

P1) a)

$$\sum_n \log P(x_n, z_n | \theta) = \sum_k \sum_n \gamma_{nk} \log w_k + \sum_k \left\{ \sum_n \gamma_{nk} \log N(x_n | \mu_k, c_k) \right\}$$

$$= \log w_1 + 2 \log w_2 + \log N(x_1 | \mu_{k=1}, c_{k=1} | \pi_{k=1})$$

$$+ \log N(x_2 | \mu_{k=2}, \pi_{k=2} | c_{k=2})$$

$$+ \log N(x_3 | \mu_{k=2}, \pi_{k=2} | c_{k=2})$$

b)  $\pi_{k=1} = \frac{1}{3}$  ;  $\pi_{k=2} = \frac{2}{3}$

$\mu_{k=1} = \left(\frac{1}{3}\right)(1)$  ;  $\mu_{k=2} = \frac{1}{2}(10+20) = 15$   $5^2$

$c_{k=1} = \left(\frac{1}{3}\right)(1-1) = 0$  ,  $c_{k=2} = \frac{1}{2} \left[ (10-15)^2 + (20-15)^2 \right]$

$= \left(\frac{1}{2}\right)(25) = 12.5$

(1.2)  $w_k = \frac{\sum_n \gamma_{nk}}{\sum_k \sum_n \gamma_{nk}}$  ,  $\mu_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} x_n$

$\pi_{k=1} = \frac{1.4}{3}$  /  $\pi_{k=2} = \frac{1.6}{3}$

$\mu_{k=1} = \frac{1}{(1.4)} \left[ (1)(1) + (0.4)(10) \right] = \left(\frac{1}{1.4}\right)(5)$   $6.120$

$\mu_{k=2} = \frac{1}{(1.6)} \left[ (0.6)(10) + (1)(20) \right] = \left(\frac{1}{1.6}\right)(26)$

Problem 2)

a)

$$P_1(B) = P(Q_t = B | q_t = 1) = 0.01$$

$$P_2(A) = P(Q_t = A | q_t = 2) = 0.41$$

b) The most frequent symbol is one with largest  $P(X|\theta)$

$$P(X=A|\theta) = (0.49) \cdot 0.99 + 0.51 \cdot 0.49$$

$$P(X=B|\theta) = (0.49)(0.01) + (0.51)(0.51)$$

$$P(X=A|\theta) > P(X=B|\theta)$$

c) It is clear by calculating the probabilities the symbol with highest probability

AAA

### P3)a)

The data is split in two part using `sklearn.train_test_split ()`

---

### P3)b)

Using the `cross_validation` function, using `sklearn`, the error found was **0.14495188989709984**.

