

# Kernel Methods for Pattern Analysis

## Assignment Solution 2

Rachit Garg  
CS14B050

Keerthana S  
CS13B041

Sphoorti K  
CS13B042

January 25, 2019

### 1 Linearly separable data set

#### 1.1 SVM

##### Linear kernel

**Parameter estimation:** Kernel gram matrix is used to choose kernel parameters. Since linear kernel does not have any kernel parameters, the figure looks like:

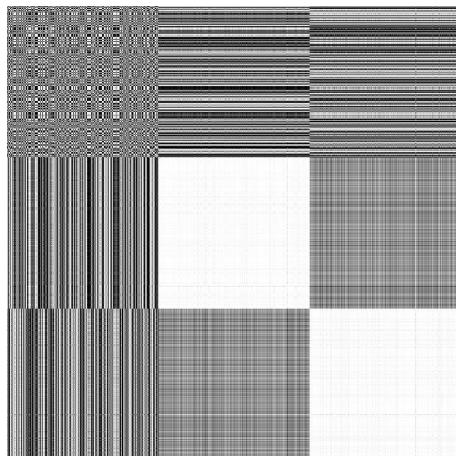


Figure 1: Grey level image of kernel gram matrix for dataset 1a training set using linear kernel

Values in kernel gram matrices are values of normalised kernel function for all pair of training examples which is measure of similarity of points. In the above figure we see that the region near diagonal has higher values since we expect examples of same class to be similar, on the other hand, region away from diagonal is grey/black since they belong to different class.

The parameters to be estimated in linear kernel is "c" which is tradeoff parameter. It is done by doing a search for "c" (in the range of  $2^{-10}$  to  $2^{10}$ ) which leads to maximum classification accuracy. The plot for accuracy on validation set vs  $\log(C)$  looks like:

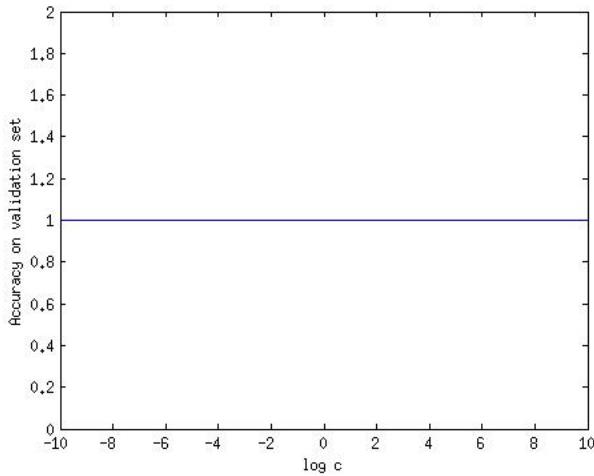


Figure 2: Figure showing accuracy for validation set for different values of  $c$

We see that accuracy=1 for all values of  $c$ .The reason for high accuracy is because the dataset is linearly separable and linear kernel is enough to fit a hyperplane.

#### Decision region plot:

The decision region plot for  $c=1$  looks like:

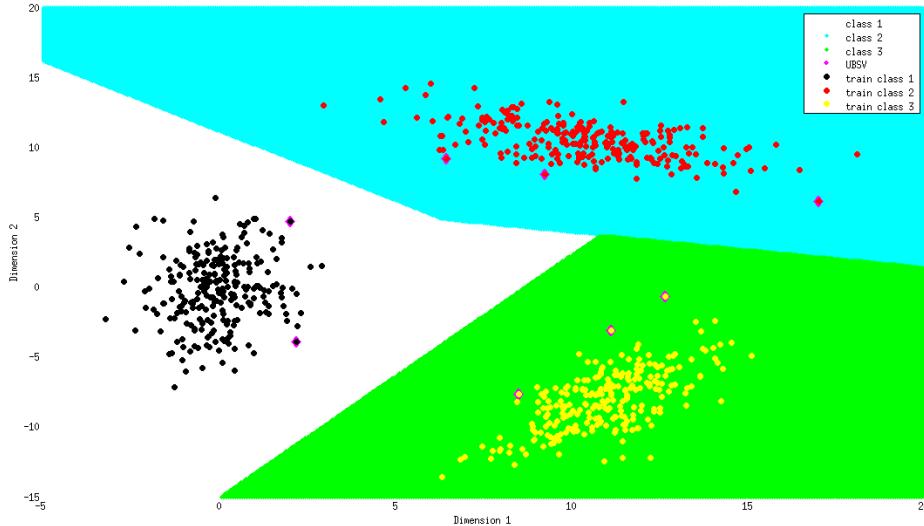


Figure 3: Figure showing decision region plot for dataset 1a using linear kernel, $c=1$

From the figure we see that, the decision boundaries separating each pair of classes are lines(linear),and the support vectors obtained by the training model are only unbounded i.e they lie on the margin because the points from each pair of classes are well separated.Decreasing  $c$  value may lead to more support vectors(especially bounded),because it gives more weight for margin term(maximising),which may end up having more points within margin.

#### Confusion matrix:

The confusion matrix for test set look like:

		Confusion Matrix		
		1	2	3
Output Class	1	100 33.3%	0 0.0%	0 0.0%
	2	0 0.0%	100 33.3%	0 0.0%
3	0 0.0%	0 0.0%	100 33.3%	100% 0.0%
	100% 0.0%	100% 0.0%	100% 0.0%	100% <b>0.0%</b>

Figure 4: Figure showing confusion matrix for test set

We see that  $accuracy = 100\%$  on test set.

## Polynomial kernel

### Parameter estimation:

The kernel parameters for polynomial function are p,a and b which are chosen using kernel gram matrix. We find that any  $p \geq 1$ , gave kernel gram matrix as expected(block diagonal matrix), hence we choose least "p" i.e  $p=1$  and we see that  $b=0$  is enough(2 dimensions in kernel space) which is enough to fit bivariate data. Hence the value of "a" chosen doesn't have any effect on normalised kernel function values, hence a is chosen to be 1.

The kernel gram matrix for above configuration looks like:

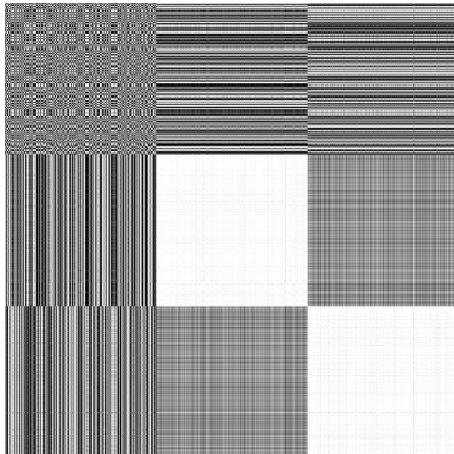


Figure 5: Figure showing grey level image of kernel gram matrix for kernel parameters  $a=1, b=0, p=1$

We see that it is the same one as obtained in linear kernel.  
For choosing c, we see accuracy on validation set which is found to be 1 for any value of "c". Hence

the below plots are plotted for  $c=1$ .

Since the model used is same the decision region plots and confusion matrix will be same  
**Decision region plot**

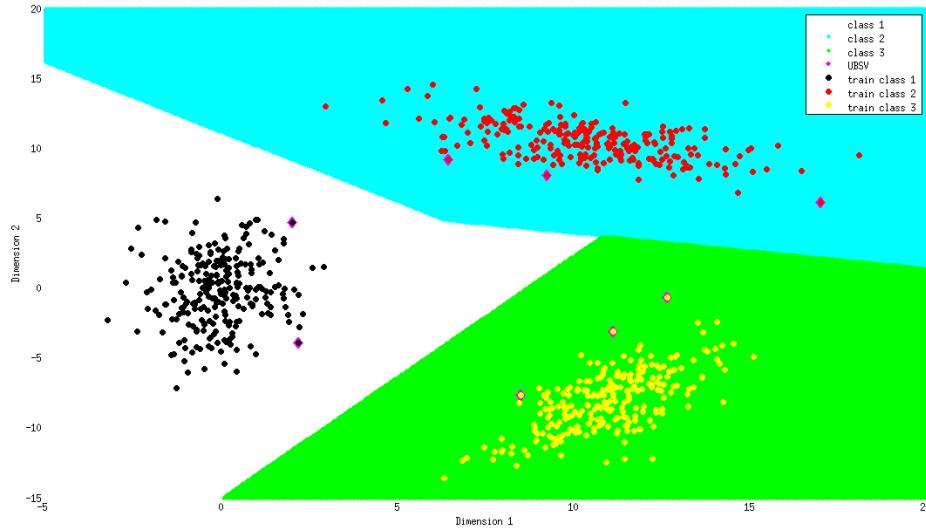


Figure 6: Figure showing decision region plot for dataset 1a polynomial kernel,  $c=1$

### Confusion matrix

The confusion matrix for dataset 1a using polynomial kernel and parameters  $a=1, b=0, p=1, c=1$  looks like:

		Confusion Matrix		
		1	2	3
Output Class	1	100 33.3%	0 0.0%	0 0.0%
	2	0 0.0%	100 33.3%	0 0.0%
3	0 0.0%	0 0.0%	100 33.3%	100% 0.0%
	1	100% 0.0%	100% 0.0%	100% 0.0%
	2	100% 0.0%	100% 0.0%	100% 0.0%
	3	100% 0.0%	100% 0.0%	100% 0.0%

Figure 7: Figure showing confusion matrix for test set for dataset 1a polynomial kernel,  $c=1$

We see that  $accuracy = 100\%$  for test set using polynomial kernel.

## Gaussian kernel

### Parameter estimation:

The kernel parameter to be estimated for gaussian kernel is  $\sigma$  where kernel function in gaussian kernel is given by:

$$K(x_m, x_n) = e^{-\|x_m - x_n\|^2/\sigma}$$

where  $\sigma$  is width parameter.

The  $\sigma$  value is chosen based on kernel gram matrix and the value obtained is  $\sigma=100$ .

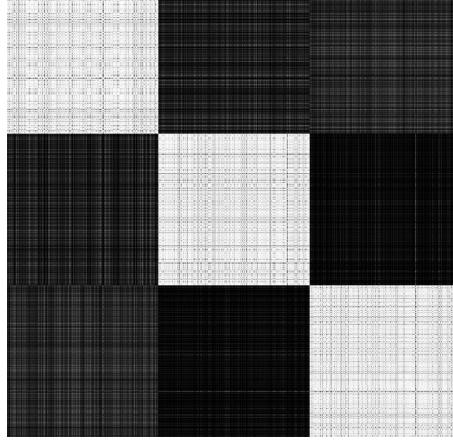


Figure 8: Grey level image of kernel gram matrix for dataset 1a training set using gaussian kernel for  $\sigma = 100$

From the figure of kernel gram matrix ,we see that kernel function values are much more high for examples of same class(especially class 1) and less when examples between classes are considered.This is indicated by more darker /black areas off the block diagonal and white area in block of first class.This is because we see that class1 examples are distributed in a circular fashion where cosine similarity has wider cone hence cosine similarity values vary whereas in gaussian kernel it takes euclidean distance into account.Hence the shape of data can also influence the choice of kernel to be used.

To find the hyperparameter.error on validation set is plotted across logc.

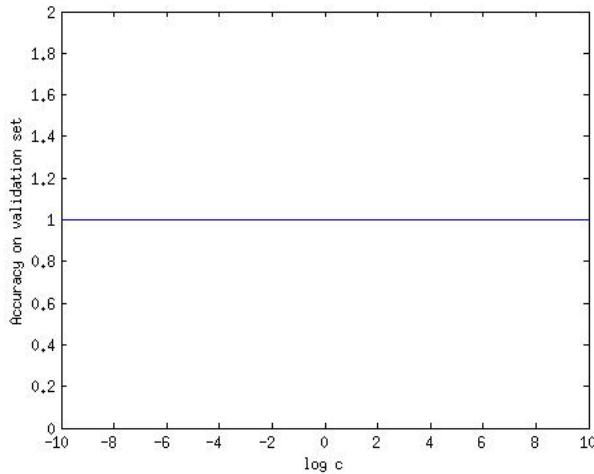


Figure 9: Figure showing accuracy for validation set for different values of c

We observe that accuracy is 100% for all values of c.Hence decision region plot is plotted for  $c=1$ .

### Decision region plot:

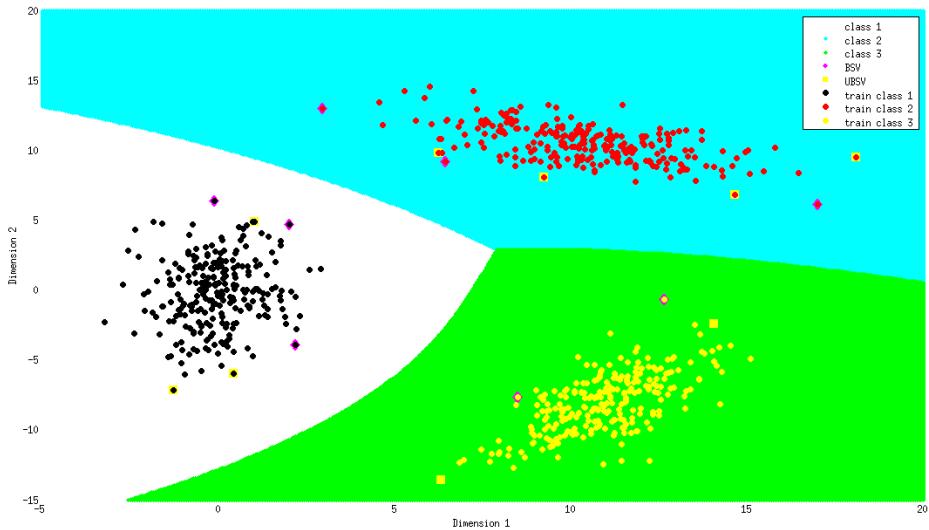


Figure 10: Figure showing decision region plot for dataset 1a gaussian kernel for  $\sigma = 100, c=1$

The shape of decision boundary is non linear in case of gaussian kernel unlike linear or polynomial kernel where it was linear. We also see the difference in number of support vectors in linear and gaussian kernel for same  $c$ .

### Confusion matrix:

Confusion Matrix				
Output Class	Target Class			
	1	2	3	
1	100 33.3%	0 0.0%	0 0.0%	100% 0.0%
2	0 0.0%	100 33.3%	0 0.0%	100% 0.0%
3	0 0.0%	0 0.0%	100 33.3%	100% 0.0%
	100% 0.0%	100% 0.0%	100% 0.0%	100% 0.0%

Figure 11: Figure showing confusion matrix for test set for dataset 1a gaussian kernel for  $\sigma = 100, c=1$

The accuracy is observed to be 100% for test data using gaussian kernel.

