

---

## CS5691: PRML Assignment #2

**Instructor :** Prof. B. Ravindran

**Topics:** ANN, Ensemble Method, Kernel and SVM.

**Deadline:** 12th November 2021

**Teammate 1:** (your name here)

**Roll number:** CS21YZZZ

**Teammate 2:** (your name here)

**Roll number:** CS21YZZZ

---

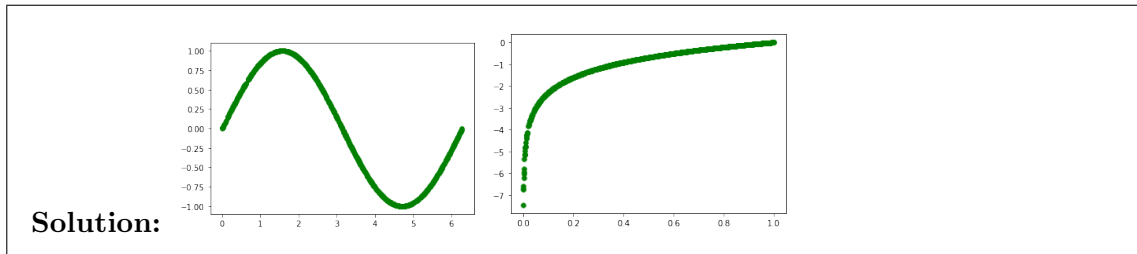
- This assignment has to be completed in teams of 2. Collaborations outside the team are strictly prohibited.
  - Be precise with your explanations. Unnecessary verbosity will be penalized.
  - Check the Moodle discussion forums regularly for updates regarding the assignment.
  - Type your solutions in the provided  $\text{\LaTeX}$  template file.
  - We highly recommend using **Python 3.6+** and standard libraries like **numpy**, **Matplotlib**, **pandas**. You can choose to use your favourite programming language however the TAs will only be able to assist you with doubts related to Python.
  - You are supposed to write your own algorithms, any library functions which implement these directly are strictly off the table. Using them will result in a straight zero on coding questions, **import wisely!**
  - **Please start early and clear all doubts ASAP.**
  - Please note that the TAs will **only** clarify doubts regarding problem statements. The TAs won't discuss any prospective solution or verify your solution or give hints.
  - Post your doubt only on Moodle so everyone is on the same page.
  - Posting doubts on Moodle that reveals the answer or gives hints may lead to penalty
  - **Only one team member will submit the solution**
  - For coding questions paste the link to your Colab Notebook of your code in the  $\text{\LaTeX}$  solutions file as well as embed the result figures in your  $\text{\LaTeX}$  solutions. Make sure no one other than TAs have access to the notebook. And do not delete the outputs from notebook before final submission . Any update made to notebook after deadline will result in standard late submission penalty.
  - Late submission per day will attract a penalty of 10 percent of the total marks.
- 

1. **[ANN]** In this Question, you will code a single layer ANN with Sigmoid Activation function and appropriate loss function from scratch. Train the ANN for the [Dataset1](#) and [Dataset2](#)