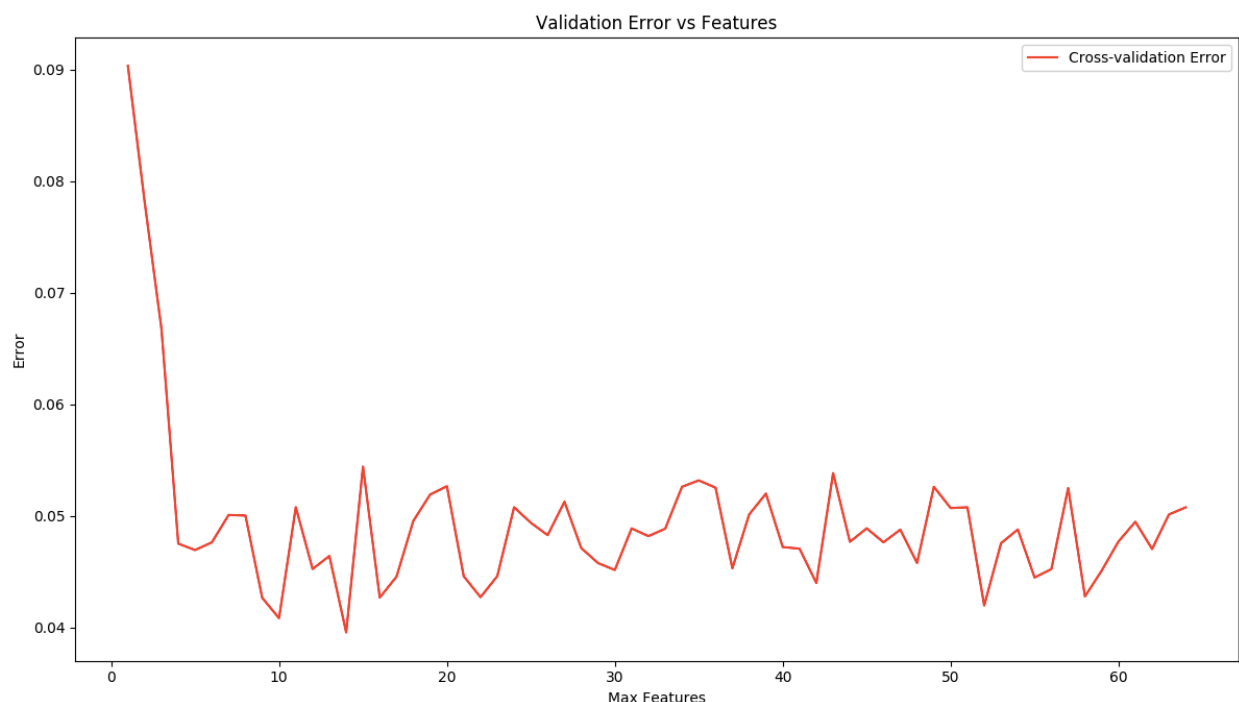
---

### P3)c)

Using the bagging classifier from `sklean`, the CV score (from `cross_val_score`) is found to be **0.050778474279623564**. We notice that using the bagging that uses sampling from the original data set to construct a classifier works much better than simply just using a weak learner such as a decision tree as the classifier.

---

### P3)d)



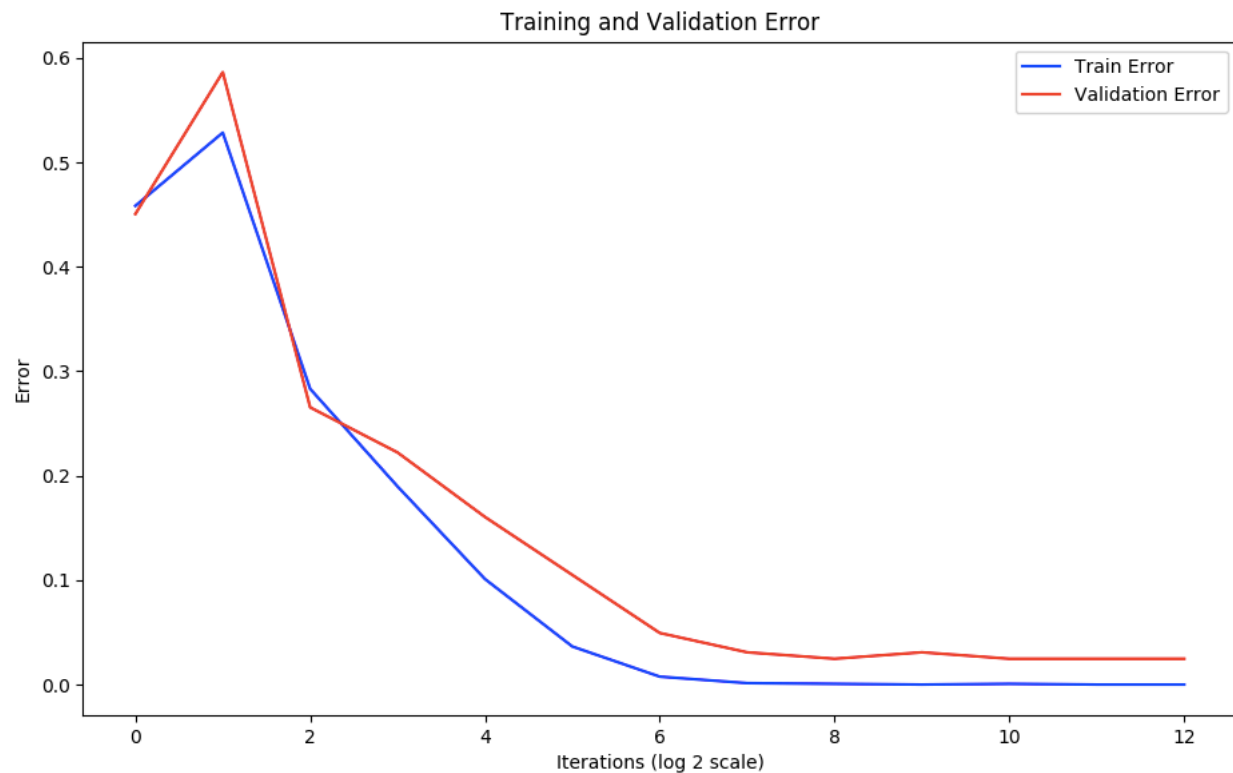
Random Forest classifier is a subclass of bagging classifier where we set the `max_features` for each node. This means that we set the `max_features` to be the number of samples that we **randomly** consider on each nodes and then split on one of those nodes (the best one), in this way, we introduce randomness into the decision tree classifier and reduce the variance.