## 1.2 Logistic Regression

In logistic regression the posterior probabilities of the class are ensured to be between 0 and 1 and the sum over all classes is 1 by using a softmax acting on a linear function. Cross entropy was used as the error function. And the python library sci-kit learn was used for this implementation.

**Parameter estimation:** Logistic Regression was done on the dataset with polynomial basis function. The parameter that was estimated was degree of the basis expansion. We tried varying the degree between 0 and 3, since this dataset is a linearly separable dataset, we got perfect results on the test set with degree 1. Here we present our findings.

The error on the validation set vs the degree was plot and that was used to select what is the best model.

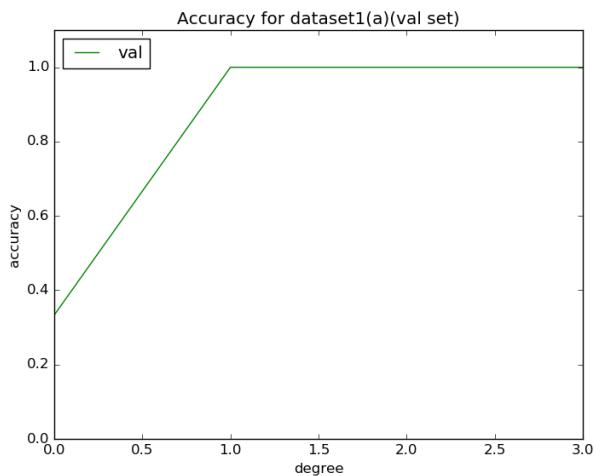


Figure 12: Figure showing accuracy for validation set for different values of degree

We see that accuracy=1 for all values of degree greater than equal to 1. The reason for high accuracy is because the dataset is linearly separable and basis expansion is not needed.

### Decision region plot:

The decision region plot for degree=1 looks like:

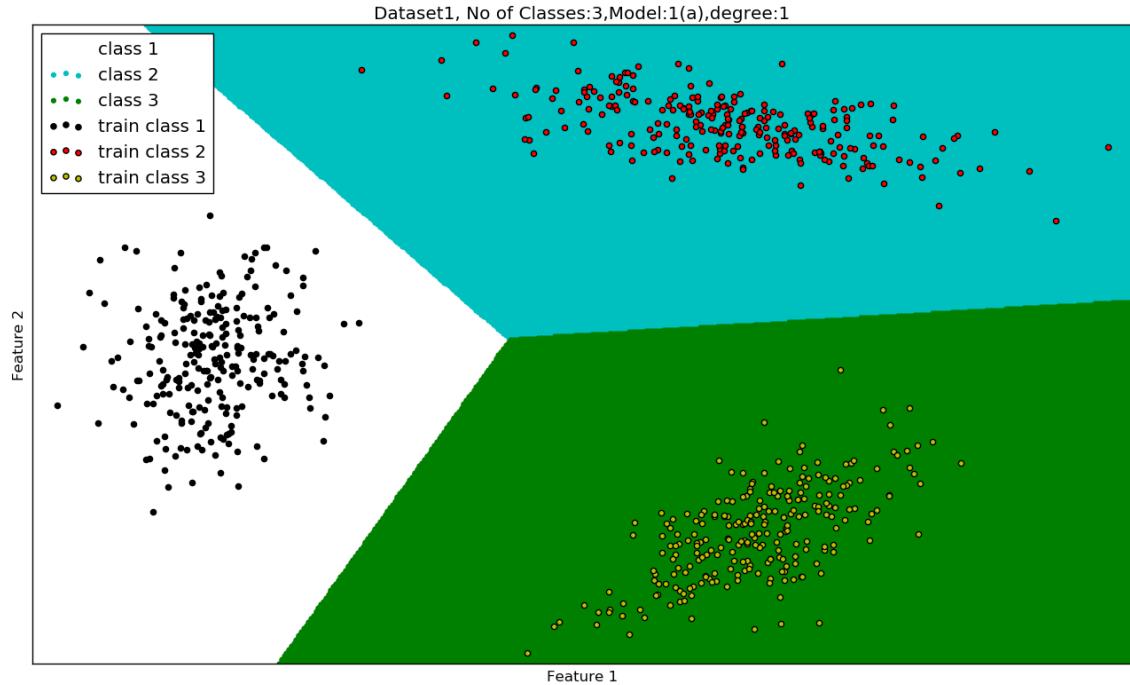


Figure 13: Figure showing decision region plot for dataset 1a using logistic regression with polynomial expansion of degree equal to 1

From the figure we see that, the decision boundaries separating each pair of classes are lines(linear),and the points in the train set are all classified correctly.

#### Confusion matrix:

The confusion matrix for test set look like:

		Confusion Matrix		
		1	2	3
Output Class	1	100 33.3%	0 0.0%	0 0.0%
	2	0 0.0%	100 33.3%	0 0.0%
3	0 0.0%	0 0.0%	100 33.3%	100% 0.0%
	1	100% 0.0%	100% 0.0%	100% 0.0%

Figure 14: Figure showing confusion matrix for test set

We see that *accuracy* = 100% on test set.

### 1.3 MLFFNN with 2 hidden layers

Multilayer feed forward with 2 hidden layers was trained. The hidden layers uses tanh activation function, and the output layer uses softmax. Cross entropy error function is used.

**Model selection:** The no of hidden layer nodes were varied from 2 to 10, for the 2 hidden layers, and the best model according to cross validation was obtained for the two hidden layers having 3 nodes each. The neural network model is able to learn the decision boundary with less no of nodes, because the boundary is linear which is less complex.

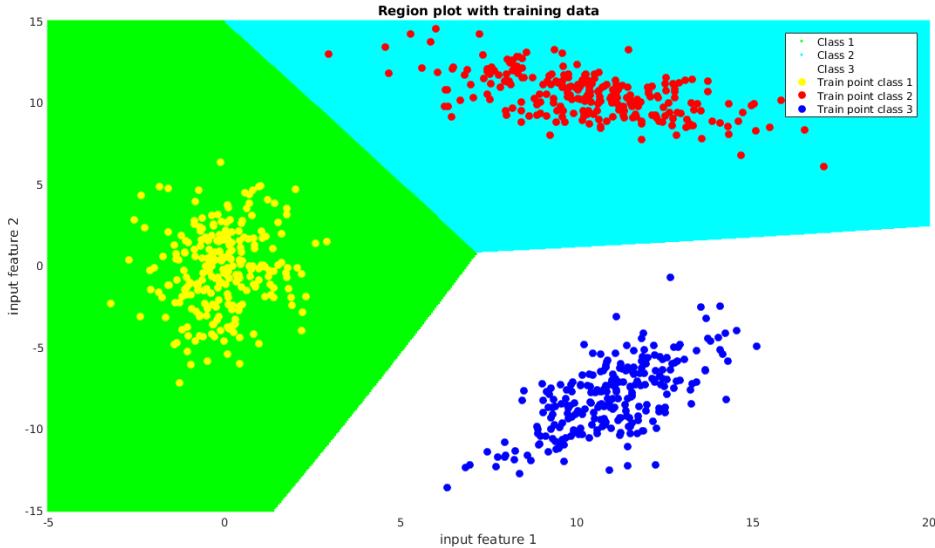


Figure 15: Region plot for dataset linearly separable data set - 1(a)

The model clearly separates the different classes as we can see from the surface plots.

Confusion Matrix				
Output Class	Target Class			
	1	2	3	
1	100 33.3%	0 0.0%	0 0.0%	100% 0.0%
2	0 0.0%	100 33.3%	0 0.0%	100% 0.0%
3	0 0.0%	0 0.0%	100 33.3%	100% 0.0%
	100% 0.0%	100% 0.0%	100% 0.0%	100% 0.0%

Figure 16: Confusion Matrix plot for dataset 1(a)

We get 100% accuracy for the data set, because it is linearly separable and since neural network is capable learning linear decision boundary.

## 2 Non-linearly separable data set

### 2.1 SVM

#### Polynomial kernel

##### Parameter estimation

To choose the kernel parameters we use kernel gram matrices. As we plotted kernel gram matrices for various values of "p" it was observed that it was performing well for even values of "p". The data across is actually separated by a circular surface when visualized. Hence "p" is chosen to be 2. Similarly for  $a=1, b=0.3$  the kernel gram matrix looks best. Hence the above are chosen as kernel parameters.

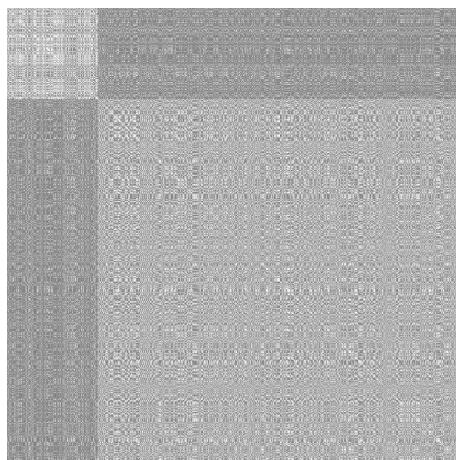


Figure 17: Figure showing grey level image of kernel gram matrix for kernel parameters  $a=1, b=0.3, p=2$ , dataset 1b

The plot for accuracy on validation vs  $\log c$  looks like:

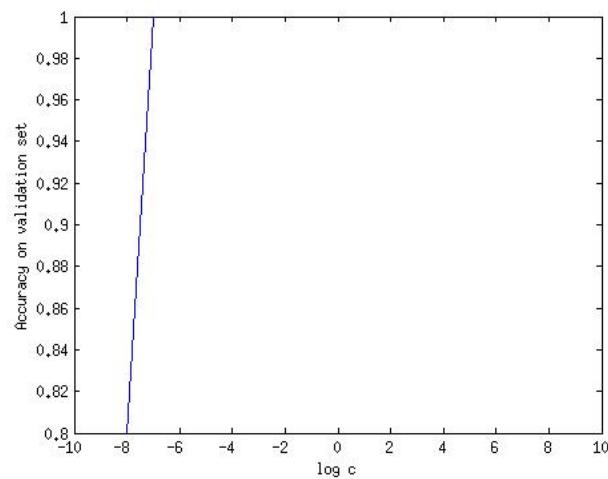


Figure 18: Figure showing accuracy for validation set for different values of  $c$

We see that accuracy is 1 for any value of  $c$  above  $2^{-6}$ , lower values of " $c$ " leads to too many support vector as margin is large. Hence we choose " $c$ " to be 1.

#### Decision region plot

The decision region plot for  $c=1$  looks like:

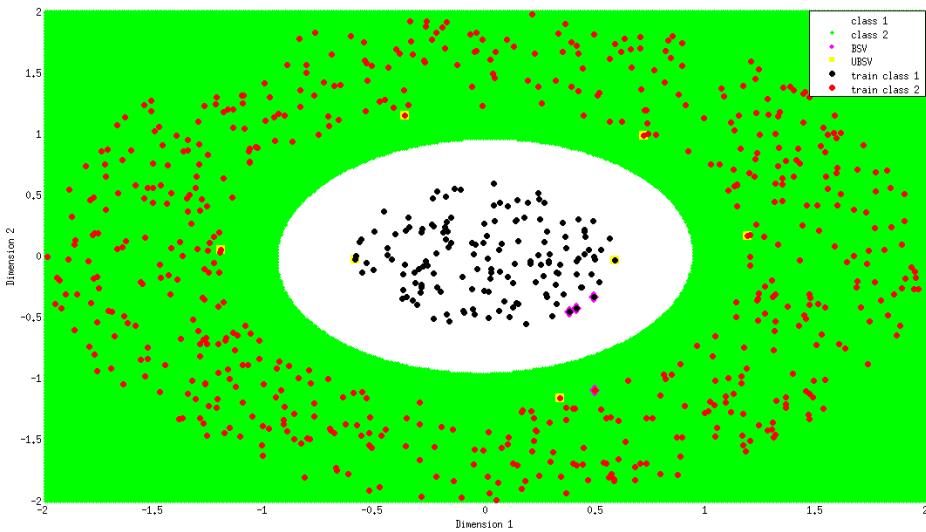


Figure 19: Figure showing decision region plot for dataset 1b using polynomial kernel,c=1

We see that unbounded support vectors are those which lie on margin. Since margin is not so large("c" is 1), we have only few points which are within margin i.e bounded support vectors.

#### Confusion matrix

The confusion matrix for dataset 1a using polynomial kernel and parameters  $a=1, b=0.3, p=2, c=1$  looks like:

Confusion Matrix		
Output Class	Target Class	
	1	2
1	60 20.0%	0 0.0%
2	0 0.0%	240 80.0%
	100% 0.0%	100% 0.0%
	100% 0.0%	100% 0.0%

Figure 20: Figure showing confusion matrix for test set for dataset 1b polynomial kernel, $a=1, b=0.3, p=2, c=1$

$$\text{accuracy} = 100\%$$

## Gaussian Kernel

### Parameter estimation:

$\sigma$  is chosen by kernel gram matrices.Best value of  $\sigma$  is found to be 3.

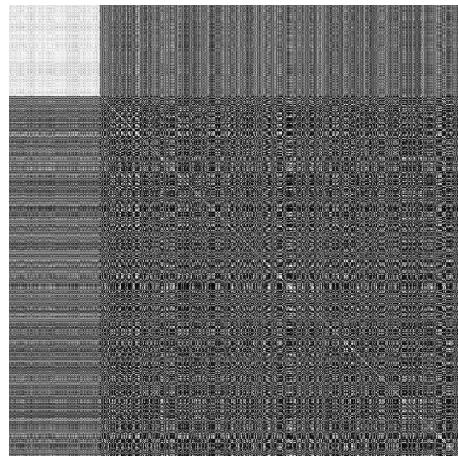


Figure 21: Grey level image of kernel gram matrix for dataset 1b training set using gaussian kernel for  $\sigma = 3$

We can see the class imbalance in the above figure.

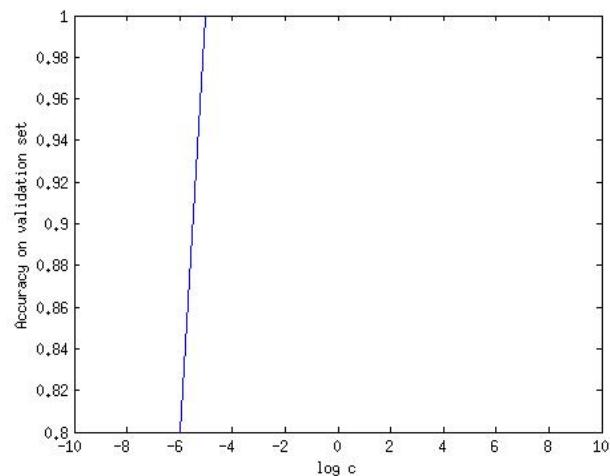


Figure 22: Figure showing accuracy for validation set for different values of c

We see that accuracy on validation set is 100% for wide range of "c".We choose a c such that number of support vectors is less."c" chosen is 1 here.

