we are going to use Relu and Tanh only.

HINT 1: When calculating the tanh and relu function, make sure you are not modifying the values in the original passed in matrix. You may find `np.copy()` helpful.

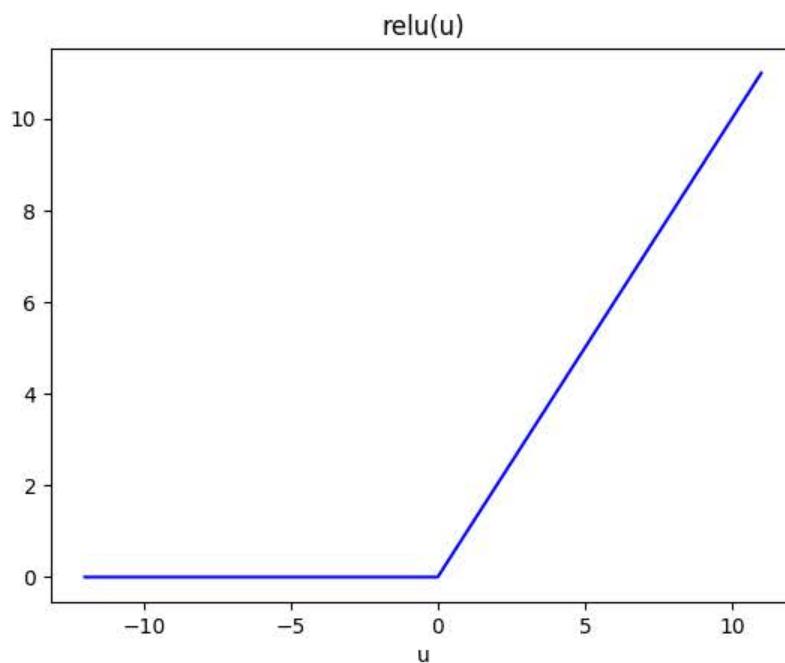
Relu

The rectified linear unit (Relu) is one of the most commonly used activation functions in deep learning models. The mathematical form is

$$o = \phi(u) = \max(0, u)$$

One of the advantages of Relu is that it is a fast nonlinearity.

The derivative of relu function is given as $o' = \phi'(u) = \begin{cases} 0 & u \leq 0 \\ 1 & u > 0 \end{cases}$



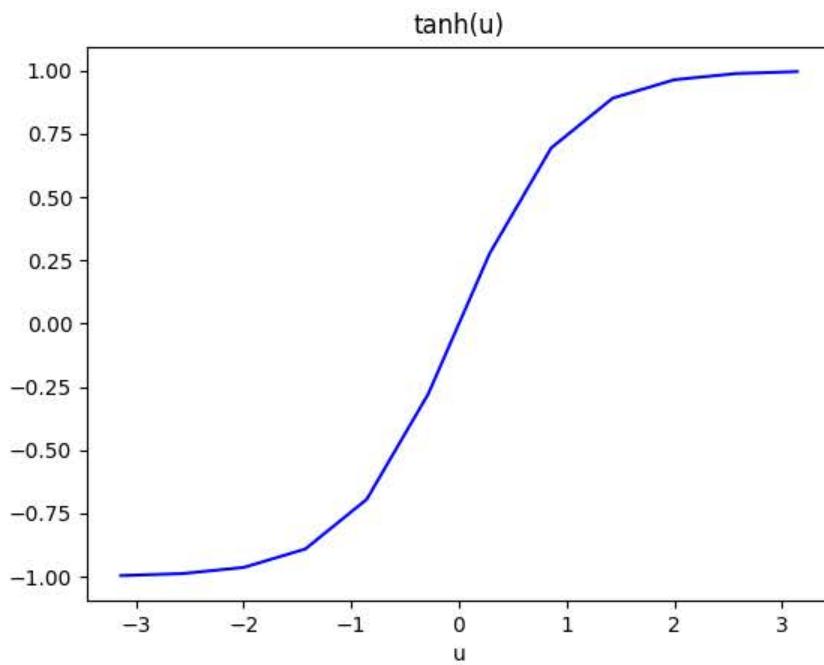
Tanh

Tanh also known as hyperbolic tangent is like a shifted version of sigmoid activation function with its range going from -1 to 1. Tanh almost always proves to be better than the sigmoid function since the mean of the activations are closer to zero. Tanh has an effect of centering data that makes learning for the next layer a bit easier. The mathematical form of tanh is given as

$$o = \phi(u) = \tanh(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}}$$

The derivative of tanh is given as

$$o' = \phi'(u) = 1 - \left(\frac{e^u - e^{-u}}{e^u + e^{-u}} \right)^2 = 1 - o^2$$



Sigmoid

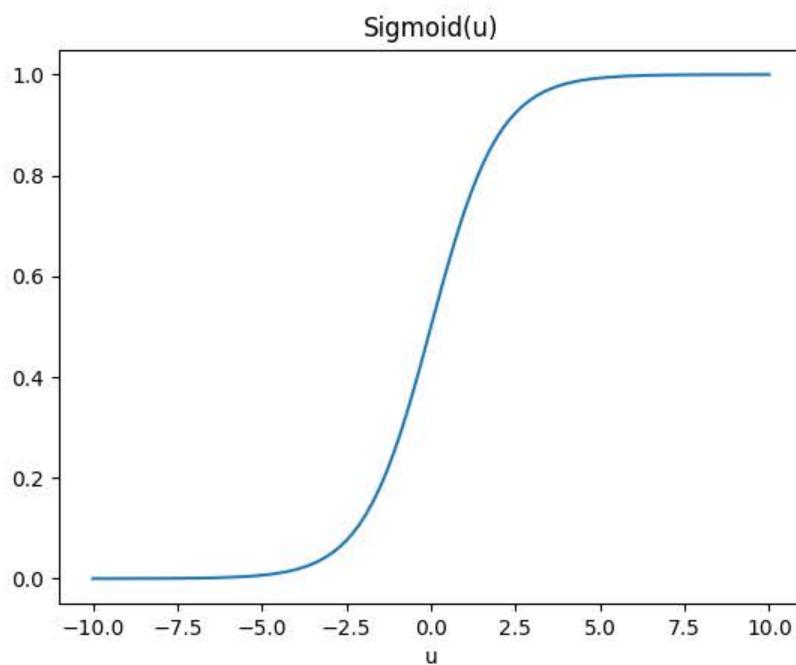
The sigmoid function is another non-linear function with S-shaped curve. This function is useful in the case of binary classification as its output is between 0 and 1. The mathematical form of the function is

$$o = \phi(u) = \frac{1}{1 + e^{-u}}$$

The derivation of the sigmoid function has a nice form and is given as

$$o' = \phi'(u) = \frac{1}{1 + e^{-u}} \left(1 - \frac{1}{1 + e^{-u}} \right) = \phi(u)(1 - \phi(u))$$

Note: We will not be using sigmoid activation function for this assignment. This is included only for the sake of completeness.



Mean Squared Error

It is an estimator that measures the average of the squares of the errors i.e. the average squared difference between the actual and the estimated values. It estimates the quality of the learnt hypothesis between the actual and the predicted values. It's non-negative and closer to zero, the better the learnt function is.

Implementation details

For regression problems as in this exercise, we compute the loss as follows:

$$MSE = \frac{1}{2N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

where y_i is the true label and \hat{y}_i is the estimated label. We use a factor of $\frac{1}{2N}$ instead of $\frac{1}{N}$ to simply the derivative of loss function.

