

## Midterm Assignment

Sarthak Kumar Maharana

1.

- (a) For a 2-class nearest-means classifier, the discriminant function is as follows :

$$\text{For class 1 : } g_1(\underline{x}) = \|\underline{x} - \mu_1\|_2^2$$

$$\text{For class 2 : } g_2(\underline{x}) = \|\underline{x} - \mu_2\|_2^2$$

$$\text{Overall, } g(\underline{x}) = \|\underline{x} - \frac{\mu_1}{2}\|_2^2 - \|\underline{x} - \frac{\mu_2}{2}\|_2^2 = 2(\underline{\mu}_1^T \underline{x} + \underline{\mu}_2^T \underline{x}) - \underline{\mu}_1^T \underline{\mu}_1 + \underline{\mu}_2^T \underline{\mu}_2$$

$$\text{and so, } \begin{cases} \text{class 1} & g(\underline{x}) \geq 0 \\ \text{class 2} & g(\underline{x}) < 0 \end{cases} \quad \leftarrow \text{prediction}$$

- (b) Let the class under consideration be 'k'.

$$\begin{aligned} \tilde{g}_K(\underline{x}) &= \|\underline{x} - \underline{\mu}_k\|_2^2 \\ &= (\underline{x} - \underline{\mu}_k)^T (\underline{x} - \underline{\mu}_k) = \underline{x}^T \underline{x} - \underline{\mu}_k^T \underline{x} - \underline{x}^T \underline{\mu}_k + \underline{\mu}_k^T \underline{\mu}_k \\ &= \underline{x}^T \underline{x} - 2\underline{\mu}_k^T \underline{x} + \underline{\mu}_k^T \underline{\mu}_k \end{aligned}$$

$$\text{Let } g_K(\underline{x}) = 2\underline{\mu}_k^T \underline{x} - \underline{\mu}_k^T \underline{\mu}_k ; \quad \tilde{g}_K(\underline{x}) = \underline{x}^T \underline{x} - g_K(\underline{x})$$

Let's include a quadratic non-linear mapping for two features :

$$[1, x_1, x_2, x_1x_2, x_1^2, x_2^2] . \quad \leftarrow \text{augmented notation}$$

$$g_K(\underline{x}) = \frac{2}{N_K} \sum_{i=1}^{N_K} x_i^{(k)T} \cdot \underline{x} - \left( \frac{1}{N_K} \sum_{i=1}^{N_K} x_i^{(k)T} \right) \left( \frac{1}{N_K} \sum_{i=1}^{N_K} x_i^{(k)} \right)$$

where  $k \rightarrow$  current class

$N_K \rightarrow$  total data of

$$= \frac{2}{N_K} \sum_{i=1}^N x_i^{(k)T} \underline{x} - \frac{1}{N_K^2} \sum_{i=1}^{N_K} \sum_{l=1}^{N_K} x_i^{(k)T} x_l^{(k)} . \quad \rightarrow ①$$

Clearly, we've represented or introduced the inner product notation  $\langle \cdot, \cdot \rangle$ . And so, this is the dual representation for non-linear mapping.

$$\begin{aligned} g(\underline{x}) &= \|\underline{x} - \mu_2\|_2^2 - \|\underline{x} - \mu_1\|_2^2 \\ &= \underline{x}^T \underline{x} - 2\mu_2^T \underline{x} + \mu_2^T \mu_2 - \underline{x}^T \underline{x} + 2\mu_1^T \underline{x} - \mu_1^T \mu_2 \\ &= -2\mu_2^T \underline{x} + \mu_2^T \mu_2 + 2\mu_1^T \underline{x} - \mu_1^T \mu_2 \\ &= \bullet (g_2(\underline{x}) \bullet + g_1(\underline{x})) \end{aligned}$$

or  $g_0(\underline{x}) = \bullet (g_1(\underline{x}) + g_2(\underline{x})) //$

(c) In equation ①, we can substitute the RBF kernel.

$$k(\underline{x}_1, \underline{x}_2) = \exp \{-\gamma \|\underline{x}_1 - \underline{x}_2\|_2^2\}$$

∴ ① becomes:

$$\frac{2}{N_K} \sum_{i=1}^{N_K} \underline{x}_i^{(K)^T} \underline{x} - \frac{1}{N_K^2} \sum_{i=1}^{N_K} \sum_{l=1}^{N_K} \exp \{-\gamma \|\underline{x}_i^{(K)} - \underline{x}_l^{(K)}\|_2^2\}; \quad k \rightarrow \text{current class}$$