### Decision region plot:

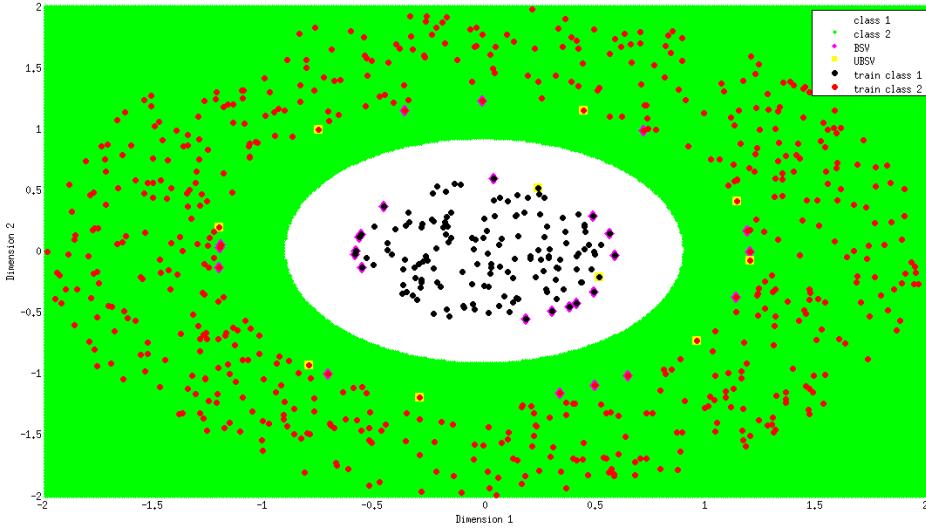


Figure 23: Figure showing decision region plot for dataset 1b gaussian kernel for  $\sigma = 3, c=1$

We see that number of support vectors are more compared to that in polynomial kernel for same value of "c" both because of change in margin and also because of change in shape of decision boundary

**Confusion matrix:**

Confusion Matrix			
Output Class	Target Class		
	1	2	
1	<b>60</b> 20.0%	<b>0</b> 0.0%	<b>100%</b> <b>0.0%</b>
2	<b>0</b> 0.0%	<b>240</b> 80.0%	<b>100%</b> <b>0.0%</b>
	<b>100%</b> <b>0.0%</b>	<b>100%</b> <b>0.0%</b>	<b>100%</b> <b>0.0%</b>

Figure 24: Figure showing confusion matrix for test set for dataset 1b gaussian kernel for  $\sigma = 3, c=1$

## 2.2 Logistic Regression

In logistic regression the posterior probabilities of the class are ensured to be between 0 and 1 and sum over all classes is 1 by using a softmax acting on a linear function. Cross entropy was used as the error function. And the python library sci-kit learn was used for this implementation.

**Parameter estimation:** Logistic Regression was done on the dataset with polynomial basis function. The parameter that was estimated was degree of the basis expansion. We tried varying

the degree between 0 and 6, since this dataset is not linearly separable dataset, we got perfect results on the test set with degree greater than 1, i.e 2. Here we present our findings.

The error on the various sets vs the degree was plot. Here we notice only one line as all the three overlap. The validation set vs degree was used to select what is the best model. We notice a maximum accuracy of 1 on validation set when degree is 2. Hence model with degree 2 is chosen as the best model.

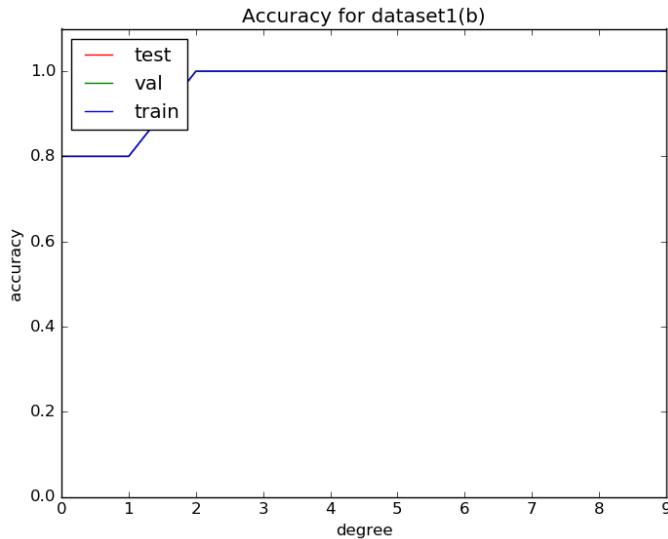


Figure 25: Figure showing accuracy for train,test and valid set for different values of degree

We see that accuracy=1 for all values of degree greater than equal to 2.

#### Decision region plot:

The decision region plot for degree=2 looks like:

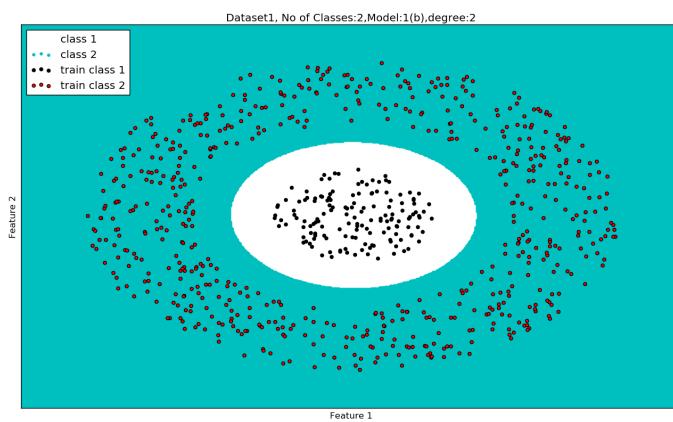


Figure 26: Figure showing decision region plot for dataset 1b using logistic regression with polynomial expansion of degree equal to 2

From the figure we see that, the decision boundaries separating each pair of classes are circular/elliptical indicating that a separation function of degree 2 is involved, and the points in the train set are all classified correctly.

#### Confusion matrix:

The confusion matrix for test set look like:

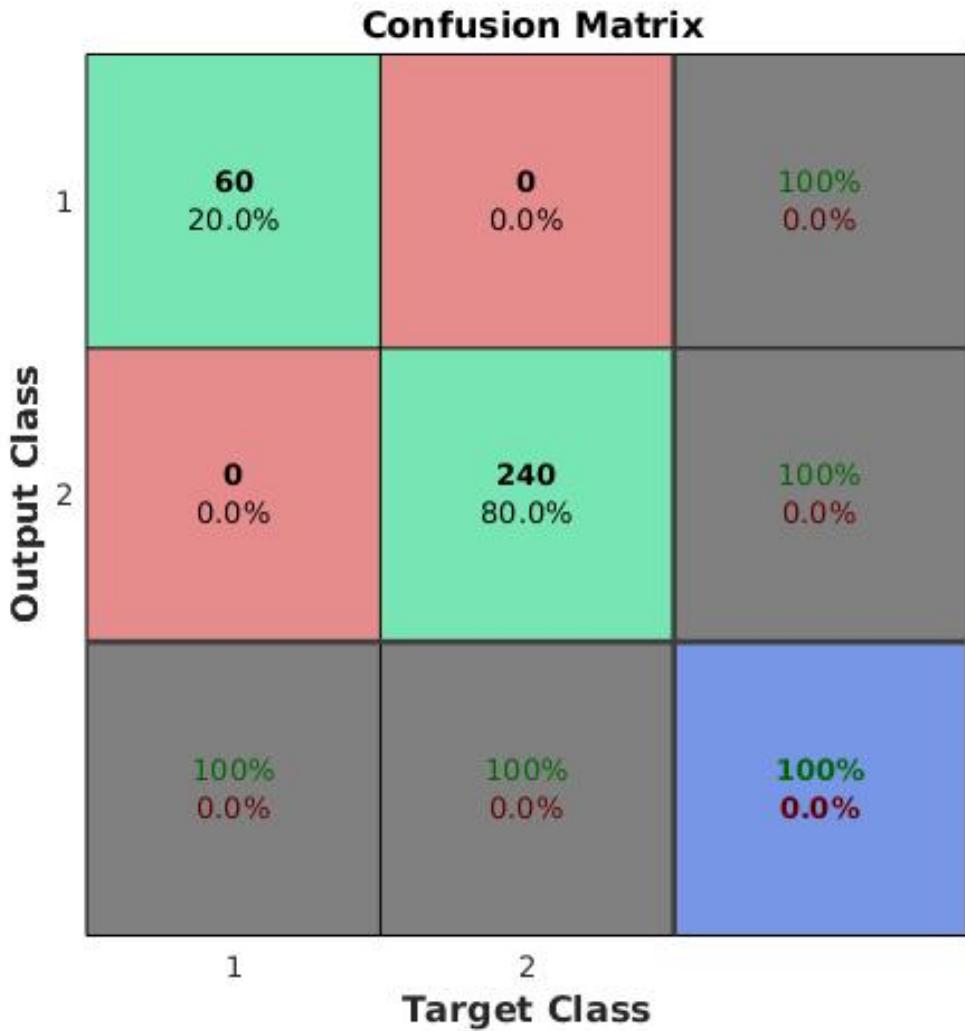


Figure 27: Figure showing confusion matrix for test set

We see that  $accuracy = 100\%$  on test set.

### 2.3 MLFFNN for 2 hidden layers

The neural network trained consists of 2 hidden layers. The hidden layers use tanh activation function, and the output layer use softmax function. Cross entropy error function is used.

**Model selection:** The no of nodes in the two hidden layers was varied from 2 to 15. The best model selected using cross validation function was for 6, 4 nodes in hidden layers 1 and 2 respectively.

**Decision region plot:**

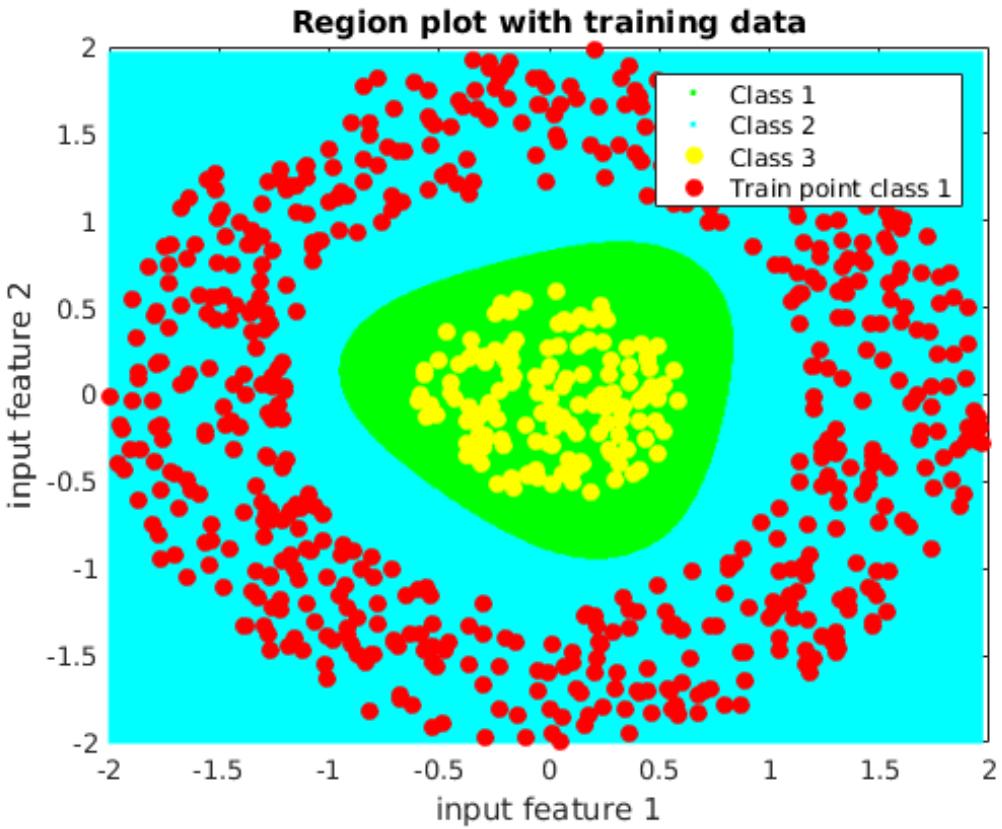


Figure 28: Decision region plot for Dataset 1(b)

As we can see from the region plot the decision boundary separates the 2 classes without any misclassification.

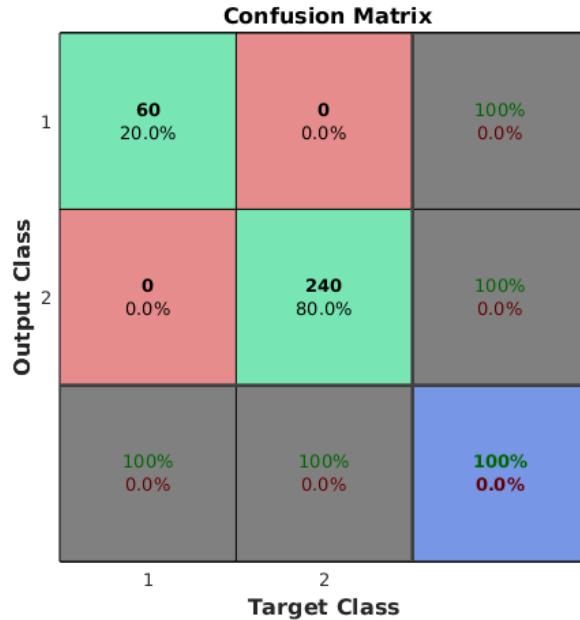


Figure 29: Confusion matrix plot for Dataset 1(b)

We get a test accuracy of 100% in this case because the dataset itself is separable and the neural network model is capable of learning non linear decision boundaries. The model leaned

here is slightly more complex (higher no of nodes in the two layers) than that learned in linearly separable dataset case. This is because the decision boundary was linear in the former case, but non linear in this case. And so requires higher complex model to learn the boundary.