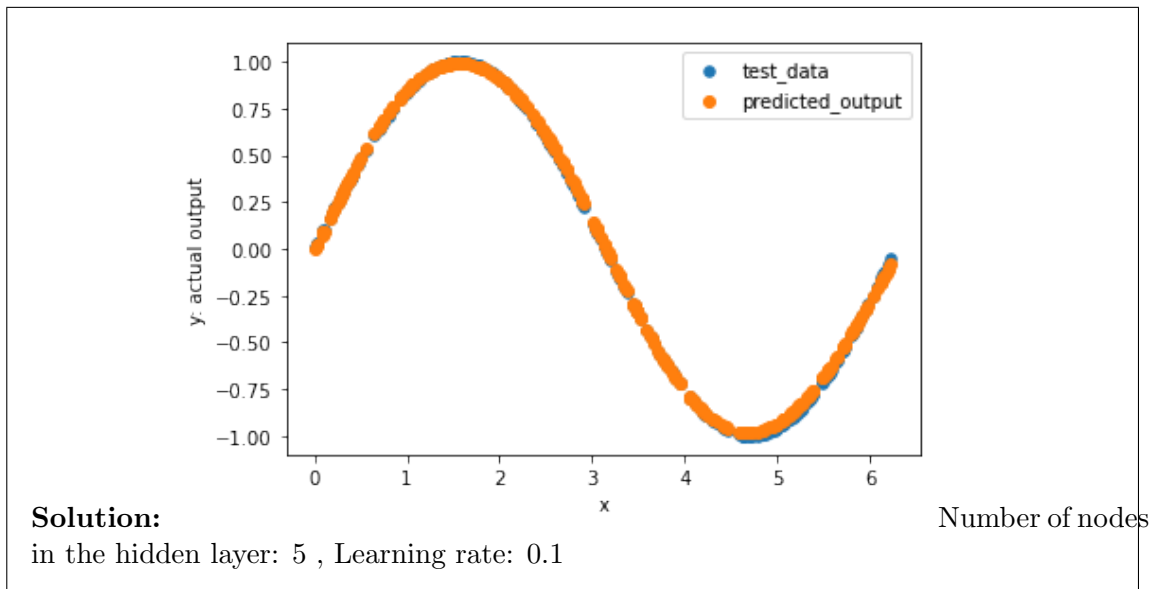
**NOTE: Test Data should not be used for training.**

**NOTE 2: You need to code from scratch.**

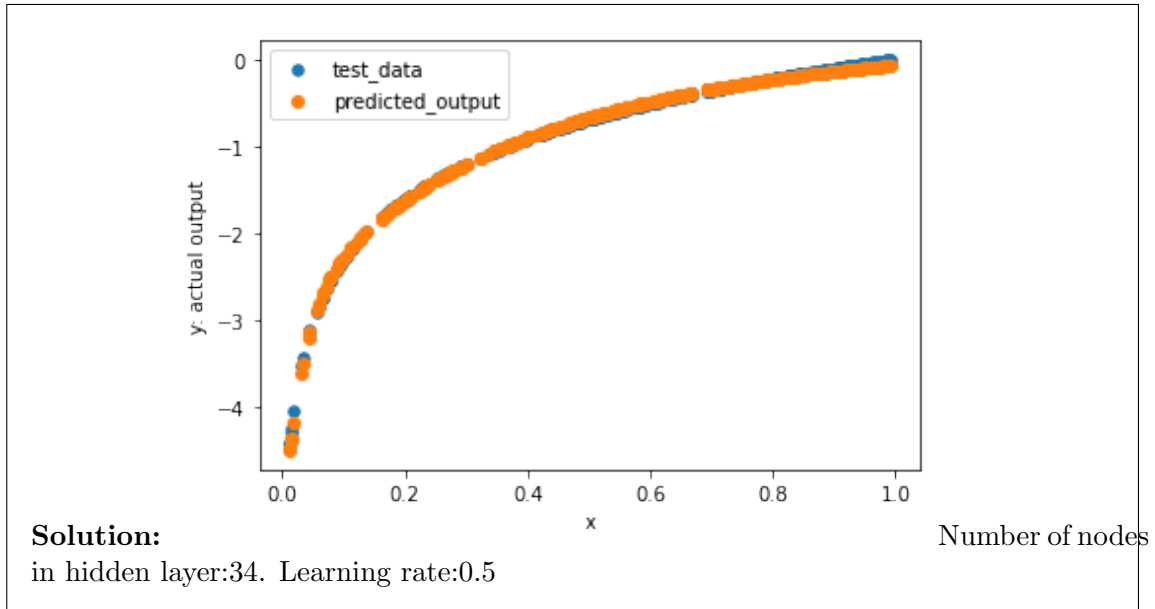
- (a) (1 mark) Plot the training Data for Dataset1 and Dataset2.



- (b) (1 mark) **For data set 1 :** (1) Write the number of nodes in the hidden layer and learning rate used (2) Plot Test Data and prediction on Test Data in the same graph with different colors and appropriate legend.



- (c) (1 mark) **For data set 2 :** (1) Write the number of nodes in the hidden layer and learning rate used (2) Plot Test Data and prediction on Test Data in the same graph with different colors and appropriate legend.



- (d) (1 mark) For each Dataset write average training Loss and average Test Loss.

