

Fall 2020 CS 4641/7641 A: Machine Learning Homework 4

Instructor: Dr. Mahdi Roozbahani

Deadline: November 18th, Wednesday, AOE

- No unapproved extension of the deadline is allowed. Late submission will lead to 0 credit.
- Discussion is encouraged on Piazza as part of the Q/A. However, all assignments should be done individually.

Instructions for the assignment

- In this assignment, we have programming and writing questions.
- To switch between cell for code and for markdown, see the menu-> Cell-> Cell Type
- You could directly type the latex equations in the markdown cell.
- Typing with latex/markdown is required for all the written questions. Handwritten answers would not be accepted.
- If a question requires a picture, you could use this syntax `<img src="" style="width: 300px;"/>` to include them within your ipython notebook.
- Questions marked with "[PI]" are programming only and should be submitted to the autograder. Questions marked with "[WI]" may require that you code a small function or generate plots, but should NOT be submitted to the autograder. It should be submitted on the writing portion of the assignment on gradescope.
- The kernel of the bonus is as follows:
  - Q1 ([5(+10) points for undergrads])> Neural Network "[WI]" | "[PI]"
  - Q2 ([15 pts(bonus for all)])> Image Classification based on Convolutional Neural Network <span> "[WI]"
  - Q3 ([40 pts])> Random Forest "[PI]" 3.1, 3.2 | "[WI]" 3.3
  - Q4 ([30 pts])> SVM "[WI]"

Using the autograder

- You will find two assignments on Gradescope that correspond to HW4: "HW4 - Programming" and "HW4 - Non-programming" (and "HW4 - Bonus Programming" if you are in CS4641).
- You will submit your code for the autograder on "HW4 - Programming" in the following format:
  - random\_forest.py
  - neural\_network.py
- You will submit your code for the autograder on "HW4 - Bonus-Programming" in the following format:
  - neural\_network.py
- All you will have to do is to copy your implementations of the classes "dinet" and "RandomForest" onto the respective files. 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 have 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 "HW4 - Non-programming" part, you will download your jupyter notebook as HTML, print it as a PDF from your browser and submit it on the autograder. To download the notebook as html, click on "File" on the top left corner of this page and select "Download as HTML". The non-programming part corresponds to Q1, Q2, Q3,3, Q4

Environment Setup

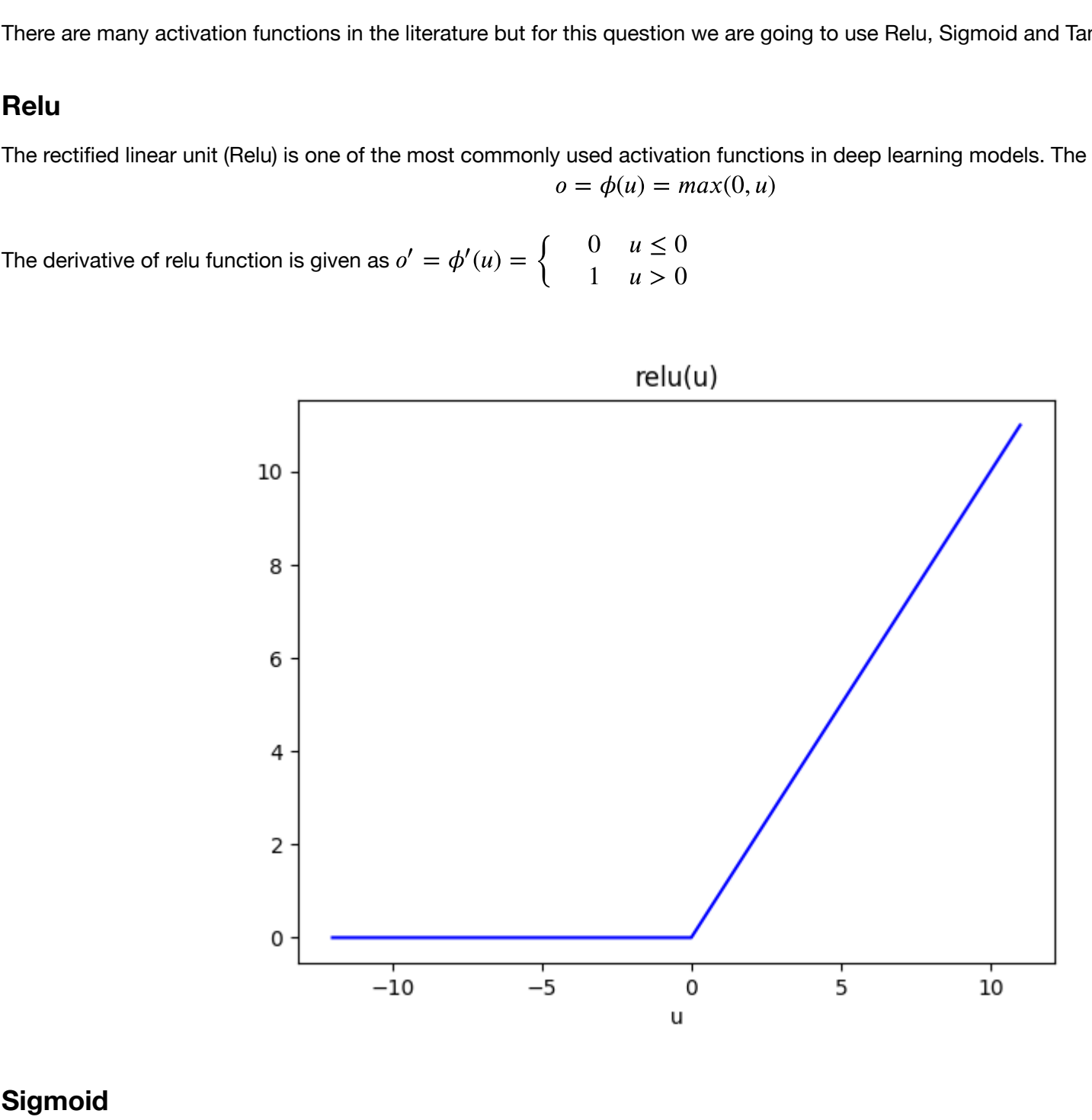
```
In [1]: import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load_breast_cancer
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report
from sklearn.metrics import plot_confusion_matrix

from collections import Counter
from scipy import stats
from math import log2, sqrt
import pandas as pd
import time

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
from sklearn.tree import DecisionTreeClassifier
from sklearn.datasets import make_moons
from sklearn.metrics import roc_auc_score
from sklearn import svm
```

1. Two Layer Neural Network [65pts] \*\*[PI]\*\*

Perceptron



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

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

where  $x$  is a  $d$ -dimensional vector i.e.  $x \in \mathbb{R}^d$ . It is one datapoint with  $d$  features.  $\theta_i \in \mathbb{R}^d$  is the weight vector for the  $i^{th}$  hidden unit,  $b_i \in \mathbb{R}$  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  $x$  weighted by  $\theta_i$ , 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 layer 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  $l-1$ .

where  $u^{[l]} \in \mathbb{R}^m$  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.  $\theta^{[l]} \in \mathbb{R}^{m \times n}$  is the weight matrix of the  $l^{th}$  layer where each row of  $\theta^{[l]}$  is analogous to  $\theta_i$  described in the previous section i.e. each row corresponds to one hidden unit of the  $l^{th}$  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  $l^{th}$  layer. This is followed by element wise non-linear activation function  $\phi^{[l]} = \phi(u^{[l]})$ . 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.