### Hidden Layer 2 outputs in between training:

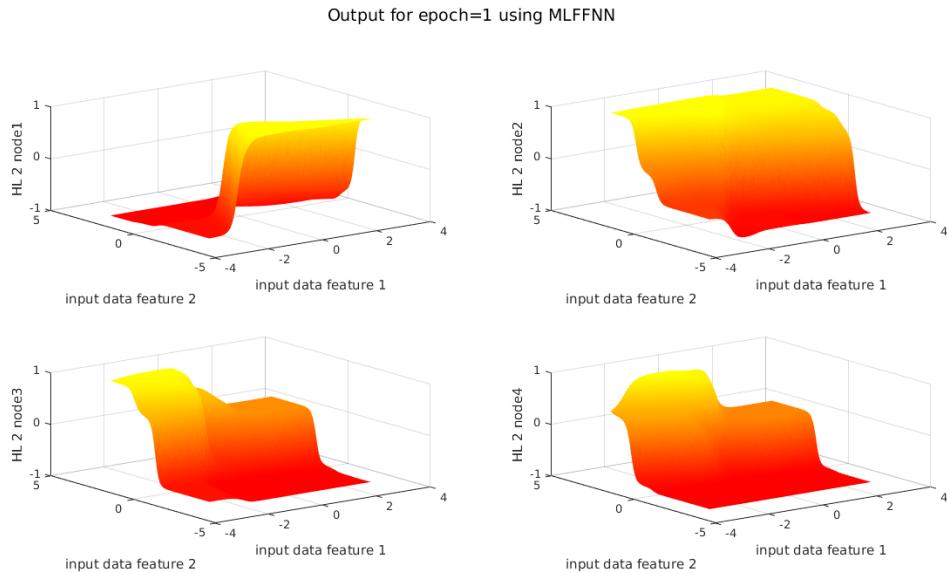


Figure 30: Hidden layer 2 outputs after first epochs

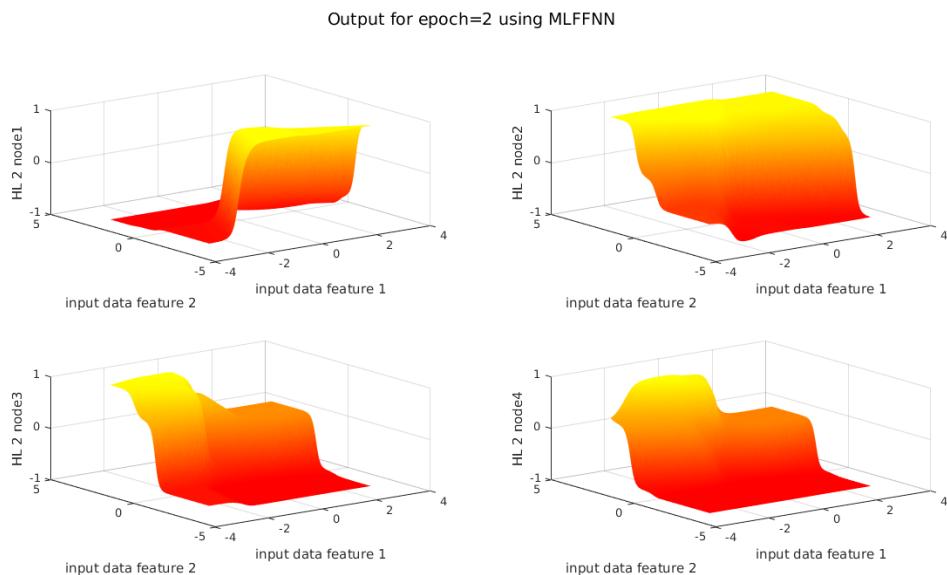


Figure 31: Hidden layer 2 outputs after 2 epochs

Output for epoch=10 using MLFFNN

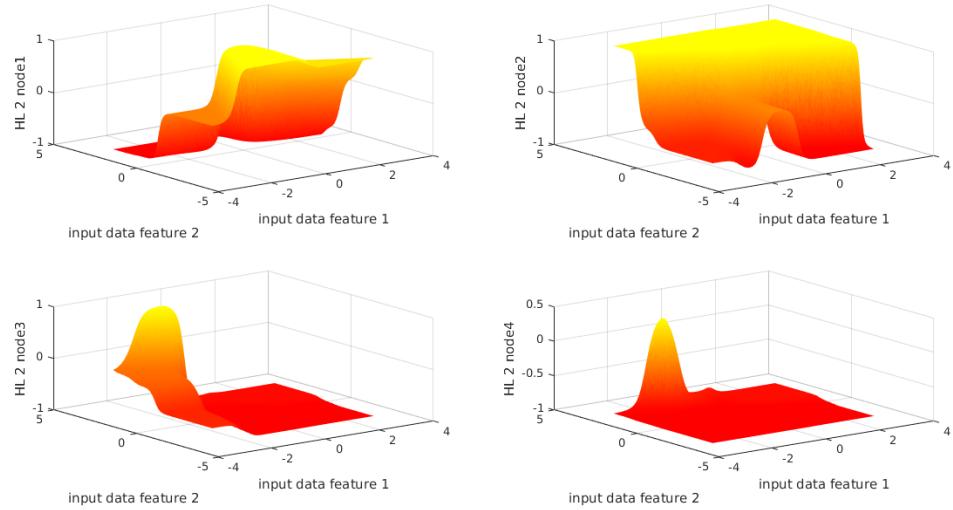


Figure 32: Hidden layer 2 outputs after 10 epochs

Output for epoch=50 using MLFFNN

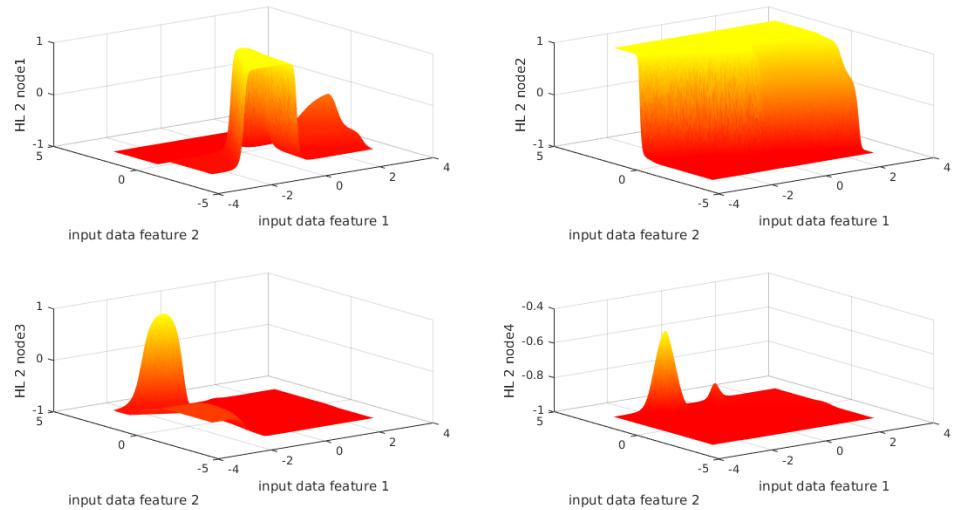


Figure 33: Hidden layer 2 outputs after 50 epochs

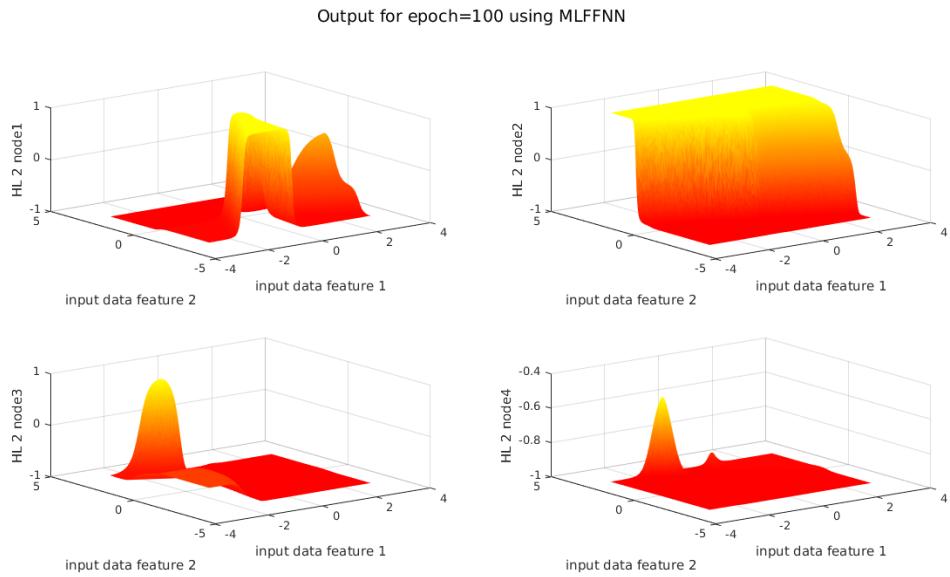


Figure 34: Hidden layer 2 outputs after 100 epochs

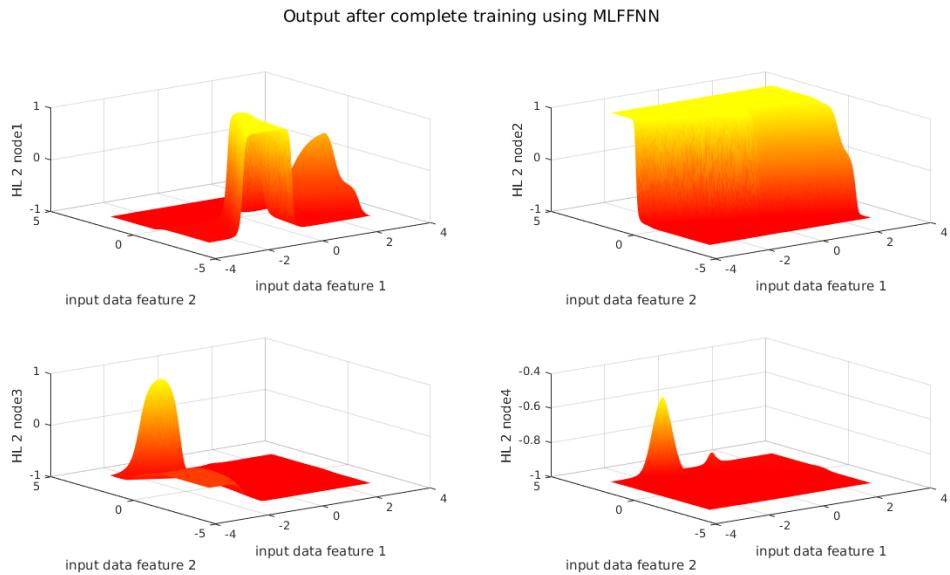


Figure 35: Second Hidden layer outputs after complete training of model

From the plots of hidden layer 2, we can see that each of the nodes slowly learn the combination of different s-shape surfaces.