**Solution:** Dataset1: Average training loss = 0.00017 Average Test Loss = 0.00016  
Dataset2: Average training loss = 0.0156 Average Test Loss = 0.00077

- (e) (1 mark) What Loss function did you use and why?

**Solution:** Have used Mean Squared Error loss function. Because the output of the datasets is continuous values hence it falls under regression problem.

- (f) (3 marks) Paste the link to your Colab Notebook of your code. Make sure that your notebook is private and give access to all the TAs. Your Notebook must contain all the codes that you used to generate the above results. **Note :** Do not delete the outputs.

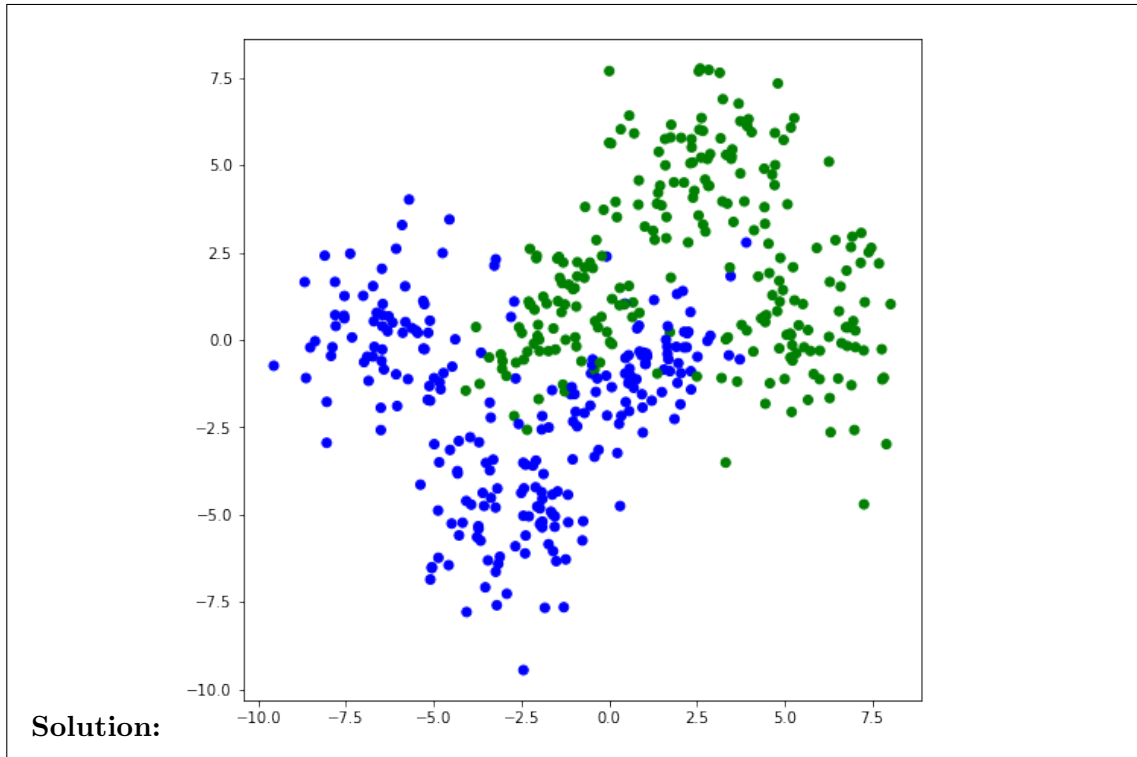
**Solution:** [https://colab.research.google.com/drive/15\\_TEvCLTI4DGLHq81GgwL3UijdaQ5jrd?usp=sharing](https://colab.research.google.com/drive/15_TEvCLTI4DGLHq81GgwL3UijdaQ5jrd?usp=sharing)

2. [AdaBoost] In this question, you will code the AdaBoost algorithm. Follow the instructions in this [Jupyter Notebook](#) for this question. Find the dataset for the question [here](#).