By looping over the hyperparameter `max_features` we found that best value for the number of features to sample from randomly is **14**, and the cross validation error for random forest is **0.03957504383281463** better than the error we saw with just the bagging algorithm.

---

### P3)e)

The graph for error vs number of iterations is shown below. As shown in the plot, the minimum value of error for the `AdaBoostAlgorithm` is found to be **0.024691358024691357** and the optimal number of iterations is **256**.



---

**P3)f)**

Using the best hyperparameters we found for num\_features for the random forest algorithm and best number of estimators we found for the AdaBoost algorithm, we obtained the following results:

('Bagging test error', 0.06666666666666665)

('Random forest test error', 0.05555555555555558)

('Boosted tree test error', 0.038888888888888886)

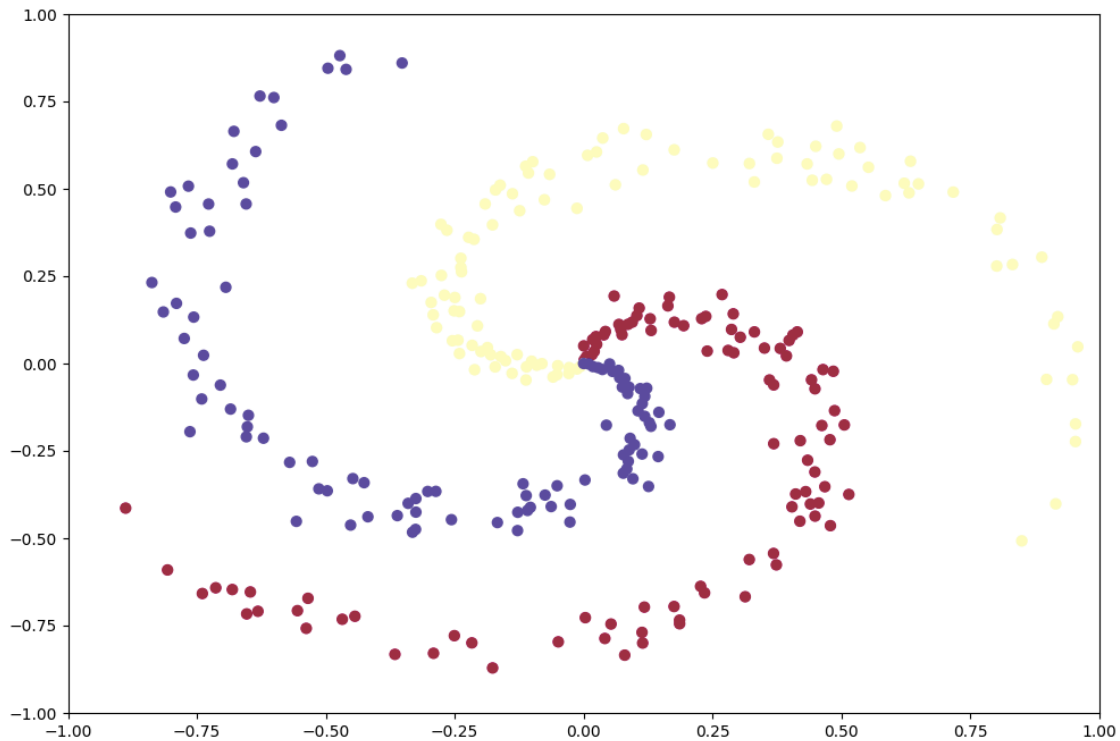
**Clearly boosted algorithm (with a tuned number of estimators) outperforms the other 2 classification algorithm.**

---

**P4)a)**



The visualization of the data is shown below. As seen from plots, the dimensionality of the data is 2 (2\_D), number of classes is 3, and data is not linearly separable which means in order to obtain better results, we can use a non-linear classifier such a neural network.




---

**P4)b)**

Let  $s = \begin{bmatrix} s_0 \\ s_1 \end{bmatrix}$ , now find  $y = \sigma(s) = \frac{e^{s_0}}{e^{s_0} + e^{s_1}} = \frac{1}{1 + e^{s_1 - s_0}} = \frac{1}{1 + e^{-s}}$  where we set  $s = s_1 - s_0$ . As seen softmax for 2D case is equivalent to sigmoid function.