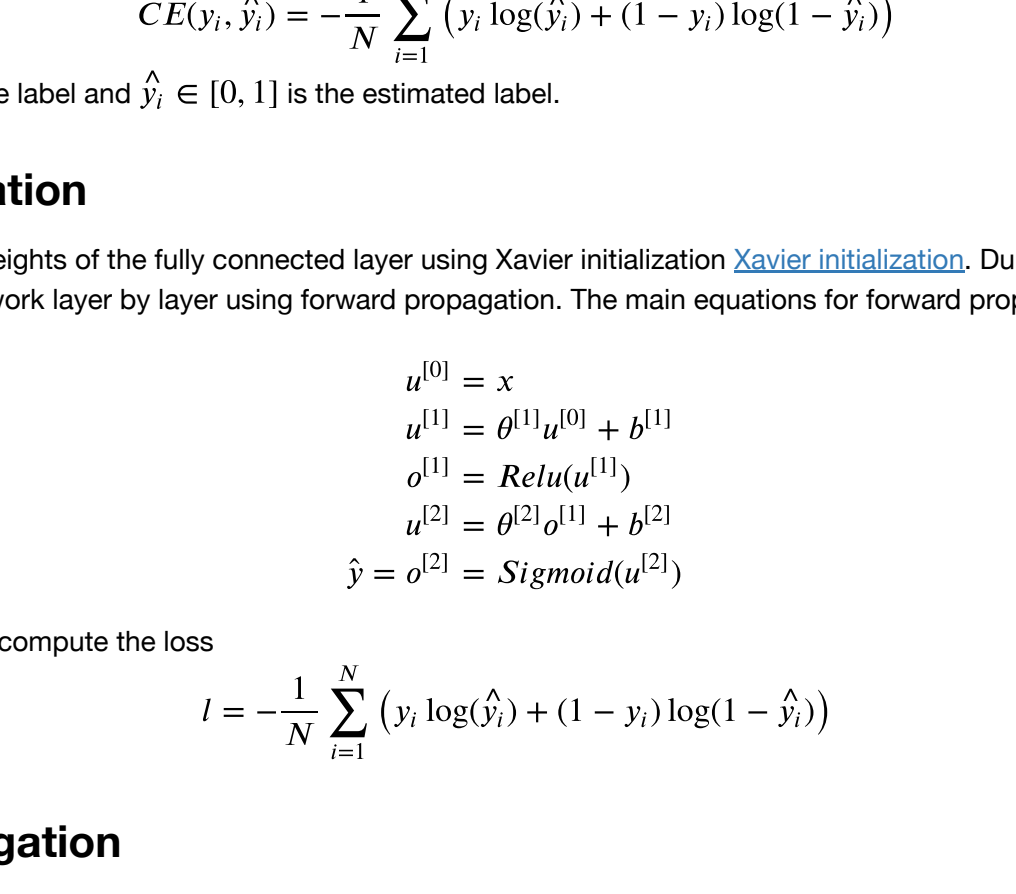
Activation Function

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

Relu

The rectified linear unit (Relu) is one of the most commonly used activation functions in deep learning models. The mathematical form is 
$$o = \phi(u) = \max(0, u)$$

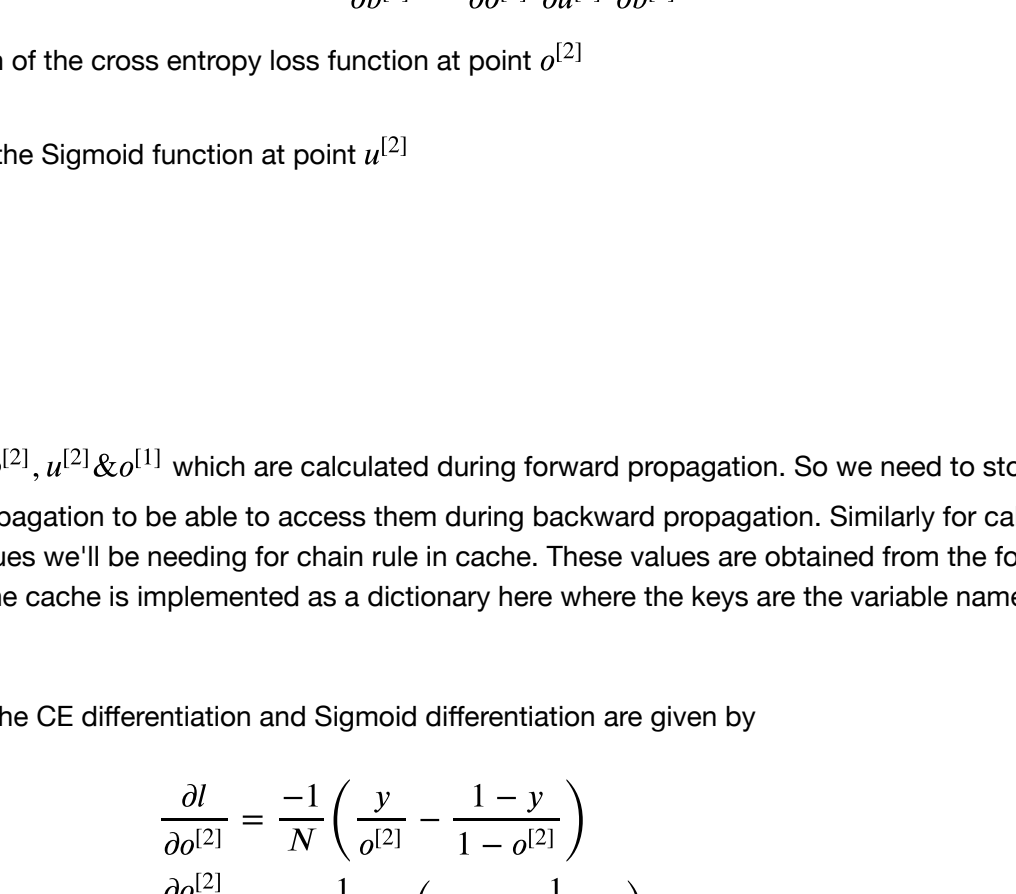
The derivative of relu function is given as 
$$\phi'(u) = \begin{cases} 0 & u \leq 0 \\ 1 & u > 0 \end{cases}$$



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

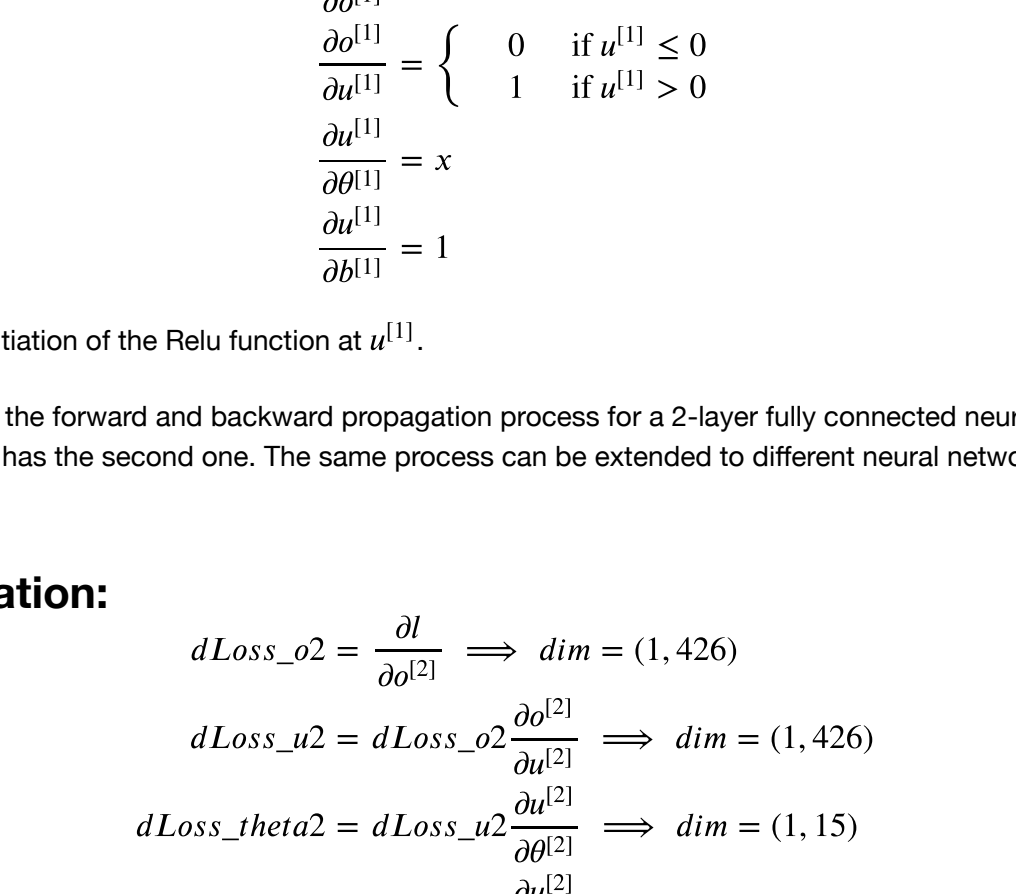
The derivative of the sigmoid function has a nice form and is given as 
$$\phi'(u) = \phi(u) \left( 1 - \frac{1}{1 + e^{-u}} \right) = \phi(u)(1 - \phi(u))$$



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

The derivative of tanh is given as 
$$\phi'(u) = \phi(u) \left( 1 - \left( \frac{e^u - e^{-u}}{e^u + e^{-u}} \right)^2 \right) = 1 - o^2$$



Cross Entropy Loss

An essential piece in training a neural network is the loss function. The main purpose of gradient descent algorithm is to find some network parameters that minimizes the loss function. In this exercise, we minimize Cross Entropy (CE) loss that represents on an intuitive level the distance between true data distribution and estimated distribution by neural network. So during training of the neural network, we will be looking for network parameters that minimize the distance between true and estimated distribution. The mathematical form of the CE loss is given by 
$$CE(p, q) = - \sum p(x_i) \log q(x_i)$$

where  $p(x)$  is the true distribution and  $q(x)$  is the estimated distribution.

Implementation details

For binary classification problems as in this exercise, we have probability distribution of a label  $y_i$  given by 
$$y_i = \begin{cases} 1 & \text{with probability } p(x_i) \\ 0 & \text{with probability } 1 - p(x_i) \end{cases}$$

A frequentist estimate of  $p(x_i)$  can be written as 
$$p(x_i) = \frac{\sum_{j=1}^N y_j}{N}$$

Therefore, the cross entropy for binary estimation can be written as 
$$CE(y_i, \hat{y}_i) = - \frac{1}{N} \sum_{i=1}^N (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$$

where  $y_i \in \{0, 1\}$  is the true label and  $\hat{y}_i \in [0, 1]$  is the estimated label.