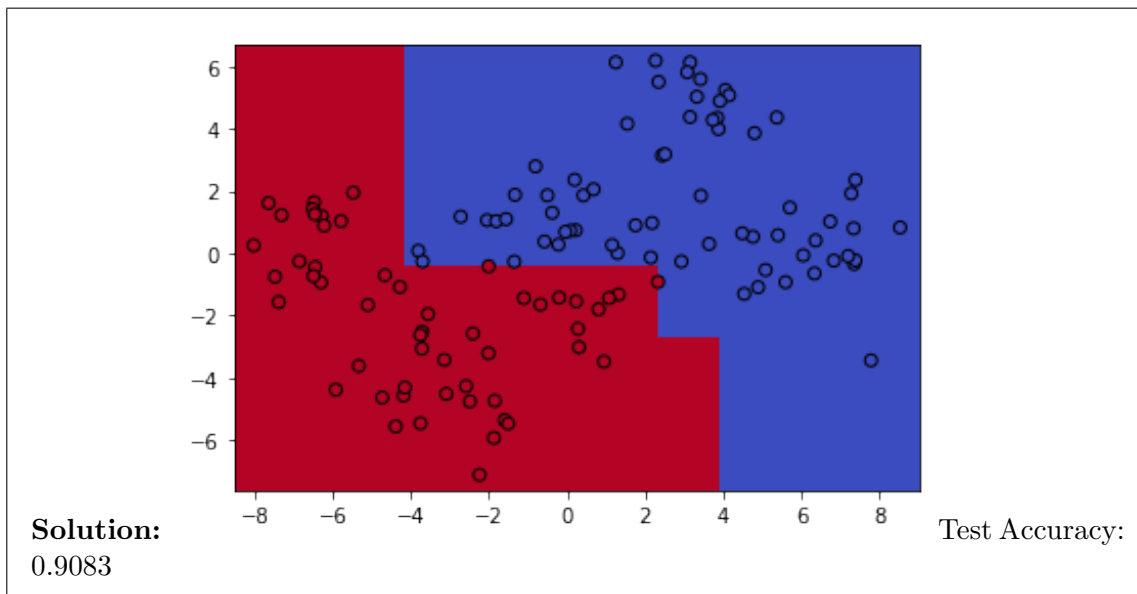
**NOTE:** Test data should not be used for training.

**NOTE 2:** You need to code from scratch. You can use the starter notebook though :)  
. Make a copy of it in your drive and start.

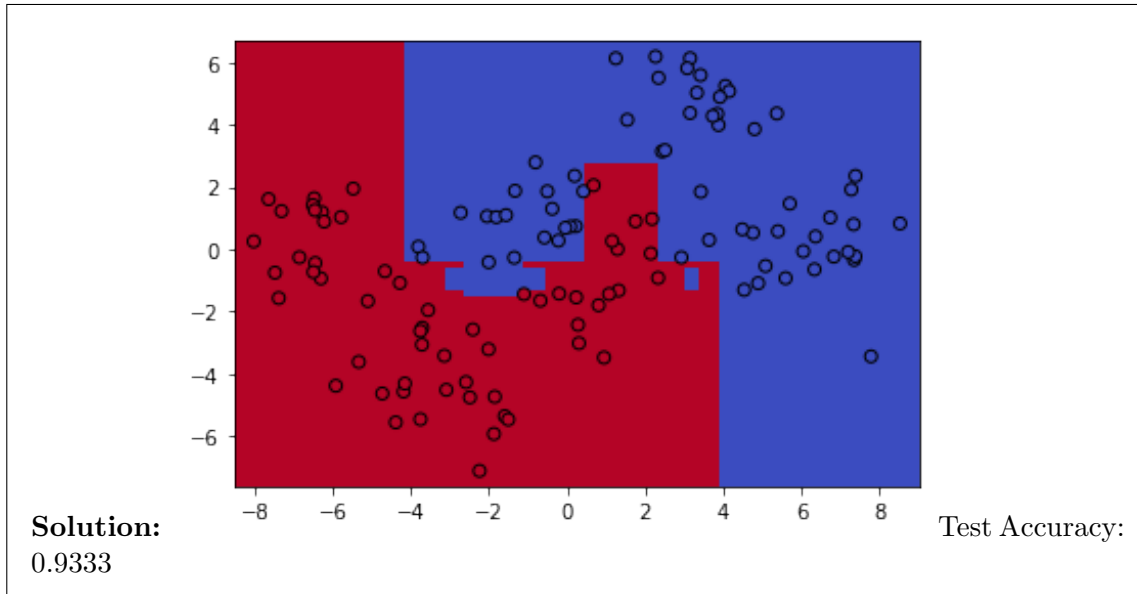
- (a) (1 mark) Plot the training data.