Forward Propagation

We start by initializing the weights of the fully connected layer using Xavier initialization [Xavier initialization](#) (<http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf>). During training, we pass all the data points through the network layer by layer using forward propagation. The main equations for forward prop have been described below.

$$\begin{aligned} u^{[0]} &= x \\ u^{[1]} &= \theta^{[1]} u^{[0]} + b^{[1]} \\ o^{[1]} &= \text{Relu}(u^{[1]}) \\ u^{[2]} &= \theta^{[2]} o^{[1]} + b^{[2]} \\ \hat{y} &= o^{[2]} = \text{Tanh}(u^{[2]}) \end{aligned}$$

Then we get the output and compute the loss

$$l = \frac{1}{2N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Backward propagation

After the forward pass, we do back propagation to update the weights and biases in the direction of the negative gradient of the loss function. So, we update the weights and biases using the following formulas

$$\begin{aligned} \theta^{[2]} &:= \theta^{[2]} - lr \times \frac{\partial l}{\partial \theta^{[2]}} \\ b^{[2]} &:= b^{[2]} - lr \times \frac{\partial l}{\partial b^{[2]}} \\ \theta^{[1]} &:= \theta^{[1]} - lr \times \frac{\partial l}{\partial \theta^{[1]}} \\ b^{[1]} &:= b^{[1]} - lr \times \frac{\partial l}{\partial b^{[1]}} \end{aligned}$$

where lr is the learning rate. It decides the step size we want to take in the direction of the negative gradient.

To compute the terms $\frac{\partial l}{\partial \theta^{[2]}}$ and $\frac{\partial l}{\partial b^{[2]}}$ we use chain rule for differentiation as follows:

$$\begin{aligned} \frac{\partial l}{\partial \theta^{[2]}} &= \frac{\partial l}{\partial o^{[2]}} \frac{\partial o^{[2]}}{\partial u^{[2]}} \frac{\partial u^{[2]}}{\partial \theta^{[2]}} \\ \frac{\partial l}{\partial b^{[2]}} &= \frac{\partial l}{\partial o^{[2]}} \frac{\partial o^{[2]}}{\partial u^{[2]}} \frac{\partial u^{[2]}}{\partial b^{[2]}} \end{aligned}$$

So, $\frac{\partial l}{\partial o^{[2]}}$ is the differentiation of the loss function at point $o^{[2]}$

$\frac{\partial o^{[2]}}{\partial u^{[2]}}$ is the differentiation of the Tanh function at point $u^{[2]}$

$\frac{\partial u^{[2]}}{\partial \theta^{[2]}}$ is equal to $\theta^{[1]}$

$\frac{\partial u^{[2]}}{\partial b^{[2]}}$ is equal to 1.

To compute $\frac{\partial l}{\partial \theta^{[2]}}$, we need $\sigma^{[2]}$, $u^{[2]}$ & $\sigma^{[1]}$ which are calculated during forward propagation. So we need to store these values in cache variables during forward propagation to be able to access them during backward propagation. Similarly for calculating other partial derivatives, we store the values we'll be needing for chain rule in cache. These values are obtained from the forward propagation and used in backward propagation. The cache is implemented as a dictionary here where the keys are the variable names and the values are the variables values.

Also, the functional form of the MSE differentiation and Relu differentiation are given by

$$\begin{aligned}\frac{\partial l}{\partial \sigma^{[2]}} &= (\sigma^{[2]} - y) \\ \frac{\partial l}{\partial u^{[2]}} &= \frac{\partial l}{\partial \sigma^{[2]}} * (1 - (\tanh^2(u^{[2]}))) \\ \frac{\partial u^{[2]}}{\partial \theta^{[2]}} &= \sigma^{[1]} \\ \frac{\partial u^{[2]}}{\partial b^{[2]}} &= 1\end{aligned}$$

On vectorization, the above equations become:

$$\begin{aligned}\frac{\partial l}{\partial \sigma^{[2]}} &= \frac{1}{n}(\sigma^{[2]} - y) \\ \frac{\partial l}{\partial \theta^{[2]}} &= \frac{1}{n} \frac{\partial l}{\partial u^{[2]}} \sigma^{[1]} \\ \frac{\partial l}{\partial b^{[2]}} &= \frac{1}{n} \sum \frac{\partial l}{\partial u^{[2]}}\end{aligned}$$

HINT 2: Division by N only needs to occur ONCE for any derivative that requires a division by N . Make sure you avoid cascading divisions by N where you might accidentally divide your derivative by N^2 or greater.

This completes the differentiation of loss function w.r.t to parameters in the second layer. We now move on to the first layer, the equations for which are given as follows:

$$\begin{aligned}\frac{\partial l}{\partial \theta^{[1]}} &= \frac{\partial l}{\partial \sigma^{[2]}} \frac{\partial \sigma^{[2]}}{\partial u^{[2]}} \frac{\partial u^{[2]}}{\partial \theta^{[1]}} \frac{\partial \sigma^{[1]}}{\partial u^{[1]}} \frac{\partial u^{[1]}}{\partial \theta^{[1]}} \\ \frac{\partial l}{\partial b^{[1]}} &= \frac{\partial l}{\partial \sigma^{[2]}} \frac{\partial \sigma^{[2]}}{\partial u^{[2]}} \frac{\partial u^{[2]}}{\partial \theta^{[1]}} \frac{\partial \sigma^{[1]}}{\partial u^{[1]}} \frac{\partial u^{[1]}}{\partial b^{[1]}}\end{aligned}$$

Where

$$\begin{aligned}\frac{\partial u^{[2]}}{\partial \theta^{[1]}} &= \theta^{[2]} \\ \frac{\partial \sigma^{[1]}}{\partial u^{[1]}} &= 1(\sigma^{[1]} > 0) \\ \frac{\partial u^{[1]}}{\partial \theta^{[1]}} &= x \\ \frac{\partial u^{[1]}}{\partial b^{[1]}} &= 1\end{aligned}$$

Note that $\frac{\partial \sigma^{[1]}}{\partial u^{[1]}}$ is the differentiation of the ReLU function at $u^{[1]}$.

The above equations outline the forward and backward propagation process for a 2-layer fully connected neural net with Tanh as the first activation layer and Relu has the second one. The same process can be extended to different neural networks with different activation layers.

Code Implementation:

$$\begin{aligned}
dLoss_o2 &= \frac{\partial l}{\partial o^{[2]}} \Rightarrow \text{dim} = (1, 331) \\
dLoss_u2 &= dLoss_o2 \frac{\partial o^{[2]}}{\partial u^{[2]}} \Rightarrow \text{dim} = (1, 331) \\
dLoss_theta2 &= dLoss_u2 \frac{\partial u^{[2]}}{\partial \theta^{[2]}} \Rightarrow \text{dim} = (1, 15) \\
dLoss_b2 &= dLoss_u2 \frac{\partial u^{[2]}}{\partial b^{[2]}} \Rightarrow \text{dim} = (1, 1) \\
dLoss_o1 &= dLoss_u2 \frac{\partial o^{[2]}}{\partial o^{[1]}} \Rightarrow \text{dim} = (15, 331) \\
dLoss_u1 &= dLoss_o1 \frac{\partial o^{[1]}}{\partial u^{[1]}} \Rightarrow \text{dim} = (15, 331) \\
dLoss_theta1 &= dLoss_u1 \frac{\partial u^{[1]}}{\partial \theta^{[1]}} \Rightarrow \text{dim} = (15, 10) \\
dLoss_b1 &= dLoss_u1 \frac{\partial u^{[1]}}{\partial b^{[1]}} \Rightarrow \text{dim} = (15, 1)
\end{aligned}$$

Note: Training set has 331 examples.