#### **Output layer plots:**

Output for epoch=1 using MLFFNN

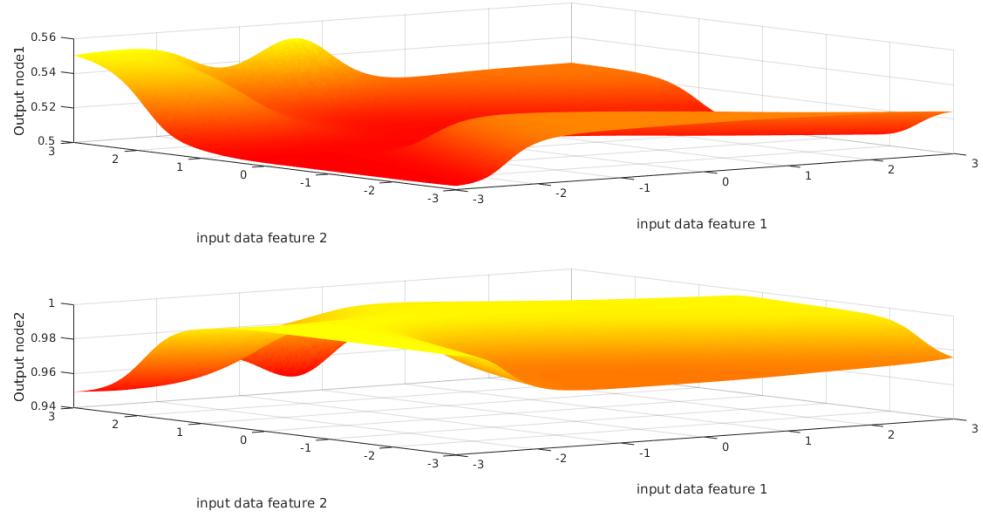


Figure 36: Output layer outputs after 1<sup>st</sup> epoch

Output for epoch=2 using MLFFNN

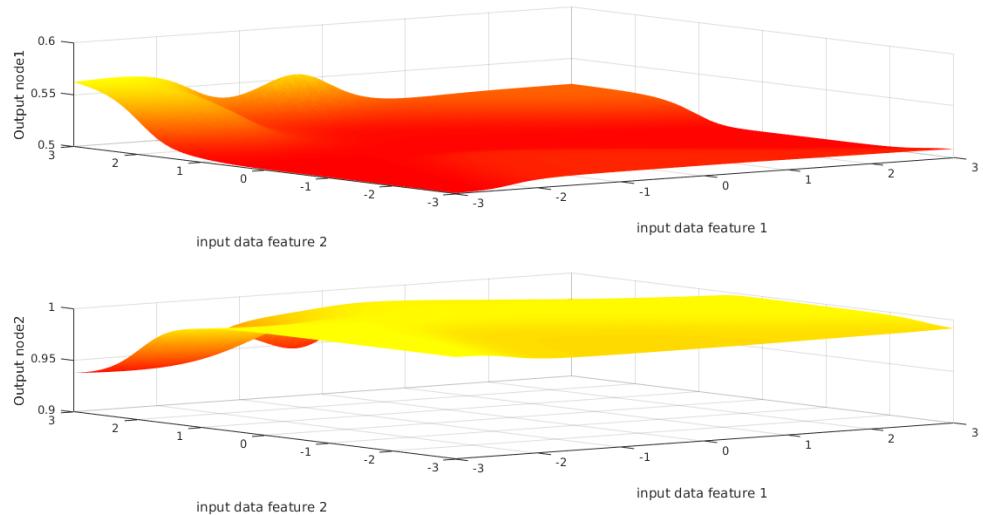


Figure 37: Output layer outputs after 2 epochs

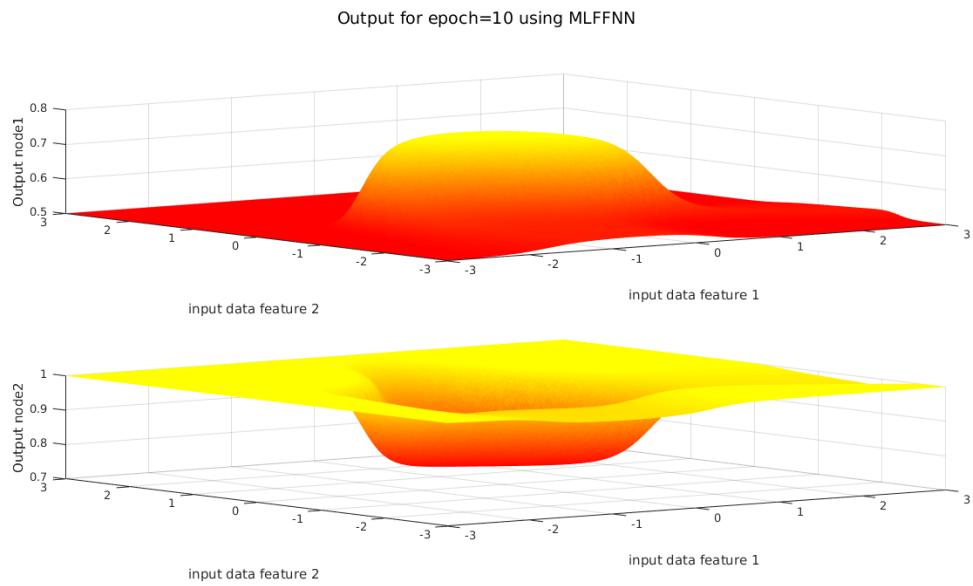


Figure 38: Output layer outputs after 10 epochs

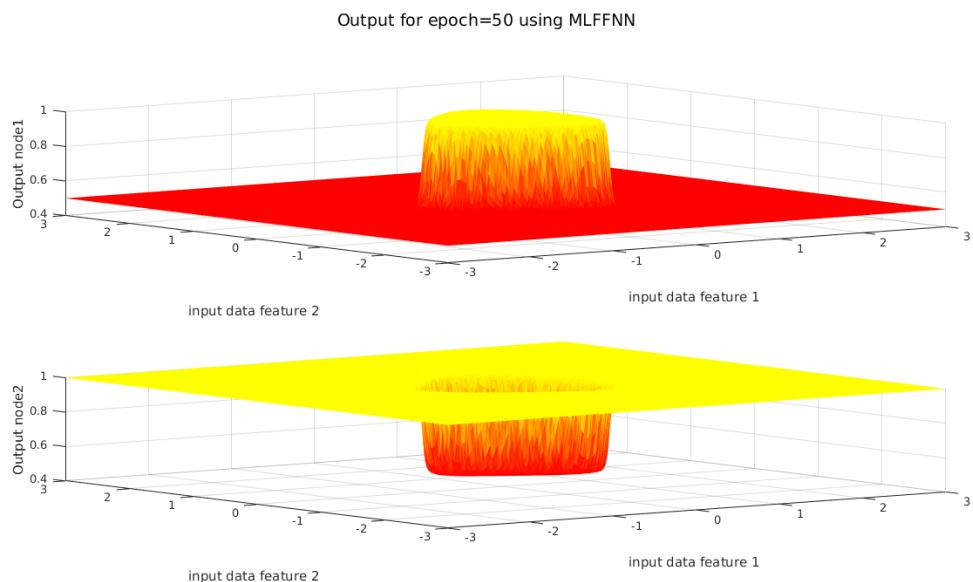


Figure 39: Output layer outputs after 50 epochs

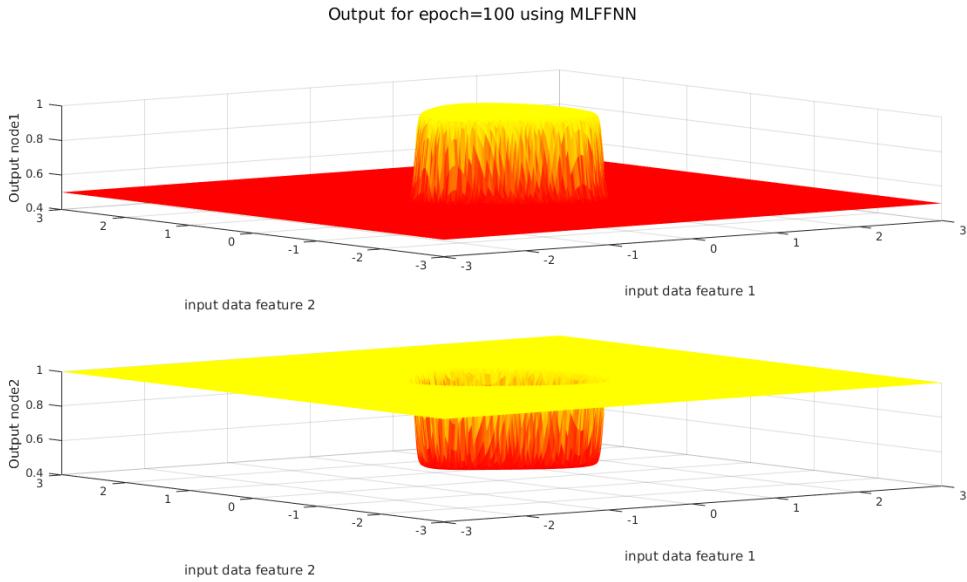


Figure 40: Output layer outputs after 100 epochs

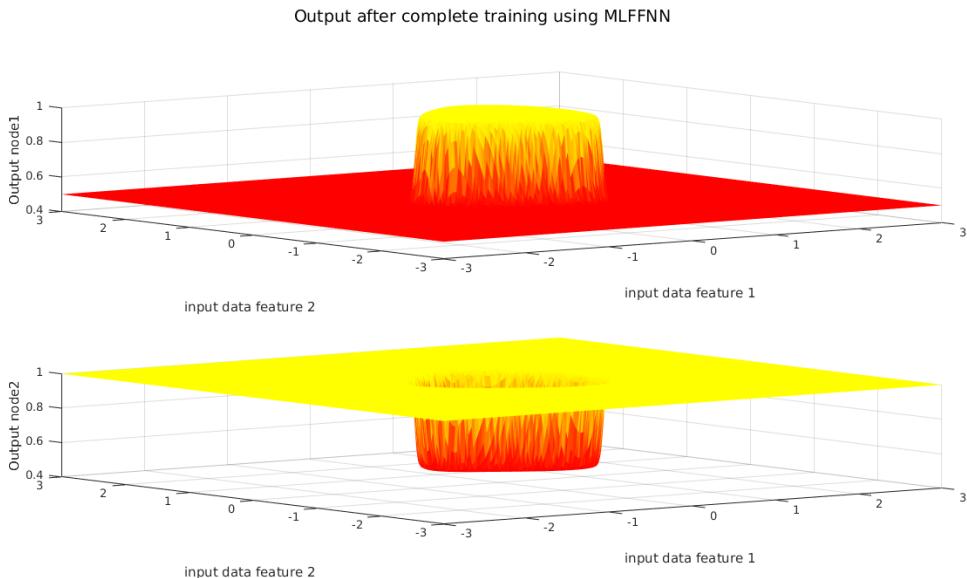


Figure 41: Output layer outputs after complete training of model

The first node outputs 1 for the points inside the circular region as observed from the plot. And the 2<sup>nd</sup> node outputs 1 for points outside the circular region. The circular region where there is a transition from 0 to 1 or 1 to 0 in the two nodes depict the decision boundary.

### Comparison between the models

From the decision region plots we can see that SVM, polynomial regression learns a perfect circular boundary whereas neural network model doesn't. This is because in SVM, maximal separating surface is learned. And in polynomial regression the surface is being tried to fit a polynomial whereas in MLFFNN there is no condition imposed on non linearity except that the decision boundary should separate the two classes.

### 3 Overlapping data set

#### 3.1 SVM

##### Polynomial kernel

###### Parameter estimation

We see that the dataset is linearly non-separable. Hence we choose  $p=1, a=1, b=0$  as kernel parameters. The kernel gram matrix for above configuration looks like:

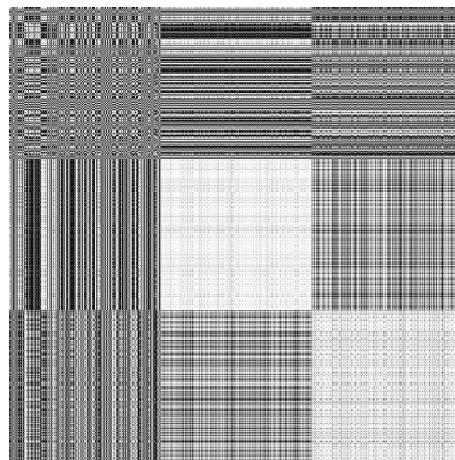


Figure 42: Figure showing grey level image of kernel gram matrix for kernel parameters  $a=1, b=0, p=1$

Since we used the same parameters for dataset 1a also, we see that, in this case for few examples same class kernel function values are lower than 1a (grey patch in diagonal blocks), for few examples kernel function values for different classes kernel function values are higher (non diagonal blocks are less darker), because of overlapping.

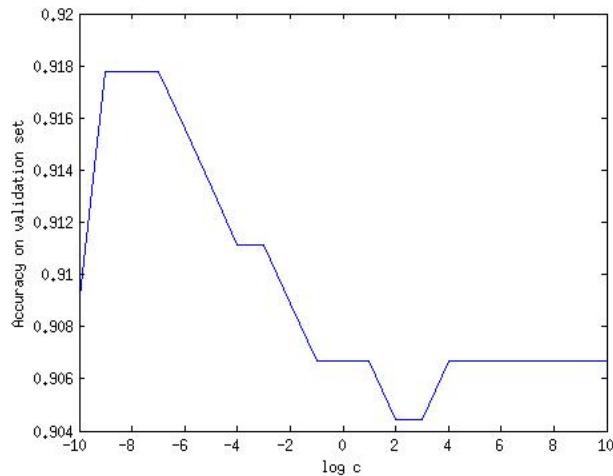


Figure 43: Figure showing accuracy for validation set for different values of  $c$

We obtain best value of  $c$  as 0.0078.

**Decision region plot:**

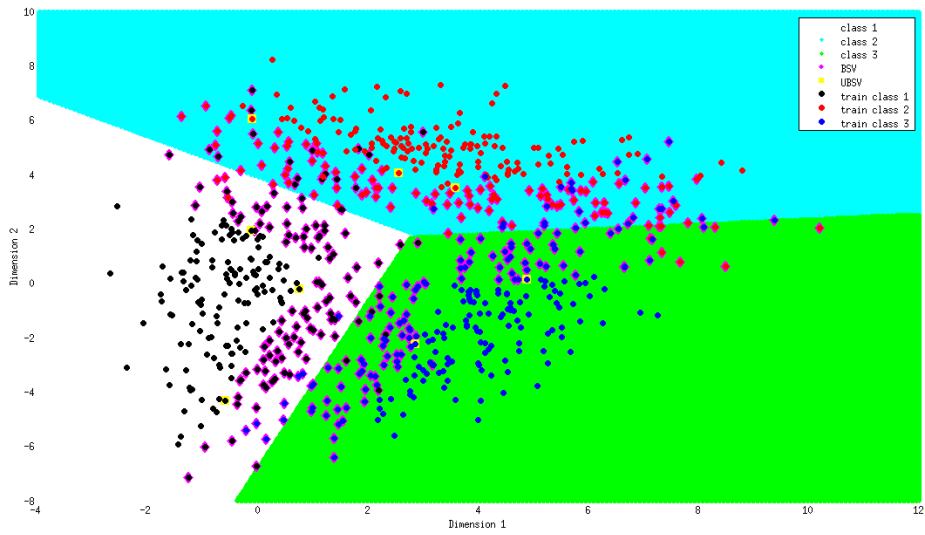


Figure 44: Figure showing decision region plot for dataset 1c using polynomial kernel,a=1 b=0 p=1 c=0.0078

Since "c" is small,we find that margin is larger,hence more points lie within the margin,number of bounded support vectors is high.

#### Confusion matrix

The confusion matrix for dataset 1c using polynomial kernel and parameters a=1,b=0,p=1 ,c=0.0078 looks like:

Confusion Matrix				
Output Class	Target Class			
	1	2	3	
1	90 30.0%	1 0.3%	5 1.7%	93.8% 6.2%
2	7 2.3%	93 31.0%	3 1.0%	90.3% 9.7%
3	3 1.0%	6 2.0%	92 30.7%	91.1% 8.9%
	90.0% 10.0%	93.0% 7.0%	92.0% 8.0%	91.7% 8.3%

Figure 45: Figure showing confusion matrix for test set for dataset 1c polynomial kernel,c=0.0078

The accuracy in this case is observed to be 91.2% which is less compared to datasets 1a and 1b.

## Gaussian kernel

### Parameter estimation:

$\sigma$  value chosen is 33 after observing kernel gram matrices.

The kernel gram matrix for  $\sigma = 33$  is shown:

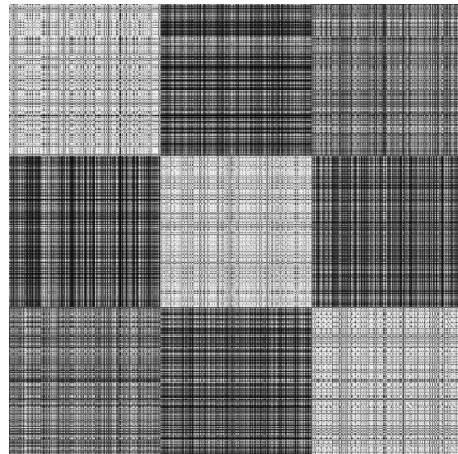


Figure 46: Grey level image of kernel gram matrix for dataset 1c training set using gaussian kernel for  $\sigma = 33$

We observe that the kernel gram matrix is better representing similarity of points of same class and dissimilarity of points of different class because it depends on distribution of data, for this distribution euclidean distance works better than cosine similarity.

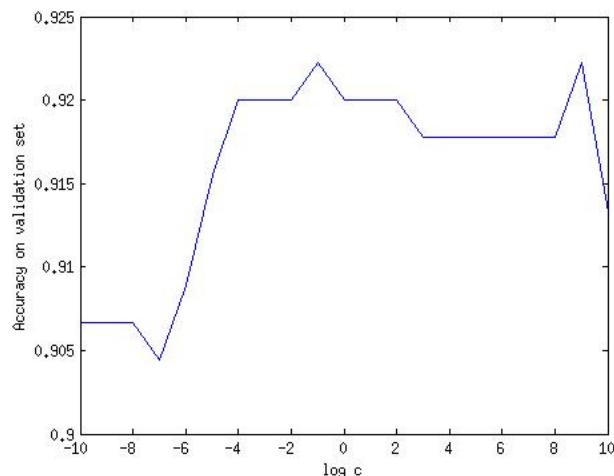


Figure 47: Figure showing accuracy for validation set for different values of  $c$

The best value for  $c$  is observed to be 0.5.