- (b) (1 mark) **For training with  $k=5$**  : (1) Plot the learnt decision surface. (2) Write down the test accuracy.



- (c) (1 mark) **For training with  $k=100$**  : (1) Plot the learnt decision surface. (2) Write down the test accuracy.

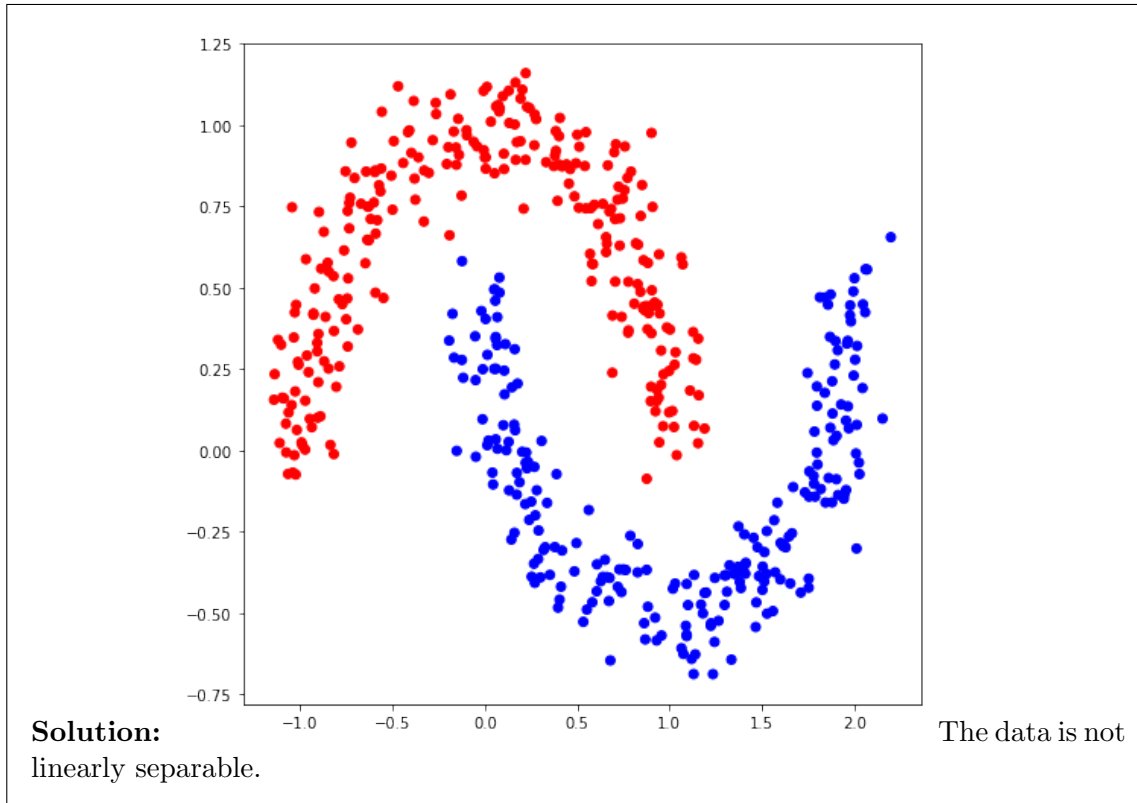


- (d) (3 marks) Paste the link to your Colab Notebook of your code. Make sure that your notebook is private and give access to all the TAs. Your notebook must contain all the codes that you used to generate the above results. **Note :** Do not delete the outputs.

**Solution:** [https://colab.research.google.com/drive/1lI7dswlvug9f\\_1RJn22Fdd0lHlcFAe\\_?usp=sharing](https://colab.research.google.com/drive/1lI7dswlvug9f_1RJn22Fdd0lHlcFAe_?usp=sharing)

3. [Kernel] Consider *Dataset\_Kernel\_Train.npy* and *Dataset\_Kernel\_Test.npy* for this question. Each row in the above matrices corresponds to a labelled data point where the first two entries correspond to its  $x$  and  $y$  co-ordinate, and the third entry  $\in \{-1, 1\}$  indicates the class to which it belongs. Find the dataset for the question [here](#). **NOTE: Test data should not be used for training.**