Question

In this question, you will implement a two layer fully connected neural network. You will also experiment with different activation functions and optimization techniques. Functions with comments "TODO: implement this" are for you to implement. We provide two activation functions here - Relu and Tanh. You will implement a neural network that would have tanh activation followed by relu layer.

You'll also implement Gradient Descent (GD) and Batch Gradient Descent (BGD) algorithms for training these neural nets. **GD is mandatory for all. BGD is bonus for undergraduate students but mandatory for graduate students.**

In the **NN.py** file, complete the following functions:

- **Relu:** Recall Hint 1
- **Tanh:** Recall Hint 1
- **nloss**
- **forward**
- **backward:** Recall Hint 2
- **gradient_descent**
- **batch_gradient_descent:** Mandatory for graduate students, bonus for undergraduate students. Please batch your data in a wraparound manner. For example, given a dataset of 9 numbers, [1, 2, 3, 4, 5, 6, 7, 8, 9], and a batch size of 6, the first iteration batch will be [1, 2, 3, 4, 5, 6], the second iteration batch will be [7, 8, 9, 1, 2, 3], the third iteration batch will be [4, 5, 6, 7, 8, 9], etc...

We'll train this neural net on sklearn's diabetes dataset. Graduate students have to use both GD and BGD to optimize their neural net. Undergraduate students have to implement GD while BGD is bonus for them. Note: it is possible you'll run into nan or negative values for loss. This happens because of the small dataset we're using and some numerical stability issues that arise due to division by zero, natural log of zeros etc. You can experiment with the total number of iterations to mitigate this.

You're free to tune hyperparameters like the batch size, number of hidden units in each layer etc. if that helps you in achieving the desired MSE values to pass the autograder tests. However, you're advised to try out the default values first.

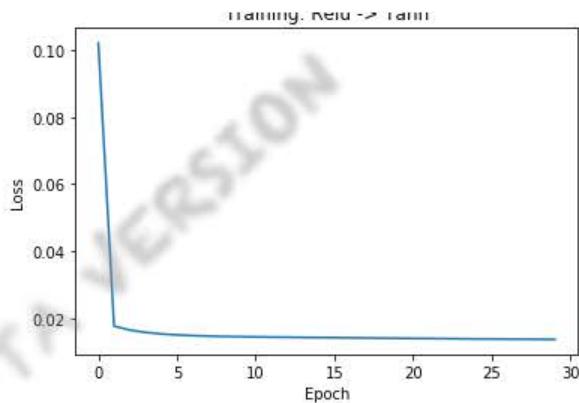
Deliverables for this question:

1. Loss plot and MSE value for neural net with gradient descent
2. Loss plot and MSE value for neural net with batch gradient descent (mandatory for graduate students, bonus for undergraduate students)

```
In [8]: 1 ##### DO NOT CHANGE THIS CELL #####
2 #####
3 #####
4
5 from NN import dlnet
6
7 dataset = load_diabetes() # Load the dataset
8 x, y = dataset.data, dataset.target
9 y = y.reshape(-1,1)
10
11 x_train, x_test, y_train, y_test = train_test_split(x, y, random_state=1) #split data
12
13 x_scale = MinMaxScaler()
14 x_train = x_scale.fit_transform(x_train) #normalize data
15 x_test = x_scale.transform(x_test)
16
17 y_scale = MinMaxScaler()
18 y_train = y_scale.fit_transform(y_train)
19 y_test = y_scale.transform(y_test)
20
21 x_train, x_test, y_train, y_test = x_train.T, x_test.T, y_train.reshape(1,-1), y_test #condition data
22
23 nn = dlnet(x_train,y_train,lr=0.01) # initialize neural net class
24 nn.gradient_descent(x_train, y_train, iter = 60000) #train
25
26 # create figure
27 fig = plt.plot(np.array(nn.loss).squeeze())
28 plt.title(f'Training: {nn.neural_net_type}')
29 plt.xlabel("Epoch")
30 plt.ylabel("Loss")
31
```

```
Loss after iteration 0: 0.102091
Loss after iteration 2000: 0.017786
Loss after iteration 4000: 0.016579
Loss after iteration 6000: 0.015885
Loss after iteration 8000: 0.015426
Loss after iteration 10000: 0.015120
Loss after iteration 12000: 0.014923
Loss after iteration 14000: 0.014776
Loss after iteration 16000: 0.014668
Loss after iteration 18000: 0.014589
Loss after iteration 20000: 0.014530
Loss after iteration 22000: 0.014479
Loss after iteration 24000: 0.014436
Loss after iteration 26000: 0.014392
Loss after iteration 28000: 0.014348
Loss after iteration 30000: 0.014307
Loss after iteration 32000: 0.014260
Loss after iteration 34000: 0.014216
Loss after iteration 36000: 0.014173
Loss after iteration 38000: 0.014130
Loss after iteration 40000: 0.014086
Loss after iteration 42000: 0.014046
Loss after iteration 44000: 0.014009
Loss after iteration 46000: 0.013974
Loss after iteration 48000: 0.013940
Loss after iteration 50000: 0.013910
Loss after iteration 52000: 0.013882
Loss after iteration 54000: 0.013855
Loss after iteration 56000: 0.013828
Loss after iteration 58000: 0.013804
```

```
Out[8]: Text(0, 0.5, 'Loss')
```



```
In [9]: 1 #####  
2 ### DO NOT CHANGE THIS CELL ###  
3 #####  
4  
5 y_predicted = nn.predict(x_test) # predict  
6 y_test = y_test.reshape(1,-1)  
7 print("Mean Squared Error (MSE)", (np.sum((y_predicted-y_test)**2)/y_test.shape[1]))  
8  
9
```

Mean Squared Error (MSE) 0.029385750726051957

In [10]:

```
1 #####  
2 ### DO NOT CHANGE THIS CELL ###  
3 #####  
4  
5 from NN import dlnet  
6  
7 dataset = load_diabetes() # load the dataset  
8 x, y = dataset.data, dataset.target  
9 y = y.reshape(-1,1)  
10  
11 x_train, x_test, y_train, y_test = train_test_split(x, y, random_state=1) #split data  
12  
13 x_scale = MinMaxScaler()  
14 x_train = x_scale.fit_transform(x_train) #normalize data  
15 x_test = x_scale.transform(x_test)  
16  
17 y_scale = MinMaxScaler()  
18 y_train = y_scale.fit_transform(y_train)  
19 y_test = y_scale.transform(y_test)  
20  
21 x_train, x_test, y_train, y_test = x_train.T, x_test.T, y_train.reshape(1,-1), y_test #condition data  
22  
23 nn = dlnet(x_train,y_train,lr=0.01) # initialize neural net class  
24 nn.batch_gradient_descent(x_train, y_train, iter = 60000) #train  
25  
26  
27 # create figure  
28 fig = plt.plot(np.array(nn.loss).squeeze())  
29 plt.title(f'Training: {nn.neural_net_type}')  
30 plt.xlabel("Epoch")  
31 plt.ylabel("Loss")
```