### Decision region plot:

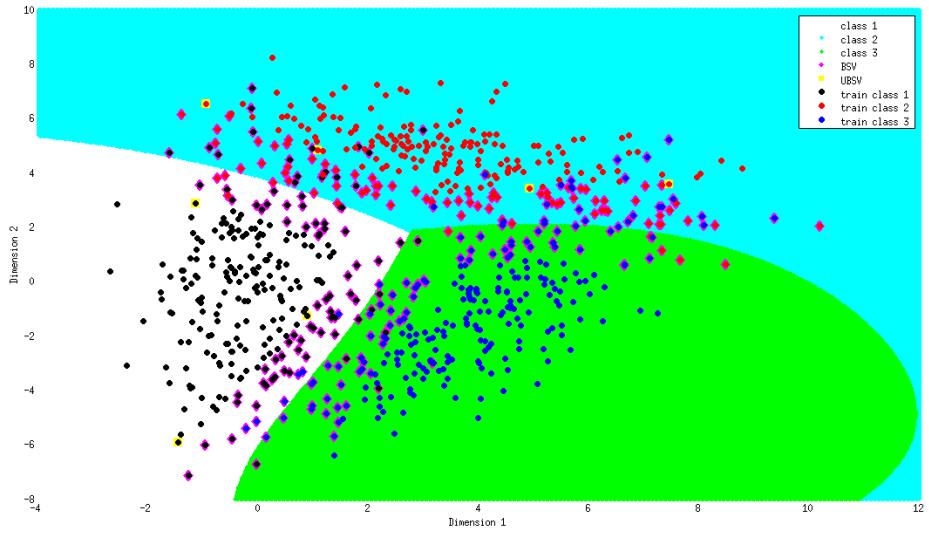


Figure 48: Figure showing decision region plot for dataset 1c gaussian kernel for  $\sigma = 33, c=1$

We see that both decision surface and number of support vectors change from polynomial to gaussian kernel. Since the value of  $c$  is observed in polynomial kernel is low, the number of support vectors is high in polynomial case.

**Confusion matrix:**

Confusion Matrix				
Output Class	Target Class			
	1	2	3	
1	<b>89</b> 29.7%	<b>1</b> 0.3%	<b>5</b> 1.7%	
2	<b>8</b> 2.7%	<b>95</b> 31.7%	<b>5</b> 1.7%	
3	<b>3</b> 1.0%	<b>4</b> 1.3%	<b>90</b> 30.0%	
	89.0% 11.0%	95.0% 5.0%	90.0% 10.0%	<b>91.3%</b> <b>8.7%</b>

Figure 49: Figure showing confusion matrix for test set for dataset 1c gaussian kernel for  $\sigma = 33, c=1$

The accuracy is found to be 91.3% which is slightly lesser than polynomial kernel which yields accuracy 91.7%

### 3.2 Logistic Regression

In logistic regression the posterior probabilities of the class are ensured to be between 0 and 1 and sum over all classes is 1 by using a softmax acting on a linear function. Cross entropy was used as the error function. And the python library sci-kit learn was used for this implementation.

**Parameter estimation:** Logistic Regression was done on the dataset with polynomial basis function. The parameter that was estimated was degree of the basis expansion. We tried varying the degree between 0 and 6, we got good results on the test set with degree 3. The best accuracy observed was 92%. Since the dataset is overlapping, a perfect accuracy seems difficult. Hence an accuracy of 92 percent implies that this model is a good model. Here we present our findings.

The error on the validation set vs degree was plot and was used to select what is the best model. We notice that all models perform almost equally between degree 1 and 4 but there is a slight maximum at 3. We notice a maximum accuracy of 92% on validation set when degree is 3. Hence model with degree 3 is chosen as the best model.

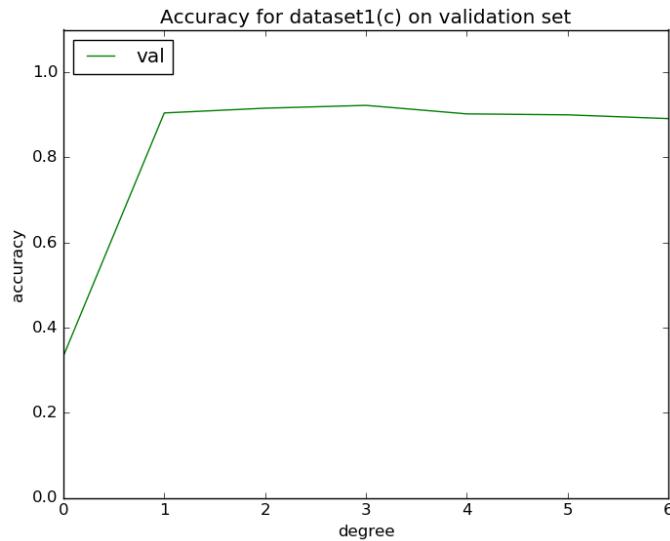


Figure 50: Figure showing accuracy of validation set for different values of degree

We see that accuracy=0.92 for degree equal to 3.

#### Decision region plot:

The decision region plot for degree=3 looks like:

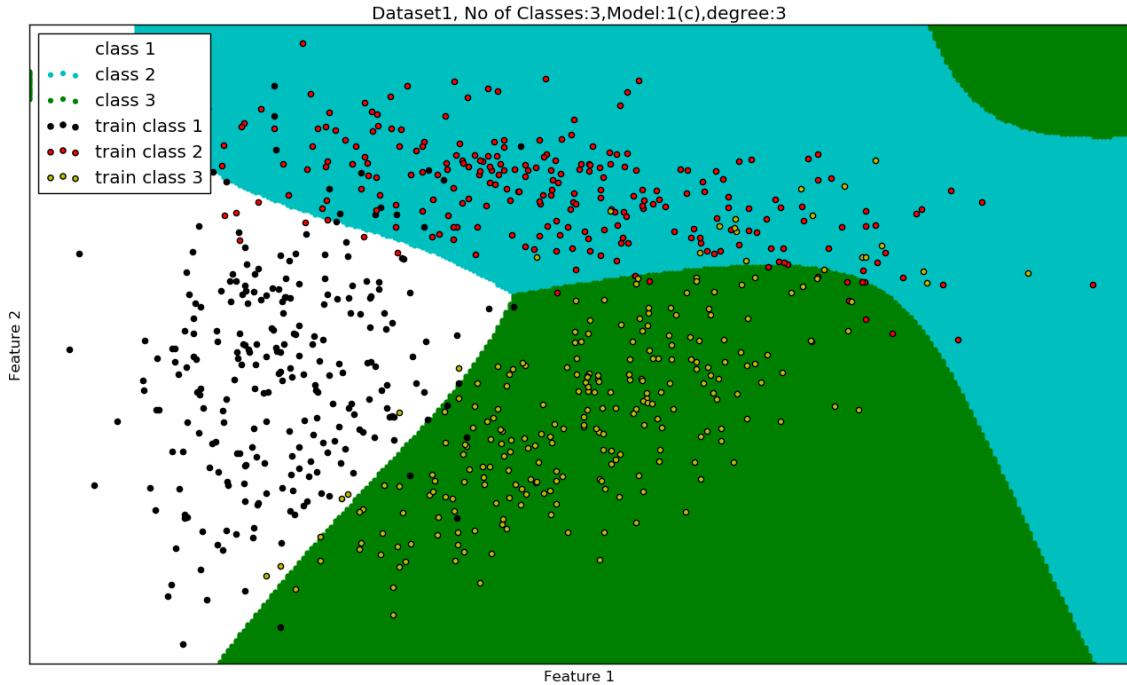


Figure 51: Figure showing decision region plot for dataset 1c using logistic regression with polynomial expansion of degree equal to 3

From the figure we see that, the decision boundaries separating each pair of classes are not linear indicating that a separation function of degree greater than 1 is involved, and the points in the train set overlap closely among them.

#### Confusion matrix:

The confusion matrix for test set look like:

Confusion Matrix			
Output Class	Target Class		
	1	2	3
1	<b>92</b> 30.7%	<b>2</b> 0.7%	<b>6</b> 2.0%
2	<b>5</b> 1.7%	<b>92</b> 30.7%	<b>2</b> 0.7%
3	<b>3</b> 1.0%	<b>6</b> 2.0%	<b>92</b> 30.7%
	92.0% 8.0%	92.0% 8.0%	92.0% 8.0%
			<b>92.0% 8.0%</b>

Figure 52: Figure showing confusion matrix for test set

We see that *accuracy* = 92% on test set.

### 3.3 MLFFNN with 2 hidden layers

MLFFNN was trained with 2 hidden layers. The hidden layers uses tanh activation function, and the output layer uses softmax activation function. Cross entropy error function was used.

**Model Selection:** For the 2 hidden layers the no of nodes were varied from 2 to 15, and the best model complexity obtained using minimum cross validation error was for 7 and 5 nodes in the hidden layers 1 and 2 respectively.

**Decision region plot:**

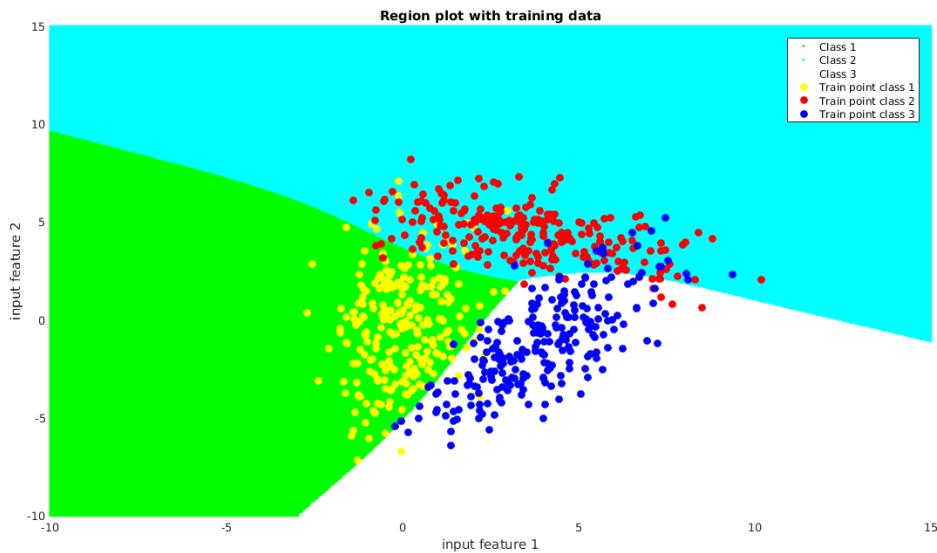


Figure 53: Decision region plot for overlapping dataset

**Confusion matrix plot:**

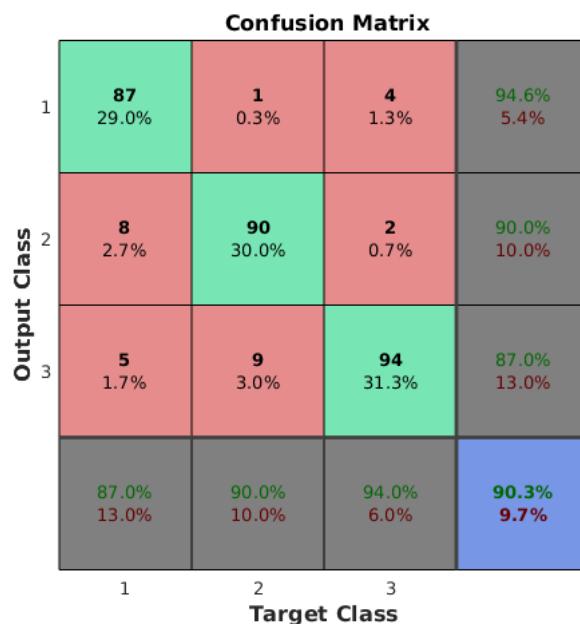


Figure 54: Confusion matrix plot for Dataset 1(c)

The classes are overlapping, neural network model learned doesn't separate the 2 classes com-

pletely. Minimizing validation error, the best model learned was for 7 nodes on first hidden layer, and 5 nodes on second hidden layer. The model gives 90.3 % accuracy on test data.

## 4 Black and white image data set

### Gaussian kernel

**Parameter estimation:**

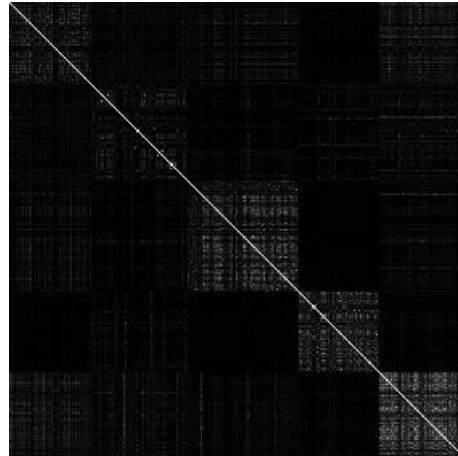


Figure 55: Figure showing grey level image of kernel gram matrix for kernel parameters , $\sigma = 60$

The best  $\sigma$  obtained is 60 from kernel gram matrix.

Here we see that some of the classes are well separable whereas other pairs of classes are not. Using CNN will have higher accuracy because it exploits local features and robust to shift errors unlike SVM.

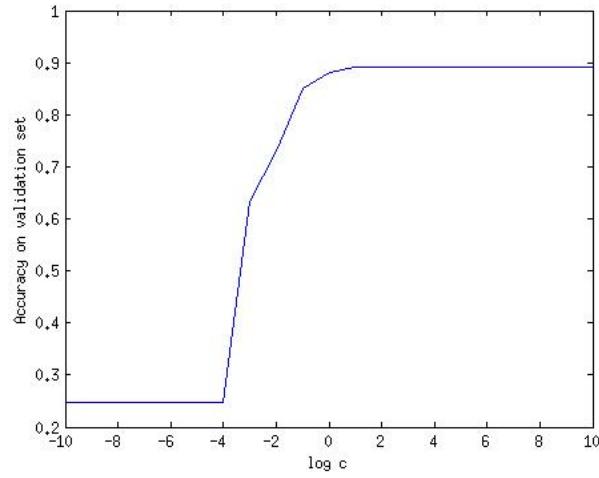


Figure 56: Figure showing accuracy for validation set for different values of  $c$

$c$  is chosen to be "1" as accuracy=100%

**Confusion matrix:**

		Confusion Matrix					
		1	2	3	4	5	
Output Class	1	16 17.4%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
	2	0 0.0%	17 18.5%	1 1.1%	0 0.0%	0 0.0%	94.4% 5.6%
	3	1 1.1%	1 1.1%	20 21.7%	0 0.0%	0 0.0%	90.9% 9.1%
	4	0 0.0%	1 1.1%	0 0.0%	16 17.4%	0 0.0%	94.1% 5.9%
	5	0 0.0%	1 1.1%	1 1.1%	0 0.0%	17 18.5%	89.5% 10.5%
		94.1% 5.9%	85.0% 15.0%	90.9% 9.1%	100% 0.0%	100% 0.0%	93.5% 6.5%
		1	2	3	4	5	Target Class

Figure 57: Figure showing confusion matrix for test set for black and white image dataset gaussian kernel for  $\sigma = 60, c=1$

The accuracy is observed to be 93.5% on test set.