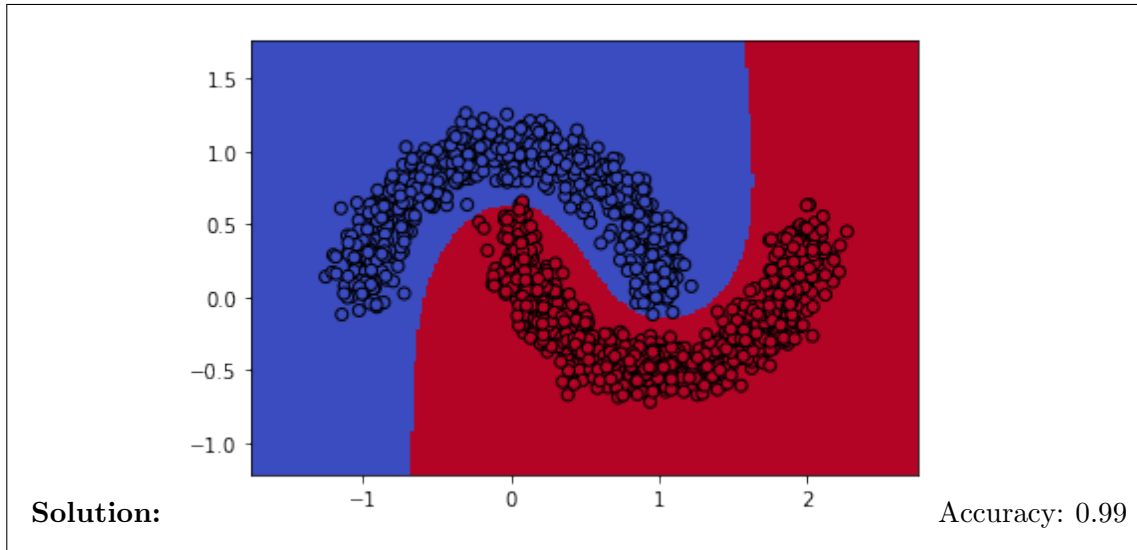
- (a) (1 mark) Plot the training data points and indicate by different colours the points belonging to the different classes. Is the data linearly separable?



- (b) (1 mark) Using *sklearn.svm* (read the documentation [here](#)), build a classifier that classifies the data points in the testing data set using the Radial Basis Function (RBF) kernel. How do you tune the involved hyperparameters?

**Solution:** Have tuned the parameter  $C$  and  $C = 1$  gave the optimal solution. The strength of regularization is inversely proportional to  $C$ , hence have tuned it. Have used different values of  $C$  in the range 1 to 10 to declare the SVM model. and then fit the SVM classifier with training data. Have evaluated the prediction of classifier with training data itself, the value of  $C$  for which I have got maximum accuracy, that value is taken.

- (c) (1 mark) Plot the separating curve and report the accuracy of prediction corresponding to the tuned hyperparameters.



- (d) (3 marks) Paste the link to your Colab Notebook of your code. Make sure that your notebook is private and give access to all the TAs. Your notebook must contain all the codes that you used to generate the above results. **Note :** Do not delete the outputs.

**Solution:** [https://colab.research.google.com/drive/1KTSa74\\_-DaYG88-GPUvpUjs1hN1pGgVI?usp=sharing](https://colab.research.google.com/drive/1KTSa74_-DaYG88-GPUvpUjs1hN1pGgVI?usp=sharing)

4. (2 marks) **[Ensemble of randomised algorithms]** Imagine we have an algorithm for solving some decision problem (*e.g.*, is a given number  $p$  a prime?). Suppose that the algorithm makes a decision at random and returns the correct answer with probability  $\frac{1}{2} + \delta$ , for some  $\delta > 0$ , which is just a bit better than a random guess. To improve the performance, we run the algorithm  $N$  times and take the majority vote. Show that for any  $\epsilon \in (0, 1)$ , the answer is correct with probability  $1 - \epsilon$ , as long as  $N > (1/2)\delta^{-2} \ln(\epsilon^{-1})$ .