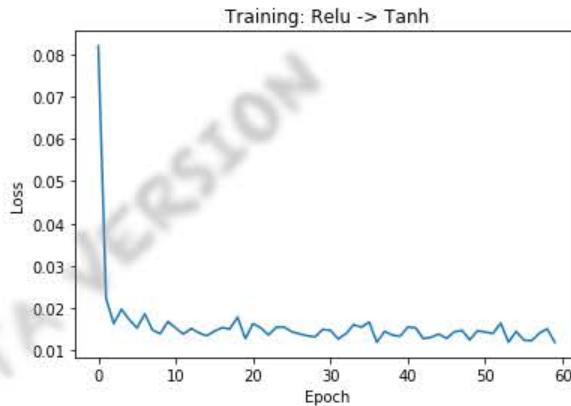
```
Loss after iteration 0: 0.081986  
Loss after iteration 1000: 0.022265  
Loss after iteration 2000: 0.016349  
Loss after iteration 3000: 0.019730  
Loss after iteration 4000: 0.017276  
Loss after iteration 5000: 0.015297  
Loss after iteration 6000: 0.018646  
Loss after iteration 7000: 0.014828  
Loss after iteration 8000: 0.013929  
Loss after iteration 9000: 0.016813  
Loss after iteration 10000: 0.015312  
Loss after iteration 11000: 0.013854  
Loss after iteration 12000: 0.015212  
Loss after iteration 13000: 0.014155  
Loss after iteration 14000: 0.013456  
Loss after iteration 15000: 0.014531  
Loss after iteration 16000: 0.015389  
Loss after iteration 17000: 0.015070  
Loss after iteration 18000: 0.017886  
Loss after iteration 19000: 0.012811  
Loss after iteration 20000: 0.016322  
Loss after iteration 21000: 0.015344  
Loss after iteration 22000: 0.013649  
Loss after iteration 23000: 0.015487  
Loss after iteration 24000: 0.015536  
Loss after iteration 25000: 0.014410  
Loss after iteration 26000: 0.013925  
Loss after iteration 27000: 0.013462  
Loss after iteration 28000: 0.013222  
Loss after iteration 29000: 0.014958  
Loss after iteration 30000: 0.014749  
Loss after iteration 31000: 0.012706  
Loss after iteration 32000: 0.013952  
Loss after iteration 33000: 0.016106  
Loss after iteration 34000: 0.015481  
Loss after iteration 35000: 0.016711  
Loss after iteration 36000: 0.011924  
Loss after iteration 37000: 0.014528  
Loss after iteration 38000: 0.013631  
Loss after iteration 39000: 0.013360  
Loss after iteration 40000: 0.015537  
Loss after iteration 41000: 0.015346  
Loss after iteration 42000: 0.012806  
Loss after iteration 43000: 0.013096
```

```

Loss after iteration 44000: 0.013884
Loss after iteration 45000: 0.012872
Loss after iteration 46000: 0.014391
Loss after iteration 47000: 0.014745
Loss after iteration 48000: 0.012497
Loss after iteration 49000: 0.014659
Loss after iteration 50000: 0.014345
Loss after iteration 51000: 0.014003
Loss after iteration 52000: 0.016489
Loss after iteration 53000: 0.011959
Loss after iteration 54000: 0.014485
Loss after iteration 55000: 0.012429
Loss after iteration 56000: 0.012323
Loss after iteration 57000: 0.014087
Loss after iteration 58000: 0.015109
Loss after iteration 59000: 0.011909

```

Out[10]: Text(0, 0.5, 'Loss')



```

In [11]: 1 #####
2 ### DO NOT CHANGE THIS CELL ####
3 #####
4
5 y_predicted = nn.predict(x_test) # predict
6 y_test = y_test.reshape(1,-1)
7 print("Mean Squared Error (MSE)", (np.sum((y_predicted-y_test)**2)/y_test.shape[1]))
8
9

```

Mean Squared Error (MSE) 0.02938477267346788

2: Image Classification based on Convolutional Neural Networks [15pts total = 12pts Bonus for Undergrad + 3pts Bonus for all] [W]

Keras is a deep learning API written in Python, running on top of the machine learning platform TensorFlow. It was developed with a focus on enabling fast experimentation. Being able to go from idea to result as fast as possible is key to doing good research. In this part, you will build a convolutional neural network based on TF/Keras to solve the image classification task for the fashion mnist dataset. If you haven't installed TensorFlow, you can install the package by [pip](#) command or train your model by uploading HW4 notebook to [Colab](#) (<https://colab.research.google.com/>) directly. Colab contains all packages you need for this section.

Hint1: [First contact with Keras \(https://keras.io/about/\)](https://keras.io/about/)

Hint2: [How to Install Keras \(https://www.pyimagesearch.com/2016/07/18/installing-keras-for-deep-learning/\)](https://www.pyimagesearch.com/2016/07/18/installing-keras-for-deep-learning/)

Hint3: [CS231n Tutorial \(Layers used to build ConvNets\) \(https://cs231n.github.io/convolutional-networks/\)](https://cs231n.github.io/convolutional-networks/)

Environment Setup

We use fashion mnist dataset to train our model. This is a dataset of 60,000 28x28 training images and 10,000 test images, labeled over 10 categories. Each example is **28 x 28** pixel grayscale image of various clothing items.

In [7]:

```
1 #####  
2 ### DO NOT CHANGE THIS CELL ###  
3 #####  
4  
5 # split data between train and test sets  
6 (x_train, y_train), (x_test, y_test) = fashion_mnist.load_data()  
7  
8 # input image dimensions  
9 img_rows, img_cols = 28, 28  
10 number_channels = 1  
11 #set num of classes  
12 num_classes = 10  
13  
14 if tf.keras.backend.image_data_format() == 'channels_first':  
15     x_train = x_train.reshape(x_train.shape[0], number_channels, img_rows, img_cols)  
16     x_test = x_test.reshape(x_test.shape[0], number_channels, img_rows, img_cols)  
17     input_shape = (number_channels, img_rows, img_cols)  
18 else:  
19     x_train = x_train.reshape(x_train.shape[0], img_rows, img_cols, number_channels)  
20     x_test = x_test.reshape(x_test.shape[0], img_rows, img_cols, number_channels)  
21     input_shape = (img_rows, img_cols, number_channels)  
22  
23 x_train = x_train.astype('float32')  
24 x_test = x_test.astype('float32')  
25 x_train /= 255  
26 x_test /= 255  
27 print('x_train shape:', x_train.shape)  
28 print('x_test shape:', x_test.shape)  
29 print(x_train.shape[0], 'train samples')  
30 print(x_test.shape[0], 'test samples')  
31  
32 fashion_classes = ["T-shirt/top", "Trouser", "Pullover", "Dress", "Coat", "Sandal", "Shirt", "Sneaker"  
33 # convert class vectors to binary class matrices  
34 y_train = tf.keras.utils.to_categorical(y_train, num_classes)  
35 y_test = tf.keras.utils.to_categorical(y_test, num_classes)
```

x_train shape: (60000, 28, 28, 1)
x_test shape: (10000, 28, 28, 1)
60000 train samples
10000 test samples

2.1.2 Load some sample images from fashion mnist

```
In [9]: 1 #####  
2 ### DO NOT CHANGE THIS CELL ###  
3 #####  
4  
5 # Show some images from fashion_mnist  
6  
7 fig = plt.figure(figsize=(20, 10))  
8 for i in range(50):  
9     random_index = np.random.randint(0, len(y_train))  
10    ax = fig.add_subplot(5, 10, i+1)  
11    ax.imshow(x_train[random_index, :].squeeze(2), cmap='Greys')  
12 plt.show()
```



As you can see from above, the fashion mnist dataset contains a selection of fashion objects. The images have been size-normalized and objects remain centered in fixed-size images.

2.1.3 Build convolutional neural network model [2pts]

In this part, you need to build a convolutional neural network as described below. All coding should be done in `cnn.py`. The architecture of the model is:

[INPUT - CONV - CONV - MAXPOOL - DROPOUT - CONV - CONV - MAXPOOL - DROPOUT - FC1 - DROPOUT - FC2]

INPUT: $[28 \times 28 \times 1]$ will hold the raw pixel values of the image, in this case, an image of width 28, height 28. This layer should give 16 filters and have appropriate padding to maintain shape.