Forward Propagation

We start by initializing the weights of the fully connected layer using Xavier initialization [Xavier initialization](#). 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.

$$u^{[0]} = x$$
$$u^{[1]} = \theta^{[0]} u^{[0]} + b^{[1]}$$
$$u^{[1]} = \text{Relu}(u^{[1]})$$
$$u^{[2]} = \theta^{[1]} u^{[1]} + b^{[2]}$$
$$\hat{y} = o^{[2]} = \text{Sigmoid}(u^{[2]})$$

Then we get the output and compute the loss

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

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

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

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:

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

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

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

$\frac{\partial l}{\partial u^{[1]}}$  is equal to  $\phi^{[1]}$

$\frac{\partial l}{\partial u^{[0]}}$  is equal to 1.

To compute  $\frac{\partial l}{\partial \theta^{[1]}}$ , we need  $\frac{\partial l}{\partial u^{[2]}}$ ,  $u^{[2]}$  &  $\phi^{[1]}$  which are calculated during forward propagation. So we need to store these values in cache variables during forward propagation so we are 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 CE differentiation and Sigmoid differentiation are given by

$$\frac{\partial l}{\partial \theta^{[2]}} = \frac{-1}{N} \left( \frac{y}{o^{[2]}} - \frac{1-y}{1-o^{[2]}} \right)$$
$$\frac{\partial l}{\partial \theta^{[2]}} = \frac{1}{1 + e^{-u^{[2]}}} \left( 1 - \frac{1}{1 + e^{-u^{[2]}}} \right) = o^{[2]}(1 - o^{[2]})$$
$$\frac{\partial l}{\partial b^{[2]}} = \phi^{[1]}$$
$$\frac{\partial l}{\partial u^{[2]}} = 1$$

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:

$$\frac{\partial l}{\partial \theta^{[1]}} = \frac{\partial l}{\partial u^{[2]}} \frac{\partial u^{[2]}}{\partial \theta^{[1]}} \frac{\partial u^{[2]}}{\partial u^{[1]}} \frac{\partial u^{[1]}}{\partial \theta^{[1]}}$$
$$\frac{\partial l}{\partial b^{[1]}} = \frac{\partial l}{\partial u^{[2]}} \frac{\partial u^{[2]}}{\partial b^{[1]}} \frac{\partial u^{[2]}}{\partial u^{[1]}} \frac{\partial u^{[1]}}{\partial b^{[1]}}$$

Where

$$\frac{\partial u^{[2]}}{\partial u^{[1]}} = \theta^{[2]}$$
$$\frac{\partial u^{[1]}}{\partial u^{[0]}} = \begin{cases} 0 & \text{if } u^{[1]} \leq 0 \\ 1 & \text{if } u^{[1]} > 0 \end{cases}$$
$$\frac{\partial u^{[1]}}{\partial u^{[0]}} = x$$
$$\frac{\partial u^{[1]}}{\partial b^{[1]}} = 1$$

Note that  $\frac{\partial u^{[1]}}{\partial u^{[0]}}$  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 relu as the first activation layer and sigmoid has the second one. The same process can be extended to different neural networks with different activation layers like tanh.

Code Implementation:

$$dLoss_o2 = \frac{\partial l}{\partial o^{[2]}} \implies dim = (1, 426)$$
$$dLoss_u2 = dLoss_o2 \frac{\partial o^{[2]}}{\partial u^{[2]}} \implies dim = (1, 426)$$
$$dLoss_theta2 = dLoss_u2 \frac{\partial u^{[2]}}{\partial \theta^{[2]}} \implies dim = (1, 15)$$
$$dLoss_b2 = dLoss_u2 \frac{\partial u^{[2]}}{\partial b^{[2]}} \implies dim = (1, 1)$$
$$dLoss_o1 = dLoss_u2 \frac{\partial u^{[2]}}{\partial u^{[1]}} \implies dim = (15, 426)$$
$$dLoss_u1 = dLoss_o1 \frac{\partial u^{[1]}}{\partial u^{[0]}} \implies dim = (15, 426)$$
$$dLoss_theta1 = dLoss_u1 \frac{\partial u^{[1]}}{\partial \theta^{[1]}} \implies dim = (15, 30)$$
$$dLoss_b1 = dLoss_u1 \frac{\partial u^{[1]}}{\partial b^{[1]}} \implies dim = (15, 1)$$

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 three activation functions here - Relu, Tanh and Sigmoid. You will implement a neural network that could have relu activation followed by sigmoid layer or tanh activation followed by sigmoid. You'll have to specify the neural net type which could be "Relu -> Sigmoid" (set by default) or "Tanh -> Sigmoid".

You'll also implement gradient descent and stochastic gradient descent algorithms for training these neural nets. SGD is bonus for undergraduate students.

We'll train these neural nets on breast cancer dataset. You're free to choose either gradient descent or SGD for training. 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.

Deliverables for this question:

1. Loss plot and classification report for any neural net type ("Relu -> Sigmoid" or "Tanh -> Sigmoid") with gradient descent
2. Loss plot and classification report for any neural net type ("Relu -> Sigmoid" or "Tanh -> Sigmoid") with stochastic gradient descent (mandatory for graduate students, bonus for undergraduate students)

```
In [6]: '''
We are going to use Breast Cancer Wisconsin (Diagnostic) Data Set provided by sklearn
https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_breast_cancer.html
to train a 2 fully connected layer neural net. We are going to build the neural network from scratch.
'''

class dinet:
    def __init__(self, x, y, lr = 0.003):
        '''
        This method initializes the class, its implemented for you.
        Args:
            x: data
            y: labels
            Yh: predicted labels
            dims: dimensions of different layers
            param: dictionary of different layers parameters
            chi: Cache dictionary to store forward parameters that are used in backpropagation
            loss: list to store loss values
            lr: learning rate
            sam: number of training samples we have

            self.Xxx # features
            self.Yy # ground truth labels

            self.Yh=np.zeros((1,self.Y.shape[1])) # estimated labels
            self.dims = [30, 15, 1] # dimensions of different layers

            self.param = {} # dictionary for different layer variables
            self.ch = {} # cache for holding variables during forward propagation to use them in back prop
            self.loss = []

            self.lr=lr # learning rate
            self.sam = self.Y.shape[1] # number of training samples we have
            self.estimator_type = 'classifier'
            self.neural_net_type = "Tanh -> Sigmoid" #can change it to "Tanh -> Sigmoid"

        def ninit(self):
            '''
            This method initializes the neural network variables, its already implemented for you.
            Check it and relate to mathematical description above.
            We are going to use these variables in forward and backward propagation.
            '''
            np.random.seed(1)
            self.param['theta2'] = np.random.randn(self.dims[1], self.dims[0]) / np.sqrt(self.dims[0])
            self.param['b2'] = np.zeros((self.dims[1], 1))
            self.param['theta1'] = self.param['theta2'] * 0.029403
            self.param['b1'] = np.zeros((self.dims[2], self.dims[1]))
            self.param['b2'] = np.zeros((self.dims[2], 1))

        def Relu(self, u):
            '''
            In this method you are going to implement element wise Relu.
            Make sure that all operations here are element wise and can be applied to an input of any dimension.
            Input: u of any dimension
            return: Relu(u)
            '''
            o = np.maximum(u, 0)
            return o

        def Sigmoid(self, u):
            '''
            In this method you are going to implement element wise Sigmoid.
            Make sure that all operations here are element wise and can be applied to an input of any dimension.
            Input: u of any dimension
            return: Sigmoid(u)
            '''
            o = 1 / (1 + np.exp(-u))
            return o

        def Tanh(self, u):
            '''
            In this method you are going to implement element wise Tanh.
            Make sure that all operations here are element wise and can be applied to an input of any dimension.
            Input: u of any dimension
            return: Tanh(u)
            '''
            o = (np.exp(u) - np.exp(-u)) / (np.exp(u) + np.exp(-u))
            return o

        def dRelu(self, u):
            '''
            This method implements element wise differentiation of Relu.
            Input: u of any dimension
            return: dRelu(u)
            '''
            u[u<=0] = 0
            u[u>0] = 1
            return u

        def dSigmoid(self, u):
            '''
            This method implements element wise differentiation of Sigmoid.
            Input: u of any dimension
            return: dSigmoid(u)
            '''
            o = 1/(1+np.exp(-u))
            do = o * (1-o)
            return do

        def dTanh(self, u):
            '''
            This method implements element wise differentiation of Tanh.
            Input: u of any dimension
            return: dTanh(u)
            '''
            o = np.tanh(u)
            return 1-o**2

        def nloss(self, y, yh):
            '''
            In this method you are going to implement Cross Entropy loss.
            Refer to the description above and implement the appropriate mathematical equation.
            Input: y 1xN: ground truth labels
            yh 1xN: neural network output after Sigmoid

            return: CE 1x1: loss value
            CE = -1 * np.mean((np.multiply(y, np.log(yh)) + np.multiply((1 - y), np.log(1 - yh))))

        def forward(self, x):
            '''
            Fill in the missing code lines, please refer to the description for more details.
            Check ninit method and variables from there as well as other implemented methods.
            Refer to the description above and implement the appropriate mathematical equations.
            do not change the lines followed by #keep.
            '''
            self.ch['X'] = x # keep
            u1 = np.dot(self.param['theta1'], self.ch['X']) + self.param['b1']
            if self.neural_net_type == "Relu -> Sigmoid":
                o1 = self.Relu(u1)
            else:
                o1 = self.Tanh(u1)
            self.ch['u1'], self.ch['o1'] = u1, o1 # keep
            u2 = np.dot(self.param['theta2'], self.ch['o1']) + self.param['b2']
            o2 = self.Sigmoid(u2)
            self.ch['u2'], self.ch['o2'] = u2, o2 # keep
            return o2 # keep

        def backward(self, y, yh):
            '''
            Fill in the missing code lines, please refer to the description for more details.
            You will need to use cache variables, as well as the implemented methods, and other variables as w
            Refer to the description above and implement the appropriate mathematical equations.
            do not change the lines followed by #keep.
            '''
            dloss_o2 = -(np.divide(y, yh) - np.divide(1 - y, 1 - yh)) / y.shape[1] # partial l by partial o2

            # Implement equations for getting derivative of loss w.r.t u2, theta2 and b2
            dloss_u2 = np.multiply(dloss_o2, np.multiply(self.ch['o2'], (1 - self.ch['o2']))) # partial l by partial u2
            dloss_theta2 = np.matmul(dloss_u2, self.ch['o1'].T) # partial l by partial theta2
            dloss_b2 = np.matmul(dloss_u2, np.ones((dloss_u2.shape[1], 1))) # partial l by partial b2

            # set dloss_u2, dloss_theta2, dloss_b2
            self.ch['dloss_u2'] = dloss_u2
            self.ch['dloss_theta2'] = dloss_theta2
            self.ch['dloss_b2'] = dloss_b2

            dloss_o1 = np.dot(self.param['theta2'].T, dloss_u2) # partial l by partial o1
            if self.neural_net_type == "Relu -> Sigmoid":
                dol_u1 = self.dRelu(self.ch['u1'])
            else:
                dol_u1 = self.dTanh(self.ch['u1'])
            dloss_u1 = np.multiply(dloss_o1, dol_u1) # partial l by partial u1
            dloss_theta1 = np.matmul(dloss_u1, self.ch['X'].T) # partial l by partial theta1
            dloss_b1 = np.matmul(dloss_u1, np.ones((dloss_u1.shape[1], 1))) # partial l by partial b1

            # set dloss_u1, dloss_theta1, dloss_b1
            self.ch['dloss_u1'] = dloss_u1
            self.ch['dloss_theta1'] = dloss_theta1
            self.ch['dloss_b1'] = dloss_b1

            # parameters update, no need to change these lines
            self.param['theta2'] = self.param['theta2'] - self.lr * dloss_theta2 # keep
            self.param['b2'] = self.param['b2'] - self.lr * dloss_b2 # keep
            self.param['theta1'] = self.param['theta1'] - self.lr * dloss_theta1 # keep
            self.param['b1'] = self.param['b1'] - self.lr * dloss_b1 # keep
            return dloss_theta2, dloss_b2, dloss_theta1, dloss_b1

        def gradient_descent(self, x, y, iter = 60000):
            '''
            This function is an implementation of the gradient descent algorithm
            '''
            self.ninit()
            for i in range(iter):
                yh = self.forward(x)
                loss = self.nloss(y, yh)
                self.loss.append(loss)
                if i % 2000 == 0:
                    print('Loss after iteration %d: %e' % (i+1, loss))
                    self.backward(y, yh)

            #bonus for undergraduate students
            def stochastic_gradient_descent(self, x, y, iter = 60000):
                '''
                This function is an implementation of the stochastic gradient descent algorithm
                Note:
                1. SGD loops over all examples in the dataset one by one and learns a gradient from each example.
                2. One iteration here is one round of forward and backward propagation on one example of the dataset.
                3. So if the dataset has 10000 examples, 1000 iterations will constitute an epoch.
                4. Append loss after every 2000 iterations for plotting loss plots.
                5. It is fine if you get a noisy plot since learning on one example at a time adds variance to the gradients learnt.
                6. You can use SGD with any neural net type
                '''
                self.ninit()
                for i in range(iter):
                    idx = i % x.shape[1]
                    xs = np.expand_dims(x[idx], axis=1)
                    ys = np.expand_dims(y[idx], axis=1)
                    yh = self.forward(xs)
                    if (i+1) % 2000 == 0:
                        loss = self.nloss(ys, yh)
                        print('Loss after iteration %d: %e' % (i+1, loss))
                        self.backward(ys, yh)

            def predict(self, x):
                '''
                This function predicts new data points
                Its implemented for you
                '''
                Yh = self.forward(x)
                return np.round(Yh).squeeze()

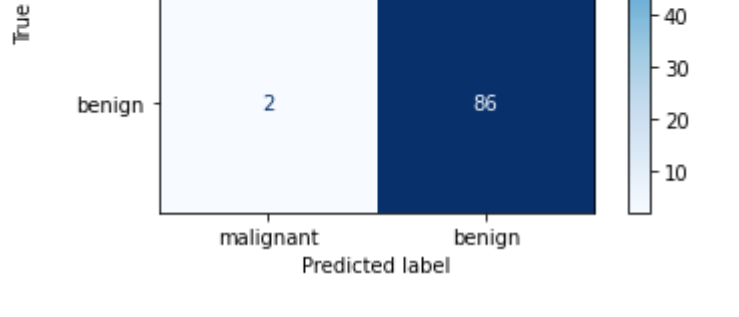
In [3]: '''
Training the Neural Network, you do not need to modify this cell
We are going to use Breast Cancer Wisconsin (Diagnostic) Data Set provided by sklearn
https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_breast_cancer.html
to train a 2 fully connected layer neural net. We are going to build the neural network from scratch.
'''

dataset = load_breast_cancer() # load the dataset
x, y = dataset.data, dataset.target
x = MinMaxScaler().fit_transform(x) #normalize data
X_train, X_test, y_train, y_test = train_test_split(x, y, random_state=1) #split data
X_train, X_test, y_train, y_test = X_train, y_test, X_train.reshape(1,-1), y_test #condition data
nn = dinet(X_train, y_train, lr=0.1) # initialize neural net class
nn.gradient_descent(X_train, y_train, iter = 60000) #train

# create figure
fig = plt.plot(nn.loss)
plt.title(f"Training: (nn.neural_net_type)")
plt.xlabel("Epoch")
plt.ylabel("Loss")

Loss after iteration 0: 0.699796
Loss after iteration 2000: 0.061302
Loss after iteration 4000: 0.052467
Loss after iteration 6000: 0.048680
Loss after iteration 8000: 0.046025
Loss after iteration 10000: 0.043946
Loss after iteration 12000: 0.042218
Loss after iteration 14000: 0.040690
Loss after iteration 16000: 0.039277
Loss after iteration 18000: 0.037926
Loss after iteration 20000: 0.036404
Loss after iteration 22000: 0.035215
Loss after iteration 24000: 0.033909
Loss after iteration 26000: 0.032485
Loss after iteration 28000: 0.030986
Loss after iteration 30000: 0.029403
Loss after iteration 32000: 0.027737
Loss after iteration 34000: 0.026004
Loss after iteration 36000: 0.024228
Loss after iteration 38000: 0.022442
Loss after iteration 40000: 0.020309
Loss after iteration 42000: 0.018405
Loss after iteration 44000: 0.016443
Loss after iteration 46000: 0.015034
Loss after iteration 48000: 0.013588
Loss after iteration 50000: 0.012028
Loss after iteration 52000: 0.011156
Loss after iteration 54000: 0.009915
Loss after iteration 56000: 0.008919
Loss after iteration 58000: 0.008131
Loss after iteration 60000: 0.007325
Loss after iteration 62000: 0.006696
Loss after iteration 64000: 0.006167

Out[3]: Text(0, 0.5, 'Loss')
```



```
In [4]: '''
Testing Neural Network
y_predicted = nn.predict(X_test) # predict

#plot
plt.plot(("Classification Report for (nn.neural_net_type)\n\n"))
print(classification_report(y_test, y_predicted, target_names=dataset.target_names))
plot_confusion_matrix(nn, X_test, y_test, cmap=plt.cm.Blues, display_labels=dataset.target_names)
plt.show()

Classification Report for Relu -> Sigmoid

              precision    recall  f1-score   support

malignant     0.96         0.95         0.95         55
benign         0.97         0.98         0.97         88

accuracy               0.97         143
macro avg              0.96         0.96         143
weighted avg           0.97         0.96         143

True label \ Predicted label
malignant      52      3
benign          2      86
```