We can also observe that examples of class a misclassified as class b have greyer patch in kernel gram matrix relative to those pairs of classes which are separable/not misclassified.(ex:class3,class4;class1,class2);

#### 4.1 CNN

A deep leaning convolutional neural network was used to train this dataset. Convolutional neural networks tend to perform well on image datasets, we use a two layer network here. Each layer included a combination of a filter layer that introduces non linearity through a relu function and a max pooling layer which combined nodes with some drop out probability. Finally after the two layers there was a fully connected layer that generated outputs based on the number of classes and a softmax cross entropy function in the end for multi class classification.

**Parameter estimation:** The parameters to be varied were the number of feature maps in the first and the second layer. Sir told us to keep the stride and the local receptor field constant. Hence we kept them the same and didn't consider them in out parameter estimation. The value of the LRF(Local Receptor field) was  $5 * 5$ , and the stride parameter was 1 in the convolutional layer, in the pooling layer the value of the LRF was  $2 * 2$ , and the stride parameter was 2. After 1 layer the output reduced by 50 percent in both dimensions. Here we present our findings when we vary the number of feature maps to be applied.

The error on the validation set vs number of feature maps in layer 1 was plot and was used to select what is the best model. We decided to keep the number of feature maps in layer 2 as twice the number of feature maps in layer 1 as the CNN was taking a long time to train. Hence instead of estimating on a huge field we varied parameters in first layer  $\in \{1, 5, 10, 50, 100\}$ . We noticed that the model with 10 performs the best and the models with 50 and 100 tend to overfit on the data. The number of epochs were set at 200 and the exit criteria was if the error changed by a value of less than 0.1 percent after 10 epochs we exited. We notice a maximum accuracy of 96.42% on validation set when number of feature maps in layer 1 are 100 and number of feature maps in layer 2 are 200. Hence this model was chosen as the best model.

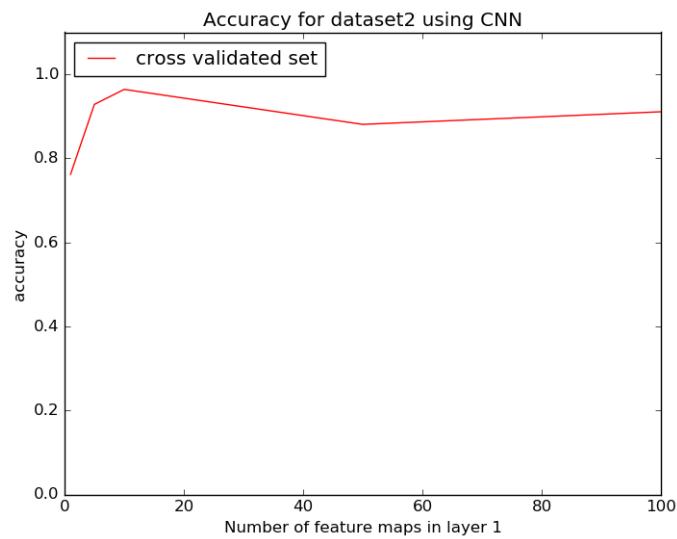


Figure 58: Figure showing accuracy on validation set for different values of feature maps in layer 1, Number of feature maps in layer 2 are twice the number of feature maps in layer 1

#### Feature Maps:

The feature maps of the two layers are shown below:



Figure 59: Input Image:Class Label 40



Figure 60: Feature Map of layer 1

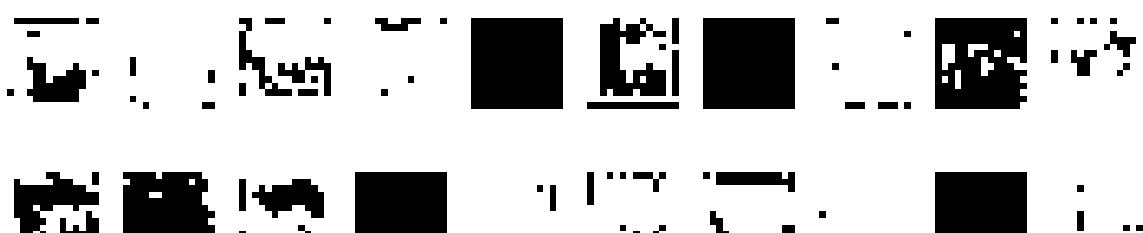


Figure 61: Feature Map of layer 2

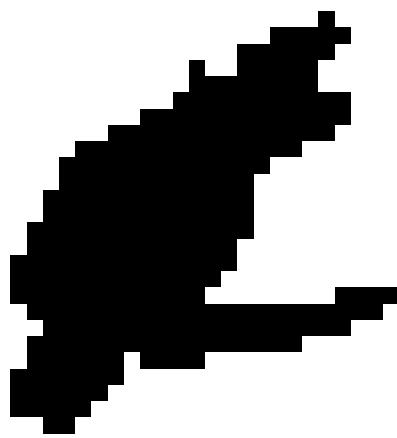


Figure 62: Input Image:Class Label 48



Figure 63: Feature Map of layer 1

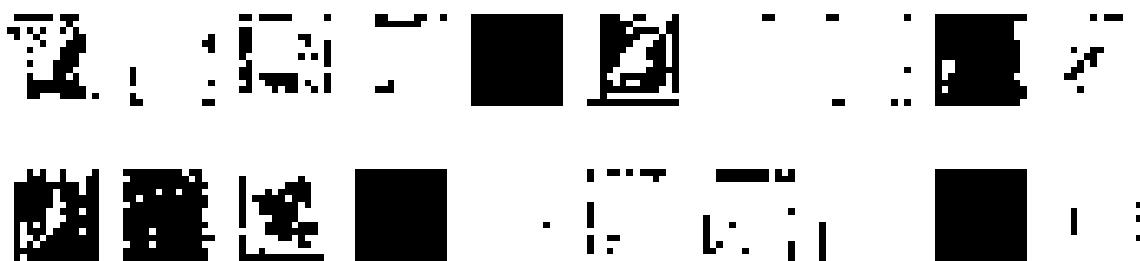


Figure 64: Feature Map of layer 2

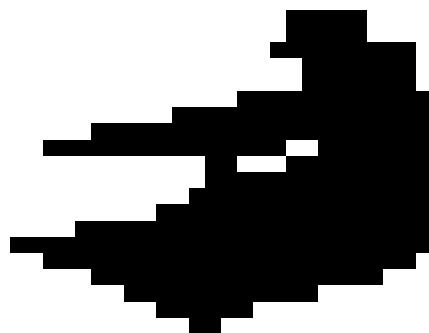


Figure 65: Input Image:Class Label 56



Figure 66: Feature Map of layer 1

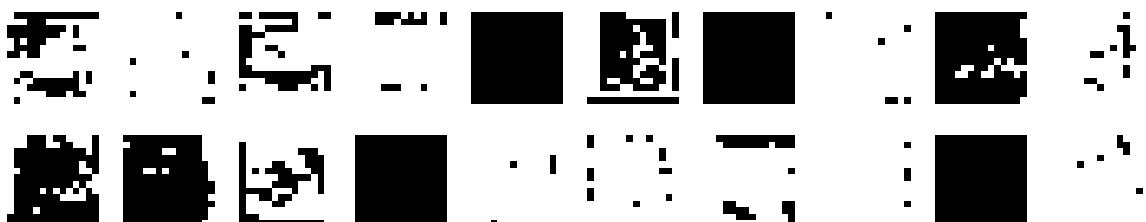


Figure 67: Feature Map of layer 2



Figure 68: Input Image:Class Label 76



Figure 69: Feature Map of layer 1

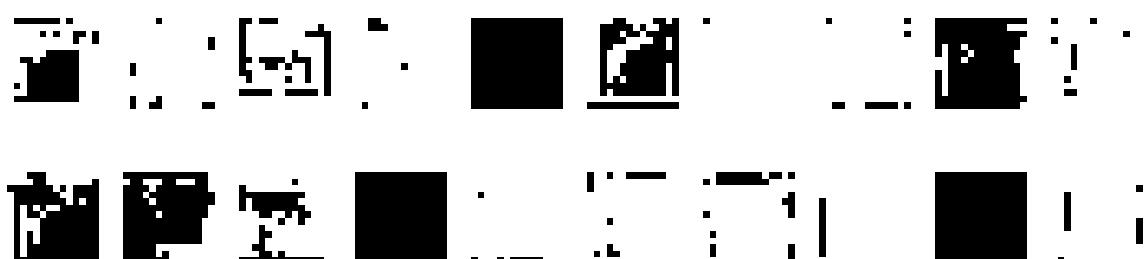


Figure 70: Feature Map of layer 2

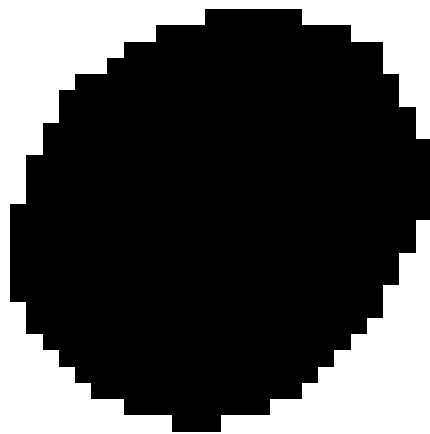


Figure 71: Input Image:Class 93



Figure 72: Feature Map of layer 1

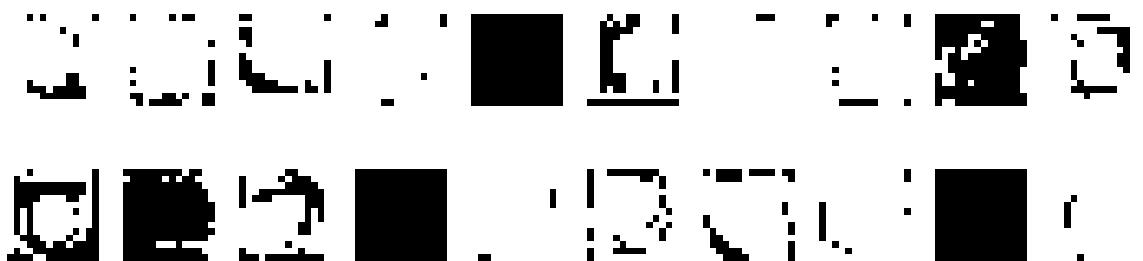


Figure 73: Feature Map of layer 2

The number of feature maps are 10 in first layer and 20 in the second layer, we display each of them for five of the output classes. It is difficult to interpret these images and output some observations from these.

**Confusion matrix:**

The confusion matrix for test set look like:

		Confusion Matrix					
		1	2	3	4	5	
Output Class	1	15 14.7%	0 0.0%	0 0.0%	0 0.0%	1 1.0%	93.8% 6.2%
	2	0 0.0%	19 18.6%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
	3	1 1.0%	3 2.9%	24 23.5%	0 0.0%	0 0.0%	85.7% 14.3%
	4	0 0.0%	1 1.0%	0 0.0%	16 15.7%	0 0.0%	94.1% 5.9%
	5	0 0.0%	0 0.0%	2 2.0%	0 0.0%	20 19.6%	90.9% 9.1%
		93.8% 6.2%	82.6% 17.4%	92.3% 7.7%	100% 0.0%	95.2% 4.8%	92.2% 7.8%
		1	2	3	4	5	

Figure 74: Figure showing confusion matrix for test set

We see that *accuracy* = 92.2% on test set.

## 4.2 MLFFNN with 2 hidden layers

MLFFNN with 2 hidden layers was trained. Cross entropy error function was used. The dataset was split into train, test and validation sets with ratio of 0.6,0.2,0.2 without any class imbalance.

**Model selection:** The hidden layer nodes were varied in high range(500- 600) ,mid range(300 - 500), and in low range(<300). The best model complexity using cross validation error was obtained for 350, and 200 nodes in the hidden layers 1 and 2 respectively. The test accuracy attained was 87%. The performance in neural network is less compared to CNN because CNN learns local feature, and the shift errors are overcome. And so the classification is better in case of CNN.

**Confusion Matrix:**

Confusion matrix for Test data							
Output Class	1	2	3	4	5		
	1	13 14.1%	1 1.1%	0 0.0%	0 0.0%	0 0.0%	92.9% 7.1%
	2	0 0.0%	14 15.2%	0 0.0%	1 1.1%	0 0.0%	93.3% 6.7%
	3	2 2.2%	1 1.1%	21 22.8%	0 0.0%	0 0.0%	87.5% 12.5%
	4	0 0.0%	1 1.1%	0 0.0%	15 16.3%	0 0.0%	93.8% 6.2%
	5	2 2.2%	3 3.3%	1 1.1%	0 0.0%	17 18.5%	73.9% 26.1%
76.5% 23.5%		70.0% 30.0%	95.5% 4.5%	93.8% 6.2%	100% 0.0%	87.0% 13.0%	

Figure 75: Confusion matrix plot for black and white dataset for test data

### 4.3 MLFFNN pre-trained with auto-encoder

Dataset: Dataset consists of 784 feature of 467 data points. The given data was split in the ratio of 60%, 20%, 20% of Train, Validation, Test data. The data was split without class imbalance. Model selection: MLFFNN with 3 hidden layers were to be leaned. To arrive at the best model configuration, some configurations in the high range(500 - 600 nodes), some in the middle(300 - 500 nodes), and some in the low range(< 200 nodes) were tried.

The model in the middle range gave minimum validation error of 86.7%, for which the test error was 91.3%.

Confusion matrix for Test data							
Output Class	1	2	3	4	5		
	1	16 17.2%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
	2	1 1.1%	20 21.5%	0 0.0%	2 2.2%	1 1.1%	83.3% 16.7%
	3	1 1.1%	0 0.0%	22 23.7%	0 0.0%	0 0.0%	95.7% 4.3%
	4	0 0.0%	0 0.0%	0 0.0%	12 12.9%	0 0.0%	100% 0.0%
	5	2 2.2%	1 1.1%	0 0.0%	0 0.0%	15 16.1%	83.3% 16.7%
80.0% 20.0%		95.2% 4.8%	100% 0.0%	85.7% 14.3%	93.8% 6.2%	91.4% 8.6%	

Figure 76: Confusion matrix of Test set of black and white data on MLFFNN pre-trained using auto-encoder

The pre-trained model is slightly better than the model of same configuration but without pre-training.