CONV: Conv. layer will compute the output of neurons that are connected to local regions in the input, each computing a dot product between their weights and a small region they are connected to the input volume. We decide to set the kernel_size 3×3 for the both Conv. layers. For example, the output of the Conv. layer may look like $[28 \times 28 \times 32]$ if we use 32 filters. Again, we use padding to maintain shape.

MAXPOOL: MAXPOOL layer will perform a downsampling operation along the spatial dimensions (width, height). With pool size of 2×2 , resulting shape takes form 14×14 .

DROPOUT: DROPOUT layer with the dropout rate of 0.25, to prevent overfitting.

CONV: Additonal Conv. layer take outputs from above layers and applies more filters. The Conv. layer may look like $[14 \times 14 \times 32]$. We set the kernel_size 3×3 and use padding to maintain shape for both Conv. layers.

CONV: Additonal Conv. layer take outputs from above layers and applies more filters. The Conv. layer may look like $[14 \times 14 \times 64]$.

MAXPOOL: MAXPOOL layer will perform a downsampling operation along the spatial dimensions (width, height).

DROPOUT: Dropout layer with the dropout rate of 0.25, to prevent overfitting.

FC1: Dense layer which takes input above layers, and has 256 neurons. Flatten operations may be useful.

DROPOUT: Dropout layer with the dropout rate of 0.5, to prevent overfitting.

FC2: Dense layer with 10 neurons, and softmax activation, is the final layer. The dimension of the output space is the number of classes.

Activation function: Use LeakyReLU unless otherwise indicated to build you model architecture.

Note that while this is a suggested model design, you may use other architectures and experiment with different layers for better results.

Use the following keras links to reference crucial layers of the model in keras API:

- Conv2d: https://keras.io/api/layers/convolution_layers/convolution2d/ (https://keras.io/api/layers/convolution_layers/convolution2d/)
- Dense: https://keras.io/api/layers/core_layers/dense/ (https://keras.io/api/layers/core_layers/dense/)
- Flatten: https://keras.io/api/layers/reshaping_layers/flatten/ (https://keras.io/api/layers/reshaping_layers/flatten/)
- MaxPool: https://keras.io/api/layers/pooling_layers/max_pooling2d/ (https://keras.io/api/layers/pooling_layers/max_pooling2d/)
- Dropout: https://keras.io/api/layers/regularization_layers/dropout/ (https://keras.io/api/layers/regularization_layers/dropout/)

And explore the keras layers API if you would like to experiment with additional layers: <https://keras.io/api/layers/> (<https://keras.io/api/layers/>)

```
In [9]: 1 ##### DO NOT CHANGE THIS CELL #####
2 ##### DO NOT CHANGE THIS CELL #####
3 ##### DO NOT CHANGE THIS CELL #####
4
5 # Show the architecture of the model
6 achi=plt.imread('data/images/Architecture.png')
7 fig = plt.figure(figsize=(10,10))
8 plt.imshow(achi)
```

Out[9]: <matplotlib.image.AxesImage at 0x7f82f2bdb390>

Model: "sequential"		
Layer (type)	Output Shape	Param #
conv2d (Conv2D)	(None, 28, 28, 16)	160
leaky_re_lu (LeakyReLU)	(None, 28, 28, 16)	0
conv2d_1 (Conv2D)	(None, 28, 28, 32)	4640
leaky_re_lu_1 (LeakyReLU)	(None, 28, 28, 32)	0
max_pooling2d (MaxPooling2D)	(None, 14, 14, 32)	0
dropout (Dropout)	(None, 14, 14, 32)	0
conv2d_2 (Conv2D)	(None, 14, 14, 32)	9248
leaky_re_lu_2 (LeakyReLU)	(None, 14, 14, 32)	0
conv2d_3 (Conv2D)	(None, 14, 14, 64)	18496
leaky_re_lu_3 (LeakyReLU)	(None, 14, 14, 64)	0
max_pooling2d_1 (MaxPooling2D)	(None, 7, 7, 64)	0
dropout_1 (Dropout)	(None, 7, 7, 64)	0
flatten (Flatten)	(None, 3136)	0
dense (Dense)	(None, 256)	803072
leaky_re_lu_4 (LeakyReLU)	(None, 256)	0
dropout_2 (Dropout)	(None, 256)	0
dense_1 (Dense)	(None, 10)	2570
activation (Activation)	(None, 10)	0
<hr/>		
Total params: 838,186		
Trainable params: 838,186		
Non-trainable params: 0		

You now need to set training variables in the `init()` function in `cnn.py`. Once you have defined variables you may use the cell below to see them.

- Recommended Batch Sizes fall in the range 16-256 (use powers of 2)
- Recommended Epoch Counts fall in the range 3-12
- Recommended Learning Rates fall in the range .0001-.01

```
In [ ]: 1 #####  
2 ### DO NOT CHANGE THIS CELL ###  
3 #####  
4  
5 # You can adjust parameters to train your model in __init__() in cnn.py  
6  
7 from cnn import CNN  
8  
9 net = CNN()  
10 batch_size, epochs, init_lr = net.get_vars()  
11 print(f'Batch Size\t: {batch_size} \nEpochs\t: {epochs} \nLearning Rate\t: {init_lr} \n')  
12
```

Defining model

You now need to complete the `create_net()` function in `cnn.py` to define your model structure. Once you have defined a model structure you may use the cell below to examine your architecture.

Your model is required to have at least 2 convolutional layers and at least 2 dense layers. ensuring that these requirements are met will earn you 2pts.

```
In [ ]: 1 #####  
2 ### DO NOT CHANGE THIS CELL ###  
3 #####  
4  
5 # model.summary() gives you details of your architecture.  
6 #You can compare your architecture with the 'Architecture.png'  
7  
8 from cnn import CNN  
9 net = CNN()  
10  
11 s = tf.keras.backend.clear_session()  
12 model=net.create_net()  
13 model.summary()
```

Compiling model

Next prepare the model for training by completing `compile_model()` in `cnn.py`. Remember we are performing 10-way classification when selecting a loss function.

```
In [13]: 1 #####  
2 ### DO NOT CHANGE THIS CELL ###  
3 #####  
4  
5 # Complete compile_model() in cnn.py.  
6 from cnn import CNN  
7  
8 net = CNN()  
9 model = net.compile_net(model)  
10 print(model)
```

```
<tensorflow.python.keras.engine.sequential.Sequential object at 0x7fbbeaae4ac10>
```

2.1.4 Train the network [8pts total (3pts, 3pts, 2pts)]

Tuning: Training the network is the next thing to try. You can set your parameter at the **Defining Variable** section. If your parameters are set properly, you should see the loss of the validation set decreased and the value of accuracy increased. It may take more than 15 minutes to train your model.

Expected Result: You should be able to achieve more than **90%** accuracy on the test set to get full 12 points. If you achieve accuracy between **70% to 80%**, you will only get 3 points. An accuracy between **80% to 90%** will earn an additional 3pts.