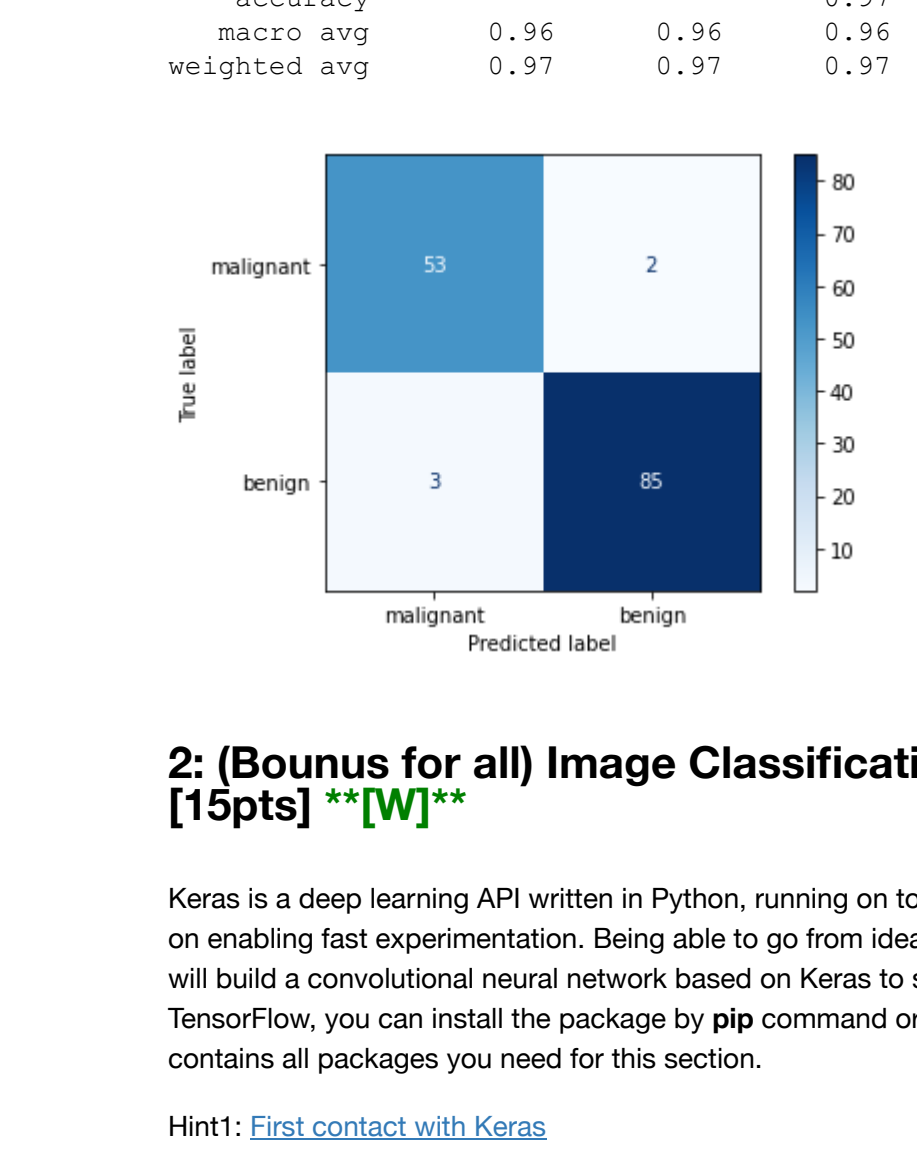
```
'''
Training the Neural Network, you do not need to modify this cell
We are going to use Breast Cancer Wisconsin (Diagnostic) Data Set provided by sklearn
https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_breast_cancer.html
'''

dataset = load_breast_cancer() # load the dataset
x, y = dataset.data, dataset.target
x = MinMaxScaler().fit_transform(x) #normalize data
x_train, x_test, y_train, y_test = train_test_split(x, y, random_state=1, split_data
x_train, x_test, y_train, y_test = x_train, x_test, y_train, x_test, y_train.reshape(1,-1), y_test #condition data

nn = dnn(x_train, y_train, lr=0.1) # initialize neural net class
nn.stochastic_gradient_descent(x_train, y_train, iter = 66000) #train

# create figure
fig = plt.plot(np.array(nn.loss).squeeze())
plt.title('Training: (nn.neural_net_type)')
plt.xlabel('Epochs')
plt.ylabel('Losses')

Loss after iteration 4000: 0.001772
Loss after iteration 4000: 0.000119
Loss after iteration 6000: 0.000119
Loss after iteration 8000: 0.000057
Loss after iteration 10000: 0.001398
Loss after iteration 12000: 0.000014
Loss after iteration 14000: 0.000057
Loss after iteration 16000: 0.000012
Loss after iteration 18000: 0.000004
Loss after iteration 20000: 0.007969
Loss after iteration 22000: 0.000218
Loss after iteration 24000: 0.001232
Loss after iteration 26000: 0.014399
Loss after iteration 28000: 0.000369
Loss after iteration 30000: 0.000063
Loss after iteration 32000: 0.001527
Loss after iteration 34000: 0.000177
Loss after iteration 36000: 0.000391
Loss after iteration 38000: 0.000000
Loss after iteration 40000: 0.000003
Loss after iteration 42000: 0.000000
Loss after iteration 44000: 0.001259
Loss after iteration 46000: 0.000075
Loss after iteration 48000: 0.000000
Loss after iteration 50000: 0.000017
Loss after iteration 52000: 0.000021
Loss after iteration 54000: 0.000017
Loss after iteration 56000: 0.019316
Loss after iteration 58000: 0.000158
Loss after iteration 60000: 0.000000
Loss after iteration 62000: 0.000437
Loss after iteration 64000: 0.000465
Loss after iteration 66000: 0.000111
```



```
In [8]: '''
Training Neural Network
'''
y_predicted = nn.predict(x_test) # predict

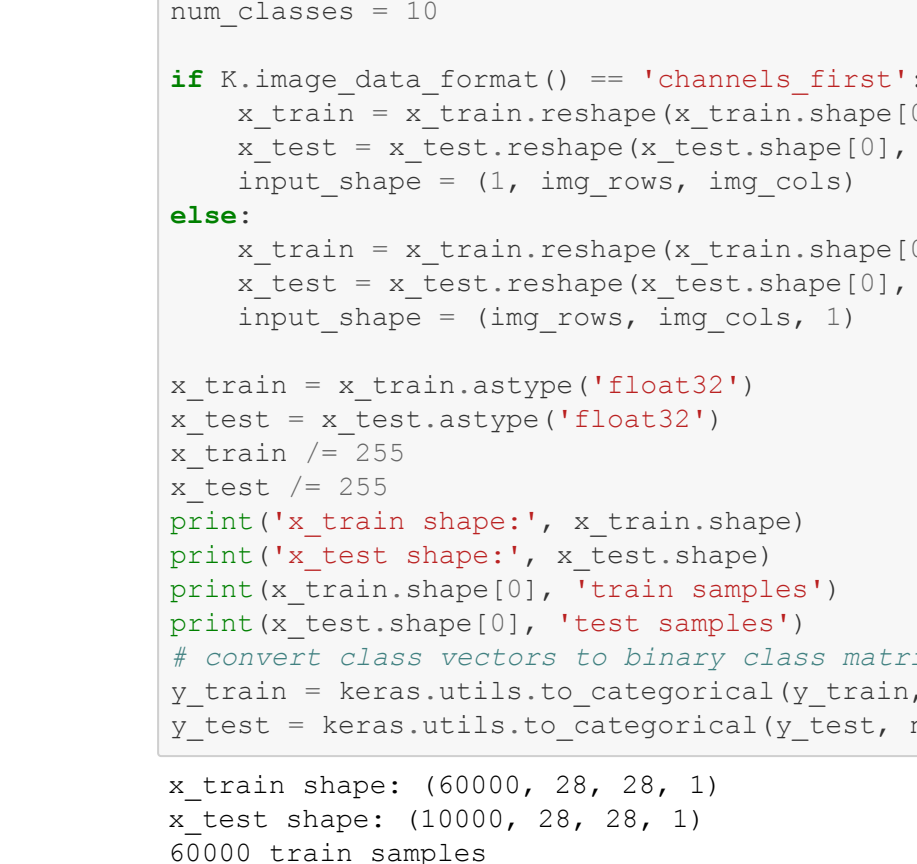
#plot
print(f"Classification Report for {nn.neural_net_type}\n\n")
print(classification_report(y_test, y_predicted, target_names=dataset.target_names))
plt.confusion_matrix(nn, x_test, y_test, cmap=plt.cm.Blues, display_labels=dataset.target_names)
plt.show()

Classification Report for Tanh -> Sigmoid

              precision    recall  f1-score   support

 malignant      0.95      0.96      0.95         55
  benign      0.98      0.97      0.97         88

 accuracy      0.96      0.96      0.97        143
 macro avg      0.96      0.96      0.96        143
 weighted avg      0.97      0.97      0.97        143
```



## 2: (Bouns for all) Image Classification based on Convolutional Neural Networks [15pts] \*\*[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 Keras to solve the image classification task for MNIST. If you haven't installed TensorFlow, you can install the package by pip command or train your model by uploading HW4 notebook to Colab directly. Colab contains all packages you need for this section.