		Confusion matrix for Test data					
		1	2	3	4	5	
Output Class	1	15 16.3%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
	2	0 0.0%	15 16.3%	1 1.1%	1 1.1%	0 0.0%	88.2% 11.8%
	3	1 1.1%	0 0.0%	20 21.7%	0 0.0%	0 0.0%	95.2% 4.8%
	4	0 0.0%	1 1.1%	0 0.0%	15 16.3%	0 0.0%	93.8% 6.2%
	5	1 1.1%	4 4.3%	1 1.1%	0 0.0%	17 18.5%	73.9% 26.1%
		88.2% 11.8%	75.0% 25.0%	90.9% 9.1%	93.8% 6.2%	100% 0.0%	89.1% 10.9%

Figure 77: Confusion matrix of Test set of black and white data on MLFFNN without any pre-training

A decrease in no of iterations for convergence is observed when pretraining is done. And the train error is slightly better in case of pretraining.

#### 4.4 MLFFNN pretrained with RBM

**Model selection:** 3 level stacked RBM was used. The 3 RBMs were leaned separately. The weights were fed to same configuration neural network and trained. Some configuration in high range, medium, and low range of nodes were tried. Out of which the best was selected using cross validation which was the one with 450, 400, 300 nodes in Hidden layers 1 , 2 and 3 respectively.

Confusion matrix for Test data							
Output Class	1	2	3	4	5		
	1	17 18.5%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%	
	2	0 0.0%	15 16.3%	1 1.1%	0 0.0%	0 0.0%	93.8% 6.2%
	3	0 0.0%	1 1.1%	19 20.7%	0 0.0%	0 0.0%	95.0% 5.0%
	4	0 0.0%	2 2.2%	0 0.0%	16 17.4%	0 0.0%	88.9% 11.1%
	5	0 0.0%	2 2.2%	2 2.2%	0 0.0%	17 18.5%	81.0% 19.0%
		100% 0.0%	75.0% 25.0%	86.4% 13.6%	100% 0.0%	100% 0.0%	91.3% 8.7%

Figure 78: Confusion matrix of Test set of black and white data on MLFFNN pre-trained using stacked RBM

Confusion matrix for Test data							
Output Class	1	2	3	4	5		
	1	14 15.2%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
	2	1 1.1%	15 16.3%	1 1.1%	0 0.0%	0 0.0%	88.2% 11.8%
	3	1 1.1%	1 1.1%	20 21.7%	0 0.0%	0 0.0%	90.9% 9.1%
	4	0 0.0%	1 1.1%	0 0.0%	16 17.4%	0 0.0%	94.1% 5.9%
	5	1 1.1%	3 3.3%	1 1.1%	0 0.0%	17 18.5%	77.3% 22.7%
		82.4% 17.6%	75.0% 25.0%	90.9% 9.1%	100% 0.0%	100% 0.0%	89.1% 10.9%

Figure 79: Confusion matrix of Test set of black and white data on MLFFNN without any pre-training

### Comparison between autoencoder and RBM

In both cases pretraining gives slightly better accuracy than without pretraining. Both methods of pretraining works similar interms of better accuracy. But when pretrained using RBM the neural network converges much faster. The reduction in no of iteration is more than 50%, which in auto encoder the decrease iterations because of pretraining is very mild. This way we can say that RBM is learning weights that are more appropriate compared autoencoder.

## 5 Multivariate data set

### Gaussian kernel

#### Parameter estimation:

the best value of  $\sigma$  is observed to be 0.5.the kernel gram matrix is observed to be:(for  $\sigma = 0.5$ )

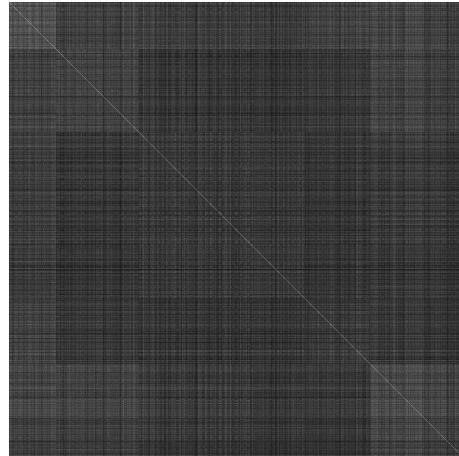


Figure 80: Figure showing grey level image of kernel gram matrix for kernel parameters  $\sigma = 0.5$

We observe from the kernel gram matrix that it is far away from block diagonal matrix which means these classes are not separable it means that features used are not sufficient to capture the class of image.

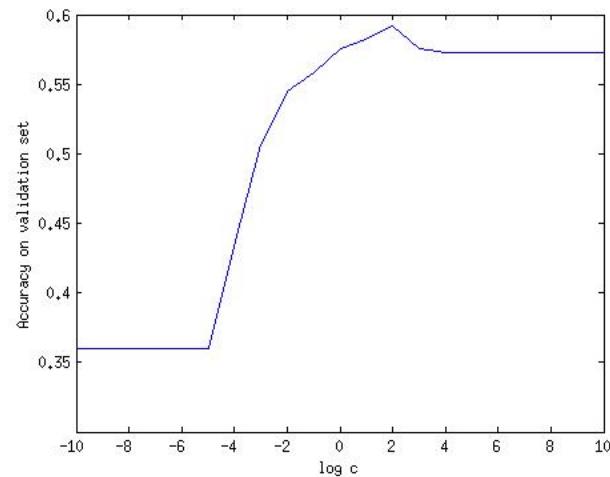


Figure 81: Figure showing accuracy for validation set for different values of  $c$

The best  $c$  is observed to be 4.

#### Confusion matrix:

		Confusion Matrix							
		1	2	3	4	5			
Output Class	1	42 6.8%	0 0.0%	10 1.6%	8 1.3%	2 0.3%	67.7% 32.3%		
	2	4 0.6%	66 10.6%	14 2.3%	10 1.6%	29 4.7%	53.7% 46.3%		
3	11 1.8%	14 2.3%	164 26.4%	29 4.7%	17 2.7%	69.8% 30.2%			
4	4 0.6%	5 0.8%	15 2.4%	40 6.4%	2 0.3%	60.6% 39.4%			
5	3 0.5%	29 4.7%	21 3.4%	6 1.0%	76 12.2%	56.3% 43.7%			
	1	2	3	4	5				
	Target Class								

Figure 82: Figure showing confusion matrix for test set for multivariate dataset gaussian kernel for  $\sigma = 0.5, c=4$

The accuracy for test set is observed to be 62.5%

We see that class1(sheep) is the class which is highest accuracy among the classes.most of the examples of class2 (dining table) are misclassified as class5(sofa) and viceversa.Similarly most of the examples in class3(car) are misclassified as class4(bicycle) and viceversa.This can also be attributed to granularity in classification.(more granular less separable).

## 5.1 Logistic Regression

In logistic regression the posterior probabilities of the class are ensured to be between 0 and 1 and sum over all classes is 1 by using a softmax acting on a linear function. Cross entropy was used as the error function. And the python library sci-kit learn was used for this implementation.

**Parameter estimation:** Logistic Regression was done on the dataset with polynomial basis function. The parameter that was estimated was degree of the basis expansion. We tried varying the degree between 1 and 3. But the number of data points were large and the number of features were 512. As the dimensions after basis expansion are given by  $\binom{d+M}{M} = d+MC_M$  where C is the combination symbol. Hence even a basis expansion of degree 2 means huge increase in the number of features and the program being too large to get results. Hence we tried reducing the number of features by running PCA and then tried to do basis expansion in order to observe non linear boundaries.

The error on the various sets vs the degree was plot. The validation set vs degree was used to select what is the best model. We notice a maximum accuracy on validation set when degree is 2. Hence model with degree 2 is chosen as the best model. But it should be noted that the PCA was done on the sets differently. The model with degree 1 was done with no PCA applied, the model with degree 2 was done with PCA of 100 features and then a basis expansion. And the model with degree 3 was done with a PCA of 10 features followed by a basis expansion. We did this because an increase in the degree even by one was leading to a huge increase in the features.

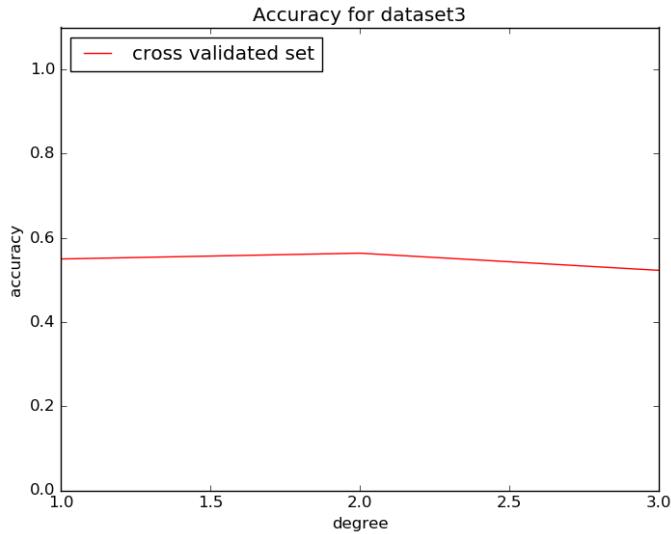


Figure 83: Figure showing accuracy for validation set for different values of degree

#### Confusion matrices:

Here we present the three confusion matrices for the three degrees the results obtained on these matrices were found on the cross validated sets where for each element in the input, the prediction was obtained for that element when it was in the test set.

Confusion Matrix							
Output Class	1	2	3	4	5		
	1	<b>123</b> 3.9%	<b>3</b> 0.1%	<b>32</b> 1.0%	<b>33</b> 1.1%	<b>7</b> 0.2%	<b>62.1%</b> 37.9%
	2	<b>18</b> 0.6%	<b>223</b> 7.2%	<b>58</b> 1.9%	<b>48</b> 1.5%	<b>126</b> 4.0%	<b>47.1%</b> 52.9%
	3	<b>114</b> 3.7%	<b>173</b> 5.6%	<b>912</b> 29.3%	<b>213</b> 6.8%	<b>172</b> 5.5%	<b>57.6%</b> 42.4%
	4	<b>34</b> 1.1%	<b>21</b> 0.7%	<b>49</b> 1.6%	<b>142</b> 4.6%	<b>14</b> 0.4%	<b>54.6%</b> 45.4%
	5	<b>32</b> 1.0%	<b>151</b> 4.8%	<b>71</b> 2.3%	<b>33</b> 1.1%	<b>312</b> 10.0%	<b>52.1%</b> 47.9%
	1	2	3	4	5		
Target Class							

Figure 84: Figure showing confusion matrix when each input was in the test set and degree of basis expansion is 1, No Feature Reduction

Confusion Matrix							
Output Class	1	2	3	4	5		
	1	<b>133</b> 4.3%	<b>4</b> 0.1%	<b>31</b> 1.0%	<b>35</b> 1.1%	<b>9</b> 0.3%	<b>62.7%</b> <b>37.3%</b>
	2	<b>15</b> 0.5%	<b>233</b> 7.5%	<b>59</b> 1.9%	<b>50</b> 1.6%	<b>125</b> 4.0%	<b>48.3%</b> <b>51.7%</b>
	3	<b>105</b> 3.4%	<b>159</b> 5.1%	<b>905</b> 29.1%	<b>203</b> 6.5%	<b>148</b> 4.8%	<b>59.5%</b> <b>40.5%</b>
	4	<b>35</b> 1.1%	<b>17</b> 0.5%	<b>46</b> 1.5%	<b>146</b> 4.7%	<b>11</b> 0.4%	<b>57.3%</b> <b>42.7%</b>
	5	<b>33</b> 1.1%	<b>158</b> 5.1%	<b>81</b> 2.6%	<b>35</b> 1.1%	<b>338</b> 10.9%	<b>52.4%</b> <b>47.6%</b>

Figure 85: Figure showing confusion matrix when each input was in the test set and degree of basis expansion is 2, Feature Reduction with 100 components

Confusion Matrix							
Output Class	1	2	3	4	5		
	1	<b>110</b> 3.5%	<b>2</b> 0.1%	<b>30</b> 1.0%	<b>38</b> 1.2%	<b>9</b> 0.3%	<b>58.2%</b> <b>41.8%</b>
	2	<b>20</b> 0.6%	<b>203</b> 6.5%	<b>65</b> 2.1%	<b>52</b> 1.7%	<b>113</b> 3.6%	<b>44.8%</b> <b>55.2%</b>
	3	<b>103</b> 3.3%	<b>179</b> 5.7%	<b>890</b> 28.6%	<b>232</b> 7.5%	<b>181</b> 5.8%	<b>56.2%</b> <b>43.8%</b>
	4	<b>52</b> 1.7%	<b>14</b> 0.4%	<b>54</b> 1.7%	<b>108</b> 3.5%	<b>11</b> 0.4%	<b>45.2%</b> <b>54.8%</b>
	5	<b>36</b> 1.2%	<b>173</b> 5.6%	<b>83</b> 2.7%	<b>39</b> 1.3%	<b>317</b> 10.2%	<b>48.9%</b> <b>51.1%</b>

Figure 86: Figure showing confusion matrix when each input was in the test set and degree of basis expansion is 3, Feature reduction with 10 components

We find a best accuracy of about 56 percent through logistic regression with basis expansion on this dataset.

## 5.2 MLFFNN using 2 hidden layers

The dataset consists of 512 features, so the hidden layer nodes were varied in the range of high (400- 500), mid range(100 - 400) and low(<100). The best accuracy was obtained for 450 and 400 nodes in hidden layers 1 and 2 respectively.

### Confusion Matrix:

		Confusion matrix for Test data					
		1	2	3	4	5	
Output Class	1	31 5.0%	1 0.2%	9 1.4%	14 2.3%	7 1.1%	50.0% 50.0%
	2	5 0.8%	35 5.6%	10 1.6%	7 1.1%	25 4.0%	42.7% 57.3%
	3	15 2.4%	30 4.8%	171 27.5%	30 4.8%	19 3.1%	64.5% 35.5%
	4	10 1.6%	9 1.4%	15 2.4%	31 5.0%	7 1.1%	43.1% 56.9%
	5	3 0.5%	39 6.3%	19 3.1%	12 1.9%	68 10.9%	48.2% 51.8%
		48.4% 51.6%	30.7% 69.3%	76.3% 23.7%	33.0% 67.0%	54.0% 46.0%	54.0% 46.0%

Figure 87: Confusion matrix for Multivariate dataset for test data

Accuracy obtained in case of neural network is around 54%. The accuracy is low because the features are pre-decided and fed to the neural network. The features might not be enough to distinguish between the different classes.