- **70% to 80%** earns 3pts
- **80% to 90%** earns 3pts more (6pts total)
- **90%+** earns 2pts more (8pts total)

Train your own CNN model

```
In [ ]: 1 #####  
2 ### DO NOT CHANGE THIS CELL ###  
3 #####  
4  
5 # Train the model  
6  
7 from cnn import CNN  
8  
9 net = CNN()  
10 batch_size, epochs, init_lr = net.get_vars()  
11  
12 def lr_scheduler(epoch):  
13     new_lr = init_lr * 0.9 ** epoch  
14     print("Learning rate:", new_lr)  
15     return new_lr  
16  
17 history = model.fit(  
18     x_train, y_train,  
19     batch_size=batch_size,  
20     epochs=epochs,  
21     callbacks=[tf.keras.callbacks.LearningRateScheduler(lr_scheduler)],  
22     shuffle=True,  
23     verbose=1,  
24     initial_epoch=0,  
25     validation_data=(x_test, y_test)  
26 )  
27 score = model.evaluate(x_test, y_test, verbose=0)  
28 print('Test loss:', score[0])  
29 print('Test accuracy:', score[1])
```

2.1.5 Examine accuracy and loss [2pts]

You should expect to see gradually decreasing loss and gradually increasing accuracy. Examine loss and accuracy by running the cell below, no editing is necessary. Having appropriate looking loss and accuracy plots will earn you the last 2pts for your convolutional net.

```
In [ ]: 1 #####  
2 ### DO NOT CHANGE THIS CELL ###  
3 #####  
4  
5 # list all data in history  
6 print(history.history.keys())  
7  
8 # summarize history for accuracy and loss  
9 plt.plot(history.history['accuracy'])  
10 plt.plot(history.history['val_accuracy'])  
11 plt.title('model accuracy')  
12 plt.ylabel('accuracy')  
13 plt.xlabel('epoch')  
14 plt.legend(['train', 'test'], loc='upper left')  
15 plt.show()  
16  
17 plt.plot(history.history['loss'])  
18 plt.plot(history.history['val_loss'])  
19 plt.title('model loss')  
20 plt.ylabel('loss')  
21 plt.xlabel('epoch')  
22 plt.legend(['train', 'test'], loc='upper left')  
23 plt.show()  
24
```

```
In [ ]: 1 #####
2 ### DO NOT CHANGE THIS CELL ####
3 #####
4
5 # make predictions
6 y_pred = model.predict_proba(x_test)
7 y_pred_classes = np.argmax(y_pred, axis=1)
8 y_pred_prob = np.max(y_pred, axis=1)
9 y_gt_classes = np.argmax(y_test, axis=1)
10
11 from sklearn.metrics import confusion_matrix, accuracy_score
12 plt.figure(figsize=(8, 7))
13 plt.imshow(confusion_matrix(y_gt_classes, y_pred_classes))
14 plt.title('Confusion matrix', fontsize=16)
15 plt.xticks(np.arange(10), fashion_classes, rotation=90, fontsize=12)
16 plt.yticks(np.arange(10), fashion_classes, fontsize=12)
17 plt.colorbar()
18 plt.show()
```

2.2 Exploring Deep CNN Architectures [3pts] (Bonus for all)

The network you have produced is rather simple relative to many of those used in industry and research. Researchers have worked to make CNN models deeper and deeper over the past years in an effort to gain higher accuracy in predictions. While your model is only a handful of layers deep, some state of the art deep architectures may include up to 150 layers. However, this process has not been without challenges.

One such problem is the problem of the vanishing gradient. The weights of a neural network are updated using the backpropagation algorithm. The backpropagation algorithm makes a small change to each weight in such a way that the loss of the model decreases. Using the chain rule, we can find this gradient for each weight. But, as this gradient keeps flowing backwards to the initial layers, this value keeps getting multiplied by each local gradient. Hence, the gradient becomes smaller and smaller, making the updates to the initial layers very small, increasing the training time considerably.

Many tactics have been used in an effort to solve this problem. One architecture, named ResNet, solves the vanishing gradient problem in a unique way. ResNet was developed at Microsoft Research to find better ways to train deep networks. Take a moment to explore how ResNet tackles the vanishing gradient problem by reading the original research paper here: <https://arxiv.org/pdf/1512.03385.pdf> (<https://arxiv.org/pdf/1512.03385.pdf>) (also included as PDF in papers directory).

Question: In your own words, explain how ResNet addresses the vanishing gradient problem in 1-2 sentences below: (Please type answers directly in the cell below.)

Answer:

3: Random Forests [50pts] [P] [W]

NOTE: Please use sklearn's DecisionTreeClassifier in your Random Forest implementation. You can find more details about this classifier here (<https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html#sklearn.tree.DecisionTreeClassifier>)

3.1 Random Forest Implementation (35 pts) [P]

The decision boundaries drawn by decision trees are very sharp, and fitting a decision tree of unbounded depth to a list of examples almost inevitably leads to **overfitting**. In an attempt to decrease the variance of a decision tree, we're going to use a technique called 'Bootstrap Aggregating' (often abbreviated 'bagging'). This stems from the idea that a collection of weak learners can learn decision boundaries as well as a strong learner. This is commonly called a Random Forest.

We can build a Random Forest as a collection of decision trees, as follows:

1. For every tree in the random forest, we're going to
 - a) Subsample the examples with replacement. Note that in this question, the size of the subsample data is equal to the original dataset.
 - b) From the subsamples in a), choose attributes at random to learn on in accordance with a provided attribute subsampling rate. Based on what it was mentioned in the class, we randomly pick features in each split. We use a more general approach here to make the programming part easier. Let's randomly pick some features (70% percent of features) and grow the tree based on the pre-determined randomly selected features. Therefore, there is no need to find random features in each split.
 - c) Fit a decision tree to the subsample of data we've chosen to a certain depth.

Classification for a random forest is then done by taking a majority vote of the classifications yielded by each tree in the forest after it classifies an example.

In RandomForest Class,

1. X is assumed to be a matrix with num_training rows and num_features columns where num_training is the number of total records and num_features is the number of features of each record.
2. y is assumed to be a vector of labels of length num_training.

NOTE: Lookout for TODOs for the parts that needs to be implemented.

3.2 Hyperparameter Tuning with a Random Forest (5pts) [W]

In machine learning, hyperparameters are parameters that are set before the learning process begins. The max_depth, num_estimators, or max_features variables from 3.1 are examples of different hyperparameters for a random forest model. In this section, you will tune your random forest model on an water potability dataset to achieve a decently high accuracy on classifying whether or not water is safe to drink.

Let's first review the dataset in a bit more detail.

Dataset Objective

Imagine that we are data scientists working for the United States Environmental Protection Agency (EPA) and we are given the task to predict whether or not a water sample is safe to drink. Water supplies are commonly contaminated with all sorts of chemicals and the EPA typically sets standards and regulations for the presence and levels for these contaminants.

After much deliberation amongst the team, you come to a conclusion that we can use data from previous lab reports on different water samples to predict whether or not a new sample of unidentified water can be drinkable.

We will use our random forest algorithm from Q3.1 to predict if a water sample is potent.

Loading the dataset

The dataset that the EPA has collected has the following features:

(each feature is followed by a quick description of what it measures)

Inputs:

1. pH value - Indicator of acidic or alkaline condition of water status. The current value ranges from 6.52 to 6.83 which are in the range of WHO standards.
2. Hardness - The capacity of water to precipitate soap caused by Calcium and Magnesium.
3. Solids - Amount of mineralization in the water.
4. Chloramines - Amount of chlorine and chloramine in the water.
5. Sulfate - Amount of sulfate concentration in the water.
6. Conductivity - Amount of ion concentration in the water. According to WHO standards, EC value should not exceed 400 $\mu\text{S}/\text{cm}$.
7. Trihalomethanes - THMs are chemicals which may be found in water treated with chlorine. The concentration of THMs in drinking water varies according to the level of organic material in the water and the temperature of the water.
8. Turbidity - Measure of light emitting properties of water and the test is used to indicate the quality of waste discharge with respect to colloidal matter.