Hint1: [First contact with Keras](#)

Hint2: [How to Install Keras](#)

Hint3: [CS231n Tutorial \(Layers used to build ConvNets\)](#)

### Environment Setup

```
In [9]: from future import print_function
import keras
from keras.datasets import mnist
from keras.models import Sequential
from keras.layers import Dense, Dropout, Flatten
from keras.layers import Conv2D, MaxPooling2D
from keras import backend as k
import matplotlib.pyplot as plt
```

### Load MNIST dataset

We use [MNIST](#) dataset to train our model. MNIST is a subset of a larger set available from NIST. MNIST database of handwritten digits has a training set of 60,000 examples, and a test set of 10,000 examples. Each example is 28 x 28 pixel grayscale image of handwritten digits between 0 to 9.

```
In [10]: # Helper function, You don't need to modify it
# split data between train and test sets
(x_train, y_train), (x_test, y_test) = mnist.load_data()

# input image dimensions
img_rows, img_cols = 28, 28
#set num of classes
num_classes = 10

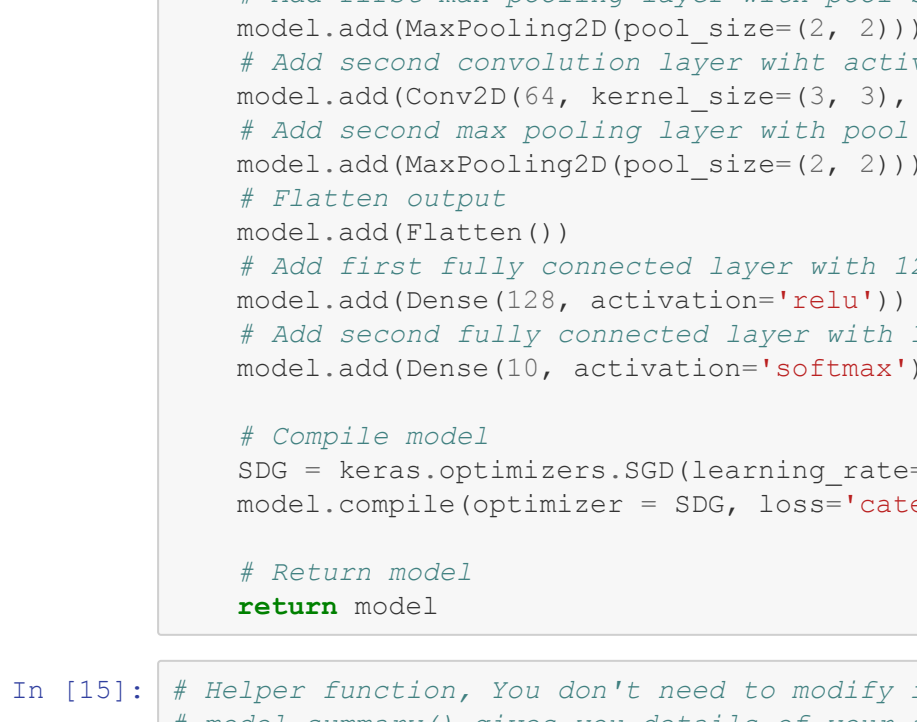
if K.image_data_format() == 'channels_first':
    x_train = x_train.reshape(x_train.shape[0], 1, img_rows, img_cols)
    x_test = x_test.reshape(x_test.shape[0], 1, img_rows, img_cols)
    input_shape = (1, img_rows, img_cols)
else:
    x_train = x_train.reshape(x_train.shape[0], img_rows, img_cols, 1)
    x_test = x_test.reshape(x_test.shape[0], img_rows, img_cols, 1)
    input_shape = (img_rows, img_cols, 1)

x_train = x_train.astype('float32')
x_test = x_test.astype('float32')
x_train /= 255
x_test /= 255
print('x_train shape:', x_train.shape)
print('x_test shape:', x_test.shape)
print(x_train.shape[0], 'train samples')
print(x_test.shape[0], 'test samples')
# convert class vectors to binary class matrices
y_train = keras.utils.to_categorical(y_train, num_classes)
y_test = 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
```

### Load some images from MNIST

```
In [11]: # Helper function, You don't need to modify it
# Show some images from MNIST
for i in range(10):
    plt.subplot(5,2,1+i)
    image=x_train[i].reshape((28,28))
    plt.imshow(image,cmap='gray')
plt.show()
```



As you can see from above, the MNIST dataset contains handwritten digits from 0 to 9. The digits have been size-normalized and centered in fixed-size images.

### Build convolutional neural network model

In this part, you need to build a convolutional neural network that contains 2 convolutional layers. The architecture of this model is:

**[INPUT - CONV - RELU - MAXPOOL - CONV - RELU - MAXPOOL - FC1 - FC2] [1]**

INPUT: [28 x 28 x 1] will hold the raw pixel values of the image, in this case, an image of width 28, height 28, and with only one color channels.

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 5x3 for the both CONV layers. For example, the output of the Conv. layer may like [26 x 26 x 32] if we use 32 filters.

RELU: As we mentioned in the previous section, the Relu layer will apply an elementwise activation function, such as the `max(0,x)` thresholding at zero, which leaves the size of the volume unchanged.

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

FC1: First Fully-Connected layer, we use **ReLU** as the activation function. The dimension of the output space is 128.

FC2: Second Fully-Connected layer will compute the class scores. We use **Softmax** as the activation function. The dimension of the output space is the number of class.

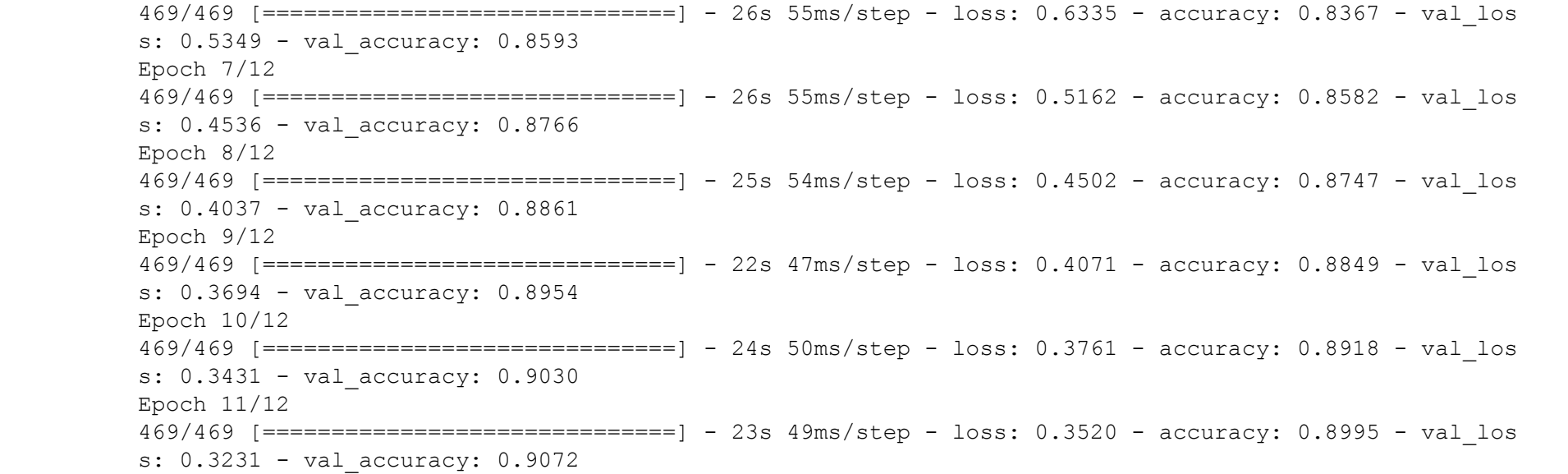
**Loss function:** Crosseentropy (mentioned in the previous section)

**optimizer:** Stochastic gradient descent(SGD)

[1] CS231n: <https://cs231n.github.io/convolutional-networks/>

```
In [12]: # Helper function, You don't need to modify it
# Show the architecture of the model
ach=plt.imshow(history.history['loss'])
fig = plt.figure(figsize=(10,10))
plt.imshow(ach)

Out[12]: <matplotlib.image.AxesImage at 0x7fa8d73ed1f0>
```



### Defining Variables

```
In [13]: # Defining Variables
# Do not change the value of num_classes.
# You can adjust of adding parameters to train your model
batch_size = 128
epochs = 12
lr = 0.001 #learning rate
```

### Defining model

```
In [14]: def create_net():
'''
In this function you are going to build a convolutional neural network based on Keras.
First, use Sequential() to build the inference model.
Then, use model.add() and model.compile() to build your own model.
Return: model you build
'''
# Initialize model
model = Sequential()

# Add first convolution layer with activation function Relu()
model.add(Conv2D(32, kernel_size=(3, 3), activation='relu', input_shape=(28, 28, 1)))
# Add first max pooling layer with pool size 2X2
model.add(MaxPooling2D(pool_size=(2, 2)))
# Add second convolution layer with activation function Relu()
model.add(Conv2D(64, kernel_size=(3, 3), activation='relu', input_shape=(13, 13, 1)))
# Add second max pooling layer with pool size 2X2
model.add(MaxPooling2D(pool_size=(2, 2)))
# Flatten output
model.add(Flatten())
# Add first fully connected layer with 128 nodes
model.add(Dense(128, activation='relu'))
# Add second fully connected layer with 10 nodes
model.add(Dense(10, activation='softmax'))

# Compile model
SGD = keras.optimizers.SGD(learning_rate=lr)
model.compile(optimizer = SGD, loss='categorical_crossentropy', metrics=['accuracy'])

# Return model
return model
```

```
In [15]: # Helper function, You don't need to modify it
# model.summary() gives you details of your architecture.
# You can compare your architecture with the 'Architecture.png'
model=create_net()
model.summary()

Model: "sequential"
Layer (type) Output Shape Param #
-----
conv2d_2 (Conv2D) (None, 26, 26, 32) 320
max_pooling2d_2 (MaxPooling2D) (None, 13, 13, 32) 0
conv2d_3 (Conv2D) (None, 11, 11, 64) 18496
max_pooling2d_3 (MaxPooling2D) (None, 5, 5, 64) 0
flatten_1 (Flatten) (None, 1600) 0
dense_2 (Dense) (None, 128) 204928
dense_3 (Dense) (None, 10) 1290
Total params: 225,034
Trainable params: 225,034
Non-trainable params: 0
```

### Train the network

**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 20 minutes to train your model.

**Expected Result:** You should be able to achieve more than 90% accuracy on the test set to get full 15 points. If you achieve accuracy between 80% to 90%, you will only get half points of this part.