---

**P4)c)**

We know

$$L_{CE} = - \sum_{k=1}^K t_k \log y_k$$

and now we take the partial derivative:

$$L = -[t_k \log(y_k)]$$

as there would be only one term where the indicator function is non-zero.

$$\frac{\partial L}{\partial s_k} = -t_k \frac{1}{y_k} \frac{\partial}{\partial s_k} y_k = -t_k \left( \frac{1}{y_k} \right) (y_k)(1 - y_k)$$

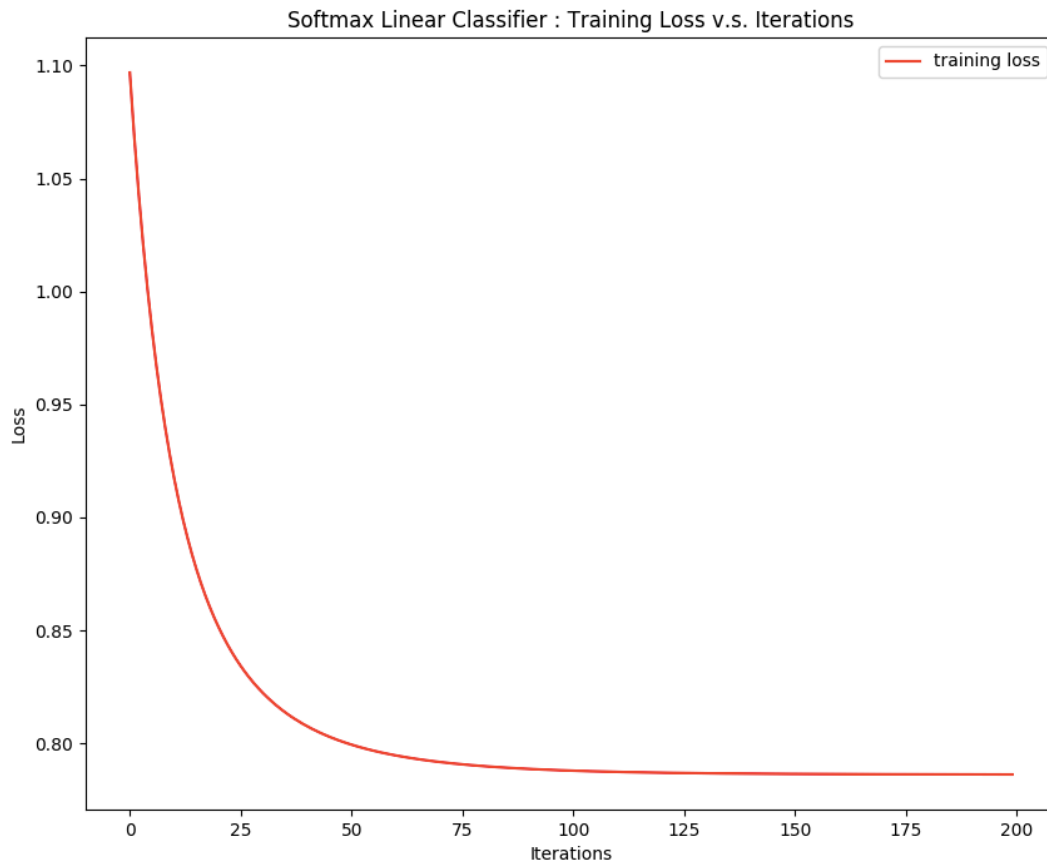
and we eventually get

$$\frac{\partial L}{\partial s_k} = y_k - t_k$$

as  $y_k * t_k = 1$ . Therefore, we have found the partial derivative and gradient of cross-entropy loss for a neural network.