Output:

Potability:

- 0 means not safe for humans to drink
- 1 means safe for humans to drink

Your random forest model will try to predict this variable.

```
In [6]: 1 #####
2 ### DO NOT CHANGE THIS CELL ####
3 #####
4 from sklearn.model_selection import StratifiedShuffleSplit
5
6 import pandas as pd
7 import numpy as np
8
9 water_data_train = pd.read_csv("data/water_potability_data.csv")
10 water_x_train = water_data_train.iloc[:, :-1]
11 water_y_train = water_data_train.iloc[:, -1]
12
13 water_x_train = np.array(water_x_train)
14 water_y_train = np.array(water_y_train)
```

```
In [7]: 1 #####
2 ### DO NOT CHANGE THIS CELL ####
3 #####
4 sss = StratifiedShuffleSplit(n_splits=2, test_size=0.2, random_state=0)
5 for train_index, test_index in sss.split(water_x_train, water_y_train):
6     X_train, X_test = water_x_train[train_index], water_x_train[test_index]
7     y_train, y_test = water_y_train[train_index], water_y_train[test_index]
```

```
In [8]: 1 X_train.shape, y_train.shape, X_test.shape, y_test.shape
```

```
Out[8]: ((2620, 9), (2620,), (656, 9), (656,))
```

In the following codeblock, train your random forest model with different values for max_depth, n_estimators, or max_features and evaluate each model on the held-out test set. Try to choose a combination of hyperparameters that maximizes your prediction accuracy on the test set (aim for 63%+). **Once you are satisfied with your chosen parameters, change the default values for max_depth, n_estimators, and max_features in the init function of your RandomForest class in random_forest.py to your chosen values, and then submit this file to Gradescope.

Expected Result: You must achieve at least a 63% accuracy against the test set to receive full credit for this section.**

The following is the breakdown of points for partial credit:

- **57% to 60%** earns 2pts
- **60% to 63%** earns 1.5pts more (3.5pts total)
- **63%+** earns 1.5pts more (5pts total)

```
In [18]: 1 """
2 TODO:
3 n_estimators defines how many decision trees are fitted for the random forest.
4 max_depth defines a stop condition when the tree reaches to a certain depth.
5 max_features controls the percentage of features that are used to fit each decision tree.
6
7 Tune these three parameters to achieve a better accuracy. While you can use the provided test set to
8 evaluate your implementation, you will need to obtain 63% on the test set to receive full credit
9 for this section.
10 """
11 from random_forest import RandomForest
12 import sklearn.ensemble
13 n_estimators = 7 #Hint: Consider values between 3-15.
14 max_depth = 12 # Hint: Consider values between 5-15
15 max_features = 0.7 # Hint: Consider values between 0.6-1.0.
16
17 random_forest = RandomForest(n_estimators, max_depth, max_features)
18
19 random_forest.fit(X_train, y_train)
20
21 accuracy=random_forest.oob_score(X_test, y_test)
22
23 print("accuracy: %.4f" % accuracy)
```

```
accuracy: 0.5801
```

3.3 Plotting Feature Importance (5pts) [W]

While building tree-based models, it's common to quantify how well splitting on a particular feature in a decision tree helps with predicting the target label in a dataset. Machine learning practitioners typically use 'Gini importance', or the (normalized) total reduction in entropy brought by that feature to evaluate how important that feature is for predicting the target variable.

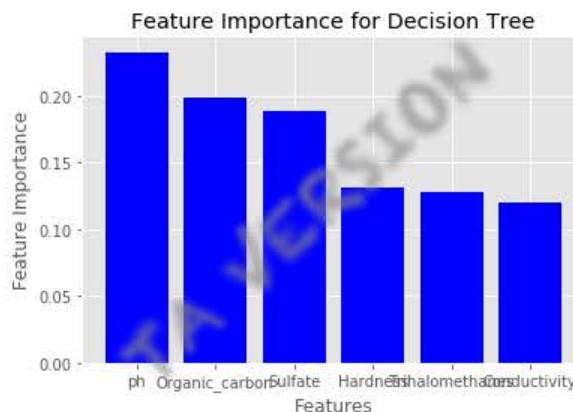
Gini importance is typically calculated as the reduction in entropy from reaching a split in a decision tree weighted by the probability of reaching that split in the decision tree. Sklearn internally computes the probability for reaching a split by finding the total number of samples that reaches it during the training phase divided by the total number of samples in the dataset. This weighted value is our feature importance.

Let's think about what this metric means with an example. A high probability of reaching a split on 'pH' in a decision tree trained on our water dataset (many samples will reach this split for a decision) and a large reduction in entropy from splitting on 'pH' will result in a high feature importance value for 'pH'. This could mean 'pH' is a very important feature for predicting the level of potability for water. On the other hand, a low probability of reaching a split on 'Sulfate' in a decision tree (few samples will reach this split for a decision) and a low reduction in entropy from splitting on 'Sulfate' will result in a low feature importance value. This could mean 'Sulfate' is not a very informative feature for predicting the level of potability in our decision tree. **Thus, the higher the feature importance value, the more important the feature is to predicting the target label.**

Fortunately for us, fitting a `sklearn.DecisionTreeClassifier` to a dataset automatically computes the Gini importance for every feature in the decision tree and stores these values in a `featureimportances` variable. [Review the docs for more details on how to access this variable \(\[https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html#sklearn.tree.DecisionTreeClassifier.feature_importances_\]\(https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html#sklearn.tree.DecisionTreeClassifier.feature_importances_\)\)](https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html#sklearn.tree.DecisionTreeClassifier.feature_importances_)

In the function below, display a bar plot that shows the feature importance values for at least one decision tree in your tuned random forest from Q3.2, and briefly comment on whether any features have noticeably higher / or lower importance weights than others. [Note that there isn't a 'correct' answer here. We simply want you to investigate how different features in your random forest contribute to predicting the target variable].

```
In [19]: 1 # TODO: Complete plot_feature_importance() in random_forest.py
2
3 random_forest.plot_feature_importance(water_data_train)
```



3.4 Improvement (5pts) [W]

For this question, we only ask that the Random Forest model have an accuracy of atleast 63%. In the real world, we would not be satisfied with this accuracy and would most likely not immediately deploy this model to determine the potability of real samples of water. What are some potential causes for why a Random Forest decision tree model may not produce high accuracies and what are some ways to improve on these potential causes?

4: SVM (30 Pts) Bonus of all [W] [P]

4.1 Fitting an SVM classifier by hand (20 Pts) [W]

Consider a dataset with the following points in 2-dimensional space:

x_1	x_2	y
0	1	1
2	2	-1
1	2	1

x_1	x_2	y
3	1	-1
-1	2	1
1	1	-1

Here, x_1 and x_2 are features and y is the label.

The max margin classifier has the formulation,

$$\min \frac{1}{2} \|\theta\|^2$$

$$s.t. y_i(\mathbf{x}_i \theta + b) \geq 1 \quad \forall i$$

Hint: \mathbf{x}_i are the support vectors. Margin is equal to $\frac{1}{\|\theta\|}$ and full margin is equal to $\frac{2}{\|\theta\|}$. You might find it useful to plot the points in a 2D plane.