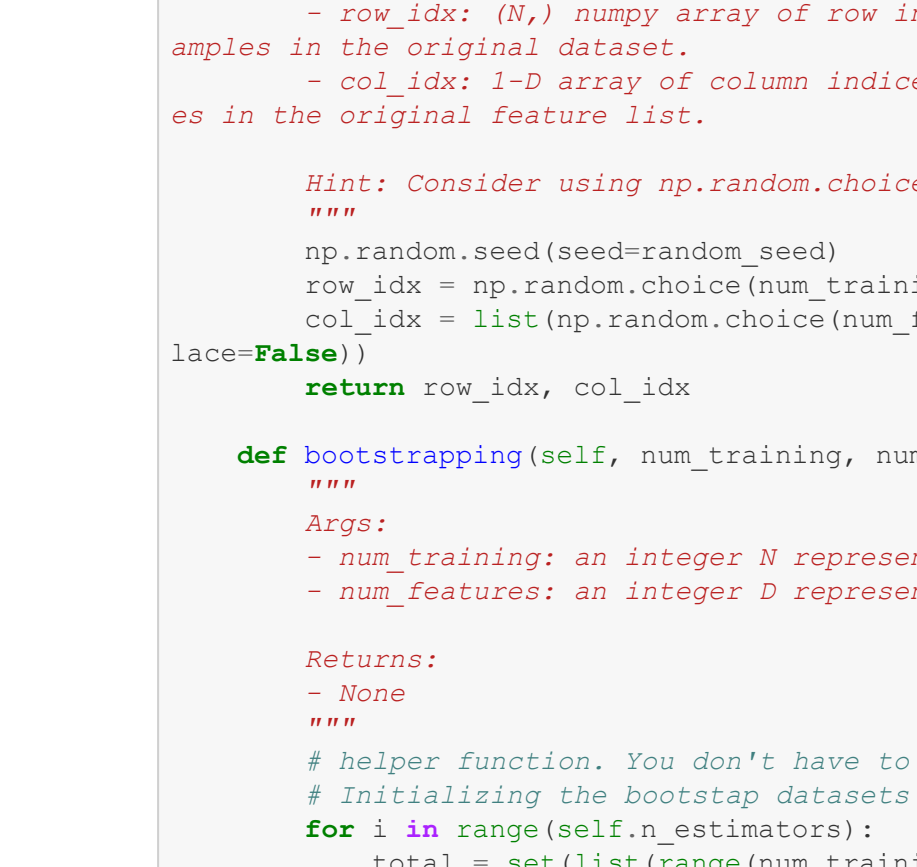
### Train your own CNN model

```
In [16]: # Helper function, You don't need to modify it
# Train the model.
history = model.fit(x_train, y_train,
                    batch_size=batch_size,
                    epochs=epochs,
                    verbose=1,
                    validation_data=(x_test, y_test))
score = model.evaluate(x_test, y_test, verbose=0)
print("Test loss:", score[0])
print("Test accuracy:", score[1])

Epoch 1/12
669/469 [=====] - 27s 58ms/step - loss: 2.2619 - accuracy: 0.2063 - val_loss: 2.2114 - val_accuracy: 0.3310
Epoch 2/12
669/469 [=====] - 28s 60ms/step - loss: 2.1491 - accuracy: 0.4555 - val_loss: 2.0391 - val_accuracy: 0.5864
Epoch 3/12
669/469 [=====] - 27s 58ms/step - loss: 1.9067 - accuracy: 0.6494 - val_loss: 2.0854 - val_accuracy: 0.7124
Epoch 4/12
669/469 [=====] - 27s 58ms/step - loss: 1.4030 - accuracy: 0.7373 - val_loss: 0.8907 - val_accuracy: 0.7819
Epoch 5/12
669/469 [=====] - 26s 56ms/step - loss: 0.8894 - accuracy: 0.7974 - val_loss: 0.6963 - val_accuracy: 0.8330
Epoch 6/12
669/469 [=====] - 26s 55ms/step - loss: 0.6335 - accuracy: 0.8367 - val_loss: 0.5349 - val_accuracy: 0.8593
Epoch 7/12
669/469 [=====] - 26s 55ms/step - loss: 0.5162 - accuracy: 0.8582 - val_loss: 0.4536 - val_accuracy: 0.8766
Epoch 8/12
669/469 [=====] - 25s 54ms/step - loss: 0.4502 - accuracy: 0.8747 - val_loss: 0.4037 - val_accuracy: 0.8861
Epoch 9/12
669/469 [=====] - 22s 47ms/step - loss: 0.4071 - accuracy: 0.8849 - val_loss: 0.3694 - val_accuracy: 0.8954
Epoch 10/12
669/469 [=====] - 24s 50ms/step - loss: 0.3761 - accuracy: 0.8918 - val_loss: 0.3431 - val_accuracy: 0.9030
Epoch 11/12
669/469 [=====] - 23s 49ms/step - loss: 0.3520 - accuracy: 0.8995 - val_loss: 0.2931 - val_accuracy: 0.9072
Epoch 12/12
669/469 [=====] - 24s 52ms/step - loss: 0.3325 - accuracy: 0.9050 - val_loss: 0.2854 - val_accuracy: 0.9072
Test loss: 0.3060(751804351807)
Test accuracy: 0.911700(10296826)
```

```
In [17]: # Helper function, You don't need to modify it
# Use all data in history
print(history.history.keys())
# summarize history for accuracy
plt.plot(history.history['accuracy'])
plt.plot(history.history['val_accuracy'])
plt.title('model accuracy')
plt.ylabel('accuracy')
plt.xlabel('epoch')
plt.legend(['train', 'test'], loc='upper left')
plt.show()

dict keys(['loss', 'accuracy', 'val_loss', 'val_accuracy'])
```



## 3: Random Forests [40pts] \*\*[P]\*\* \*\*[W]\*\* \*\*[G]\*\*

**NOTE:** Please use sklearn's DecisionTreeClassifier in your Random Forest implementation. You can find more details about this classifier [here](#).

### 3.1 Random Forest Implementation (30 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 is inevitably leads to overfitting. 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 Keras to solve the image classification task for MNIST. If you haven't installed TensorFlow, you can install the package by pip command or train your model by uploading HW4 notebook to Colab directly. Colab contains all packages you need for this section.

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

- For every tree in the random forest, we're going to
  - Subsample the examples with replacement. Note that in this question, the size of the subsample is equal to the original dataset.
  - From the subsampling 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 task easier. We randomly pick some features (70% percent of features) and grow the tree based on the pre-decided randomly selected features. Therefore, there is no need to find random features in each split.
  - 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,

- 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.
- 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.

```
In [18]: import numpy as np
import math
import sklearn
import time
import joblib

class RandomForest(object):
    def __init__(self, n_estimators=50, max_depth=None, max_features=0.7):
        # Helper function, You don't have to modify it
        # Initialization done here
        self.n_estimators = n_estimators
        self.max_depth = max_depth
        self.max_features = max_features
        self.bootstrap = True
        self.out_of_bag = []
        self.decision_trees = [sklearn.tree.DecisionTreeClassifier(max_depth=max_depth, criterion='entropy')]
        for i in range(n_estimators):
            self.bootstrap_fit(self, num_training, num_features, random_seed = None)

    def bootstrap(self, num_training, num_features, random_seed = None):
        TODO:
        - Randomly select indices of size num_training with replacement corresponding to row locations
        or
        - Randomly select indices without replacement corresponding the column locations of selected features in the original feature
        # This function computes the accuracy of the random forest model predicting y given x.
        # num_features: an integer N representing the total number of features in the training set, max_features denotes the percentage
        # of features that are used to fit each decision tree.

        Reference: https://en.wikipedia.org/wiki/Bootstrapping_(statistics)

        Args:
        - num_training: an integer N representing the total number of training instances.
        - num_features: an integer N representing the total number of features.

        Returns:
        - row_idx: (N,) numpy array of row indices corresponding to the row locations of the selected samples in the original dataset.
        - col_idx: (N,) array of column indices corresponding to the column locations of selected features in the original feature

        Hint: Consider using np.random.choice.

        np.random.seed(random_seed)
        row_idx = np.random.choice(num_training, num_training, replace=True)
        col_idx = np.random.choice(num_features, math.floor(num_features * self.max_features), replace=False)
        return row_idx, col_idx

    def bootstrap(self, num_training, num_features):
        Args:
        - num_training: an integer N representing the total number of training instances.
        - num_features: an integer N representing the total number of features.

        Returns:
        - None
        - self: decision tree in self.decision_trees
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        TODO:
        Train decision trees using the bootstrapped datasets.
        Note that you need to use the row indices and column indices.

        Args:
        -X: NxD numpy array, where N is number
        of instances and D is the dimensionality of each instance
        -y: Nx1 numpy array, the predicted labels

        Returns:
        - None
        - self: decision tree in self.decision_trees
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self.feature_indices.append(col_idx)
            self.out_of_bag.append(total)

    def fit(self, X, y):
        # Helper function, You don't have to modify it
        # Initializing the bootstrap datasets for each tree
        for i in range(self.n_estimators):
            total = set(list(range(num_training)))
            row_idx, col_idx = self.bootstrap(num_training, num_features)
            self.bootstrap_row_indices.append(row_idx)
            self
```



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

Consider a dataset with 2 points in 1-dimensional space:  $(x_1 = -3, y_1 = -1)$  and  $(x_2 = 2, y_2 = 1)$ . Here  $x$  are the point coordinates and  $y$  are the classes.