$N_K \rightarrow \text{total points of class } K$

$$g_{rbf}(\underline{x}) = g_1(\underline{x}) + g_2(\underline{x})$$

$$\begin{aligned} &= \frac{2}{N_1^2} \sum_{i=1}^{N_1} \underline{x}_i^{(1)^T} \underline{x} - \frac{1}{N_1^2} \sum_{i=1}^{N_1} \sum_{l=1}^{N_1} \exp \{-\gamma \|\underline{x}_i^{(1)} - \underline{x}_l^{(1)}\|_2^2\} \\ &\quad + \frac{2}{N_2^2} \sum_{i=1}^{N_2} \underline{x}_i^{(2)^T} \underline{x} - \frac{1}{N_2^2} \sum_{i=1}^{N_2} \sum_{l=1}^{N_2} \exp \{-\gamma \|\underline{x}_i^{(2)} - \underline{x}_l^{(2)}\|_2^2\} // \end{aligned}$$

The following questions have been answered and solved with respect to a **non-augmented feature space**.

Objective: Code up a 2-class kernel nearest means classifier that can use Radial Basis Function (RBF) or linear kernel.

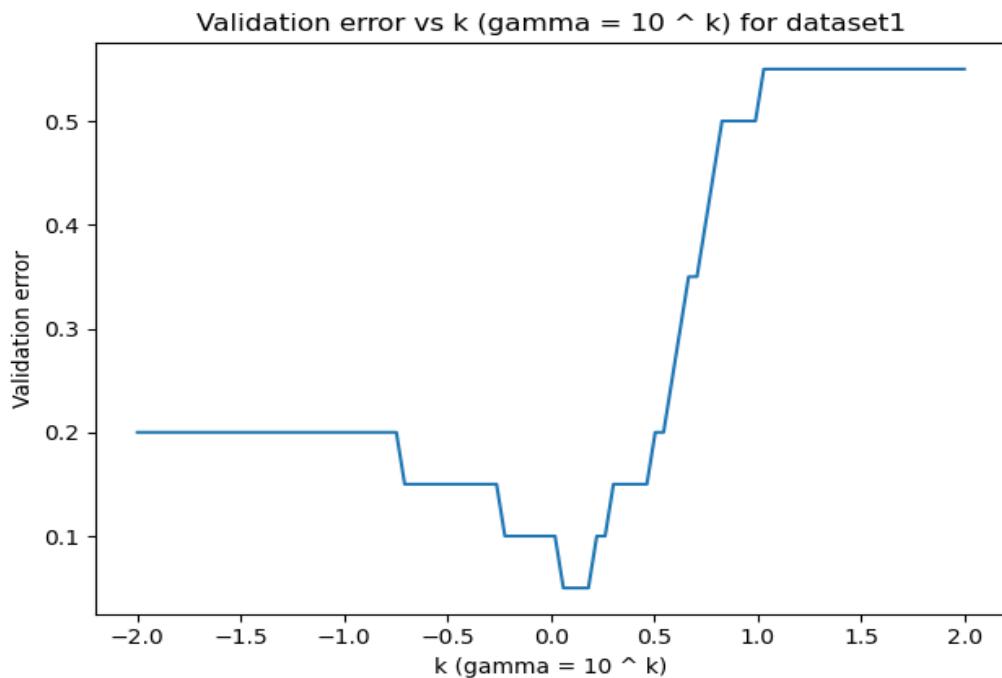
(d) The validation set, as per model selection, is used to choose the best performing model. In the current problem, the validation set is used to obtain the optimal value of gamma for the RBF kernel. Here, based on the validation set, the value of gamma that gives the lowest classification error, is chosen as the optimal gamma for the rest of the problems that follow.