- (1) Are the points linearly separable? Does adding the point $\mathbf{x} = (2, 1)$, $y = 1$ change the separability? (4 pts)
- (2) Write the lagrangian associated to the optimization problem. Can you apply the KKT conditions to obtain theta? (4pt)
- (3) Explain why only the data point on the margin will contribute to θ (4pt)
- (4) How would your lagrangian change if we want to kernelize the SVM? (4pt)
- (5) let's change the optimization problem to Soft Margin

$$\min \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^N \xi_i$$

$$s.t. y_i(\mathbf{x}_i \theta + b) \geq 1 - \xi_i \quad \forall i$$

Explain what is the role of C and how does it impact the final solution. (4pt)

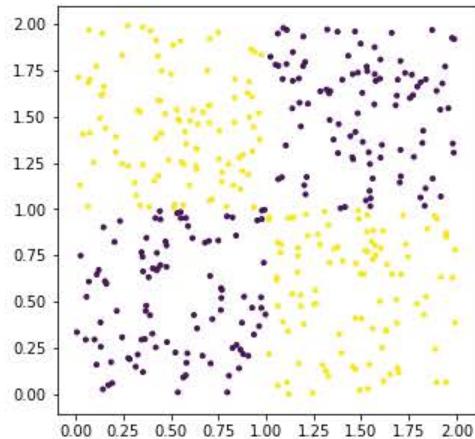
4.2 Feature Mapping (10 Pts) [P]

Let's look at a dataset where the datapoint can't be classified with a good accuracy using a linear classifier. Run the below cell to generate the dataset.

We will also see what happens when we try to fit a linear classifier to the dataset.

In [5]:

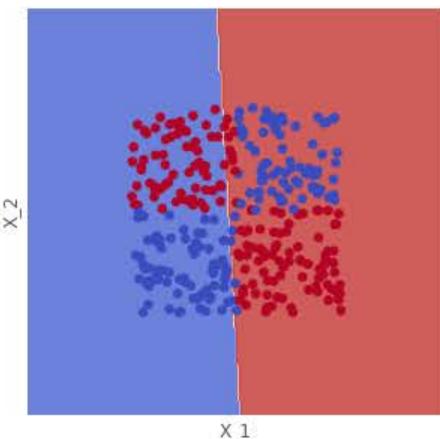
```
1 #####  
2 ### DO NOT CHANGE THIS CELL ###  
3 #####  
4 # Generate dataset  
5  
6 random_state = 1  
7  
8 X_1 = np.random.uniform(size=(100, 2))  
9 y_1 = np.zeros((100,)) - 1  
10  
11 X_2 = np.random.uniform(size=(100, 2))  
12 X_2[:, 0] = X_2[:, 0] + 1.0  
13 y_2 = np.ones((100,))  
14  
15 X_3 = np.random.uniform(size=(100, 2))  
16 X_3[:, 1] = X_3[:, 1] + 1.0  
17 y_3 = np.ones((100,))  
18  
19 X_4 = np.random.uniform(size=(100, 2))  
20 X_4[:, 0] = X_4[:, 0] + 1.0  
21 X_4[:, 1] = X_4[:, 1] + 1.0  
22 y_4 = np.zeros((100,)) - 1  
23  
24 X = np.concatenate([X_1, X_2, X_3, X_4], axis=0)  
25 y = np.concatenate([y_1, y_2, y_3, y_4], axis=0)  
26 X_train, X_test, y_train, y_test = train_test_split(X, y,  
27 test_size=0.20,  
28 random_state=random_state)  
29  
30 f, ax = plt.subplots(nrows=1, ncols=1, figsize=(5,5))  
31 plt.scatter(X[:, 0], X[:, 1], c = y, marker = '.')  
32 plt.show()
```



```
In [32]: 1 ##########
2 ### DO NOT CHANGE THIS CELL #####
3 #####
4
5 def visualize_decision_boundary(X, y, feature_new=None, h=0.02):
6     """
7         You don't have to modify this function
8
9     Function to vizualize decision boundary
10
11     feature_new is a function to get X with additional features
12     """
13     x1_min, x1_max = X[:, 0].min() - 1, X[:, 0].max() + 1
14     x2_min, x2_max = X[:, 1].min() - 1, X[:, 1].max() + 1
15     xx_1, xx_2 = np.meshgrid(np.arange(x1_min, x1_max, h),
16                             np.arange(x2_min, x2_max, h))
17
18     if X.shape[1] == 2:
19         z = svm_cls.predict(np.c_[xx_1.ravel(), xx_2.ravel()])
20     else:
21         X_cono = np.c_[xx_1.ravel(), xx_2.ravel()]
22         X_new = feature_new(X_cono)
23         z = svm_cls.predict(X_new)
24     z = z.reshape(xx_1.shape)
25
26     f, ax = plt.subplots(nrows=1, ncols=1, figsize=(5, 5))
27     plt.contourf(xx_1, xx_2, z, cmap=plt.cm.coolwarm, alpha=0.8)
28     plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.coolwarm)
29     plt.xlabel('X_1')
30     plt.ylabel('X_2')
31     plt.xlim(xx_1.min(), xx_1.max())
32     plt.ylim(xx_2.min(), xx_2.max())
33     plt.xticks(())
34     plt.yticks(())
35
36     plt.show()
```

```
In [33]: 1 #####
2 ### DO NOT CHANGE THIS CELL #####
3 #####
4 # Try to fit a linear classifier to the dataset
5
6 svm_cls = svm.LinearSVC()
7 svm_cls.fit(X_train, y_train)
8 y_test_predicted = svm_cls.predict(X_test)
9
10 print("Accuracy on test dataset: {}".format(accuracy_score(y_test,
11                                         y_test_predicted)))
12
13 visualize_decision_boundary(X_train, y_train)
```

Accuracy on test dataset: 0.45



We can see that we need a non-linear boundary to be able to successfully classify data in this dataset. By mapping the current feature x to a higher space with more features, linear SVM could be performed on the features in the higher space to learn a non-linear decision boundary. In the function below add additional features which can help classify in the above dataset. After creating the additional features

use code in the further cells to see how well the features perform on the test set.

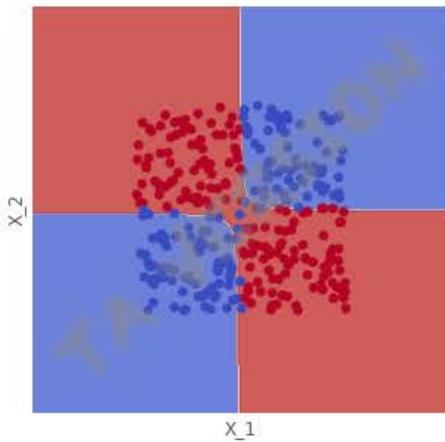
Note: You should get an accuracy above 95% correction, 90% not 95% for full credit

Hint: Think of the shape of the decision boundary that would best separate the above points. What additional features could help map the linear boundary to the non-linear one? Look at [this](https://xavierbourretsicotte.github.io/Kernel_feature_map.html) for a detailed analysis of doing the same for points separable with a circular boundary

```
In [34]: 1 #####  
2 ### DO NOT CHANGE THIS CELL ###  
3 #####  
4 from feature import create_nl_feature  
5  
6 X_new = create_nl_feature(X)  
7 X_train, X_test, y_train, y_test = train_test_split(X_new, Y,  
8 test_size=0.20,  
9 random_state=random_state)
```

```
In [35]: 1 #####  
2 ### DO NOT CHANGE THIS CELL ###  
3 #####  
4 # Fit to the new features and vizualize the decision boundary  
5 # You should get more than 90% accuracy on test set  
6  
7 svm_ols = svm.LinearSVC()  
8 svm_ols.fit(X_train, y_train)  
9 y_test_predicted = svm_ols.predict(X_test)  
10  
11 print("Accuracy on test dataset: {}".format(accuracy_score(y_test, y_test_predicted)))  
12  
13 visualize_decision_boundary(X_train, y_train, create_nl_feature)
```

Accuracy on test dataset: 0.975



```
In [ ]: 1
```