Consider mapping each point to 3-dimensional space using the feature vector  $\phi(x) = [1, 2x, x^2]$ . (This is equivalent to using a second order polynomial kernel.) The max margin classifier has the form

$$y_1(\phi(x_1)|\theta) \geq 1$$
$$y_2(\phi(x_2)|\theta + b) \geq 1$$

**Hint:**  $\phi(x_1)$  and  $\phi(x_2)$  are the support vectors. We have already given you the solution for the support vectors and you need to calculate back the parameters. Margin is equal to  $\frac{1}{\|\theta\|}$  and full margin is equal to  $\frac{2}{\|\theta\|}$

(1) Find a vector parallel to the optimal vector  $\theta$ . (4pts)

$$Z_1 = \phi(x_1) = [1, 2x_1, x_1^2] = [1, -6, 9] \quad \text{and} \quad Z_2 = \phi(x_2) = [1, 2x_2, x_2^2] = [1, 4, 4]$$
$$v = Z_2 - Z_1 \Rightarrow v = [1, 4, 4] - [1, -6, 9] \Rightarrow v = [0, 10, -5]$$
$$\hat{v} = \frac{v}{\|v\|} \Rightarrow \hat{v} = \frac{[0, 10, -5]}{\sqrt{(10)^2 + (-5)^2}} \Rightarrow \hat{v} = [0, \frac{2\sqrt{5}}{5}, -\frac{\sqrt{5}}{5}]$$

(2) Calculate the value of the margin achieved by this  $\theta$ ? (4pts)

$$margin = \frac{\|\hat{v}\|}{2} \Rightarrow margin = \frac{5\sqrt{5}}{2}$$

(3) Solve for  $\theta$ , given that the margin is equal to  $1/\|\theta\|$ . (4pts)

Because  $\theta$  and  $\hat{v}$  are parallel, we know that  $\theta = k * \hat{v}$ . Lets solve for  $\theta$

$$margin = \frac{1}{\|\theta\|} \Rightarrow \|\theta\| = \frac{1}{margin} \Rightarrow \|\theta\| = \frac{2\sqrt{5}}{25}$$
$$\theta = k * \hat{v} \Rightarrow \theta = [0, \frac{2\sqrt{5}}{5}k, -\frac{\sqrt{5}}{5}k]$$
$$\|\theta\| = \sqrt{(\frac{2\sqrt{5}}{5}k)^2 + (-\frac{\sqrt{5}}{5}k)^2} = \frac{2\sqrt{5}}{25} \Rightarrow (\frac{2\sqrt{5}}{5}k)^2 + (-\frac{\sqrt{5}}{5}k)^2 = \frac{4}{125} \Rightarrow k^2(\frac{4}{5} + \frac{1}{5}) = \frac{4}{125} \Rightarrow \sqrt{k^2} = \sqrt{\frac{4}{125}} \Rightarrow k = \pm \frac{2\sqrt{5}}{25}$$
$$\theta = [0, 0.16, -0.08] \quad \text{or} \quad \theta = [0, -0.16, 0.08]$$

(4) Solve for  $b$  using your value for  $\theta$ . (4pts)

To find  $b$ , we must plug in our two options for  $\theta$  into the two constraint inequalities. The correct one will yield a result where  $b \leq b'$  and  $b \geq b'$

Lets try  $\theta = [0, 0.16, -0.08]$  first:

$$y_1(Z_1\theta + b) \geq 1 \Rightarrow -1((0 - 6(0.16) - 9(0.08)) + b) \geq 1 \Rightarrow (-0.96 - 0.72) + b \leq -1 \Rightarrow -1.68 + b \leq -1 \Rightarrow b \leq 0.68$$
$$y_2(Z_2\theta + b) \geq 1 \Rightarrow 1((0 + 4(0.16) - 4(0.08)) + b) \geq 1 \Rightarrow (0.64 - 0.32) + b \geq 1 \Rightarrow 0.32 + b \geq 1 \Rightarrow b \geq 0.68$$

We can conclude that:  $b = 0.68$  and that the correct  $\theta$  in question 3 is:  $\theta = [0, 0.16, -0.08]$

(5) Write down the form of the discriminant function  $f(x) = \phi(x)\theta + b$  as an explicit function of  $x$ .

$$f(x) = \phi(x)\theta + b \Rightarrow f(x) = -0.08x^2 + 0.32x + 0.68$$

4.2 Feature Mapping (10 Pts) \*\*[W]\*\*

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 [33]: # DO NOT CHANGE
# Generate dataset

random_state = 1

X, y = make_moons(n_samples=1000, noise=.05)

y = np.where(y == 0, -1, y)

X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.20,
                                                    random_state=random_state)

f, ax = plt.subplots(nrows=1, ncols=1, figsize=(5,5))
plt.scatter(X[:, 0], X[:, 1], c = y, marker = '+')
plt.show()
```

```
In [34]: def visualize_decision_boundary(X, y, feature_new=None, h=0.02):
    """
    You don't have to modify this function

    Function to visualize decision boundary

    feature_new is a function to get X with additional features
    """
    x1_min, x1_max = X[:, 0].min() - 1, X[:, 0].max() + 1
    x2_min, x2_max = X[:, 1].min() - 1, X[:, 1].max() + 1
    xx_1, xx_2 = np.meshgrid(np.arange(x1_min, x1_max, h),
                              np.arange(x2_min, x2_max, h))

    if X.shape[1] == 2:
        Z = svm_clf.predict(np.c_[xx_1.ravel(), xx_2.ravel()])
    else:
        X_conc = np.c_[xx_1.ravel(), xx_2.ravel()]
        X_new = feature_new(X_conc)
        Z = svm_clf.predict(X_new)

    Z = Z.reshape(xx_1.shape)

    f, ax = plt.subplots(nrows=1, ncols=1, figsize=(5,5))
    plt.contourf(xx_1, xx_2, Z, cmap=plt.cm.coolwarm, alpha=0.8)
    plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.coolwarm)
    plt.xlabel('X_1')
    plt.ylabel('X_2')
    plt.xlim(xx_1.min(), xx_1.max())
    plt.ylim(xx_2.min(), xx_2.max())
    plt.xticks(())
    plt.yticks(())

    plt.show()
```

```
In [35]: # DO NOT CHANGE
# Try to fit a linear classifier to the dataset

svm_clf = svm.LinearSVC()
svm_clf.fit(X_train, y_train)
y_test_predicted = svm_clf.predict(X_test)

print("Accuracy on test dataset: {}".format(accuracy_score(y_test,
                                                            y_test_predicted)))

visualize_decision_boundary(X_train, y_train)
```

Accuracy on test dataset: 0.865



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.

**(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](#) for a detailed analysis of doing the same for points separable with a circular boundary)

```
In [36]: def create_n1_feature(X):
    """
    TODO - Create additional features and add it to the dataset

    returns:
        X_new - (N, d + num_new_features) array with
        additional features added to X such that it
        can classify the points in the dataset.
    """
    # Delete this line when you implement the function
    # (x1, x2, exp(-(x1)**2+(x2)**2))
    X_new = np.asarray([X[:,0], X[:,1], np.exp(-(X[:,0]**2 + X[:,1]**2))]).T
    return X_new
```

```
In [37]: # DO NOT CHANGE
# Create new features

X_new = create_n1_feature(X)
X_train, X_test, y_train, y_test = train_test_split(X_new, y,
                                                    test_size=0.20,
                                                    random_state=random_state)
```

```
In [38]: # DO NOT CHANGE
# Fit to the new features and visualize the decision boundary
# You should get more than 90% accuracy on test set

svm_clf = svm.LinearSVC()
svm_clf.fit(X_train, y_train)
y_test_predicted = svm_clf.predict(X_test)

print("Accuracy on test dataset: {}".format(accuracy_score(y_test, y_test_predicted)))

visualize_decision_boundary(X_train, y_train, create_n1_feature)
```



In [ ] :