*Hint 1: Try to calculate the probability with which the answer is not correct i.e. when the majority votes are not correct.*

*Hint 2: What value of  $N$  will you require so that the above probability is less than  $\epsilon$ . Rearrange Inequalities :-)*

**Solution:**

5. (2 marks) **[Boosting]** Consider an additive ensemble model of the form  $f_m(\mathbf{x}) = \frac{1}{2} \sum_{l=1}^m \alpha_l y_l(\mathbf{x})$ , where  $y_l$ 's are individual models and  $\alpha_l$ 's are weights. Show that the sequential minimization of the sum-of-squares error function for the above model trained in the style of boosting (i.e.  $y_m$  is trained after accounting the weaknesses of  $f_{m-1}$ ) simply involves fitting each new base classifier  $y_m$  to the residual errors  $t_n - f_{m-1}(\mathbf{x}_n)$  from previous model.

**Solution:**

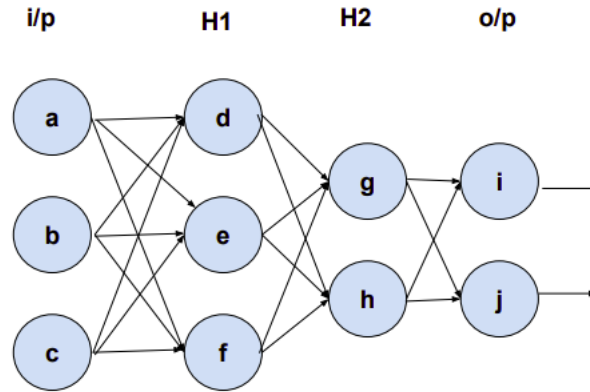
6. (2 marks) **[Backpropagation]** We are trying to train the following chain like neural network with back-propagation. Assume that the transfer functions are sigmoid activation functions i.e.  $g(x) = \frac{1}{1+e^{-x}}$ . Let the input  $x = 0.5$ , the target output  $y = 1$ , all the weights are initially set to 1 and the bias for each node is -0.5.



- (a) Give an expression that compares the magnitudes of the gradient updates for weights ( $\delta$ ) across the consecutive nodes.  
 (b) How does the magnitude of the gradient update vary across the network/chain as we move away from the output unit?

**Solution:**

7. (2 marks) **[NN & Activation Functions]** The following diagram represents a feed-forward network with two hidden layers.



A weight on connection between nodes  $x$  and  $y$  is denoted by  $w_{xy}$ , such as  $w_{ad}$  is the weight on the connection between nodes a and d. The following table lists all the weights in the network:

$w_{ad} = 0.5$	$w_{be} = -1.4$	$w_{cf} = -1.25$	$w_{eh} = -2$	$w_{gj} = -1.5$
$w_{ae} = 0.9$	$w_{bf} = 0.75$	$w_{dg} = 1$	$w_{fg} = 3$	$w_{hi} = 0.5$
$w_{af} = -2$	$w_{cd} = 0$	$w_{dh} = 3$	$w_{fh} = 1.25$	$w_{hj} = -0.25$
$w_{bd} = 1.3$	$w_{ce} = 0.3$	$w_{eg} = 2.5$	$w_{gi} = 2.5$	



Find the output of the network for the following input vectors:

$V_1 = [0.2, 1, 3]$ ,  $V_2 = [2.5, 3, 7]$ ,  $V_3 = [0.75, -2, 3]$

- (a) If *sigmoid* activation function is used in both H1 & H2
- (b) If *tanh* activation function is used in both H1 & H2
- (c) If *sigmoid* activation is used in H1 and *tanh* in H2
- (d) If *tanh* activation is used in H1 & ReLU activation in H2

Please provide all steps and explain the same.

**Solution:**

8. (2 marks) [**Decision Trees**] Consider a dataset with each data point  $x \in \{0, 1\}^m$ , i.e.,  $x$  is a binary valued feature vector with  $m$  features, and the class label  $y \in \{+1, -1\}$ . Suppose the true classifier is a majority vote over the features, such that

$$y = \text{sign}\left(\sum_{i=1}^m (2x_i - 1)\right)$$

where  $x_i$  is the  $i^{\text{th}}$  component of the feature vector. Suppose you build a binary decision tree with minimum depth, that is consistent with the data described above. What is the range of number of leaves which such a decision tree will have?

**Solution:**