---

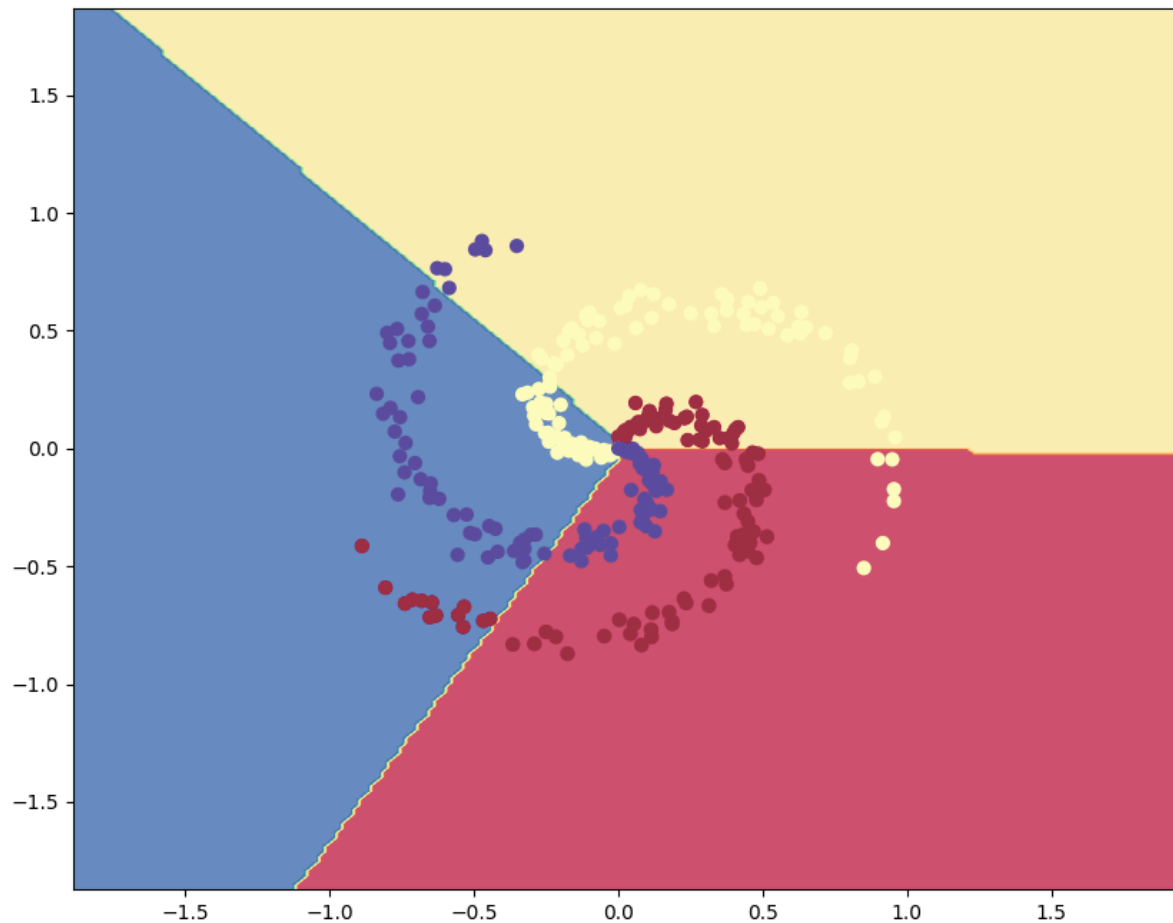
The learning curve (training error vs. number of iterations) is shown below. It is evident that as the number of iterations increase, the training error decreases. This also makes sense intuitively as with more iterations, we train the model better by using more iterations of the gradient descent to decrease the loss function. Also, the decision boundary of the Neural network is shown below. We can see that a shallow neural network produced a non-linear decision boundary. We can increase the accuracy by increasing the depth of the neural network. From the result we can see that the error for the softmax model is **0.49**.



$$\frac{\partial L}{\partial s_k} = -t_k \frac{1}{y_k} \frac{\partial}{\partial s_k} y_k = -t_k \left( \frac{1}{y_k} \right) (y_k)(1 - y_k)$$

$$\frac{\partial L}{\partial s_k} = -t_k \frac{1}{y_k} \frac{\partial}{\partial s_k} y_k = -t_k \left( \frac{1}{y_k} \right) (y_k)(1 - y_k)$$

==dc



---

#### P4)e)

The learning curve of the deeper neural network (with 2 layers) is shown below. Again, we can see that the training error reduces with the number of iteration. There are some spikes in the graph, and that can be due to a wrong direction of movement in the gradient descent but in general the learning curve is downward. Also, we can see

Accuracy is **0.98** for the case of neural network with 2 layers, which is less than 0.49 for the softmax case. We can see that the neural network classifies the points almost perfectly and with a very low error so neural network outperforms the softmax classifier and the deeper the Neural network, we would expect a higher performance as long as we can avoid overfitting.

#### P4)f)

The training accuracy for the 2 layer neural network is 0.98 (0.02 error) and the training accuracy for the linear model used in part is 0.49 (0.51 error). Evidently, we can see that the neural network performs a much better job at classifying a non-linear distribution. We see that our shallow neural network with only 2 layers perform a much better job than a linear model,

and we can improve performance even on more complex data sets by increasing the number of layers in the neural network.