For **dataset 1**: optimal gamma --- **1.1497**

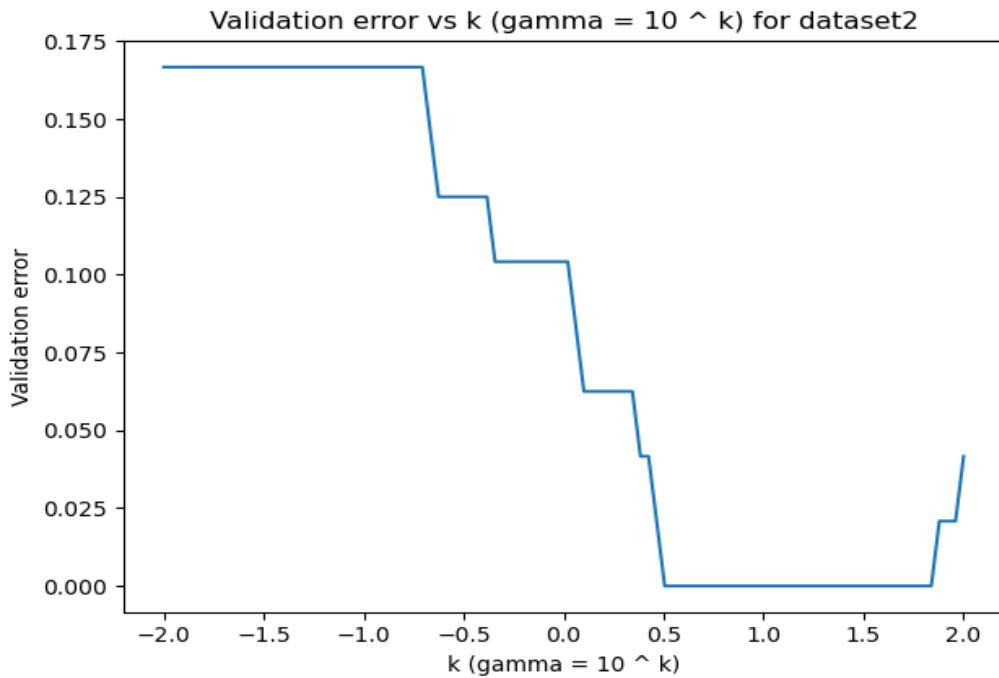
For **dataset 2**: optimal gamma --- **3.1992**

(e) To obtain the optimal gamma, a brute-force search was done with the values of k ranging from **-2 to 2**, for **100 points in between**. The following figures illustrate the validation-set classification error vs k, where  $\gamma = 10^k$ .

For dataset 1:



For dataset 2:



**(f) For dataset 1:**

The test set error using the linear kernel is : **0.43 (43%)**

The test set error using the RBF kernel is : **0.23 (23%)**

For dataset 2:

The test set error using the linear kernel is : **0.243 (24.3%)**

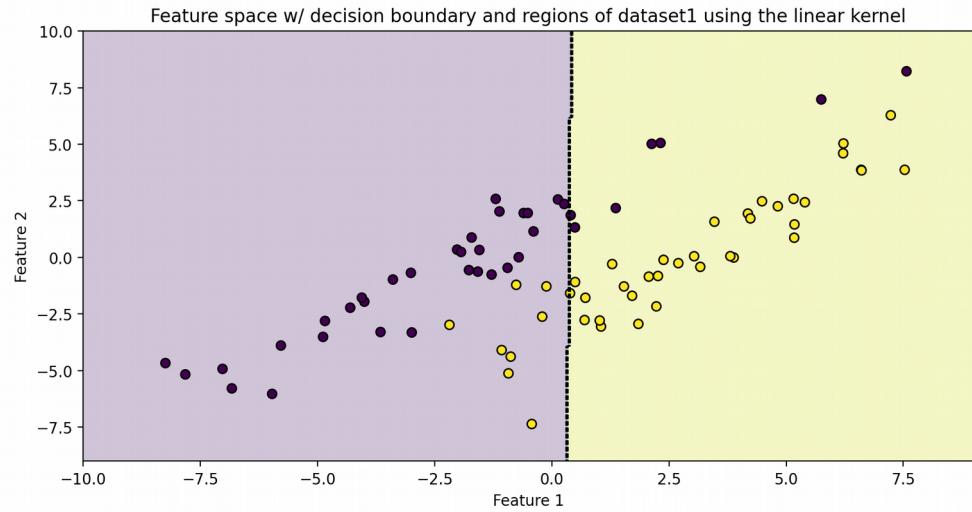
The test set error using the RBF kernel is : **0.0125 (1.25%)**

An RBF kernel's main objective is used to find a non-linear classifier or regression decision boundary. The gamma used in the kernel is used to control the degree to which two data points are to be similar to each other. This increases the chances that the decision boundary fits well and helps in allowing the dataset to be linearly separable, in a higher dimensional feature space. A linear class doesn't do well in a higher-dimensional space, and utterly fails in allowing the dataset to be linearly separable.

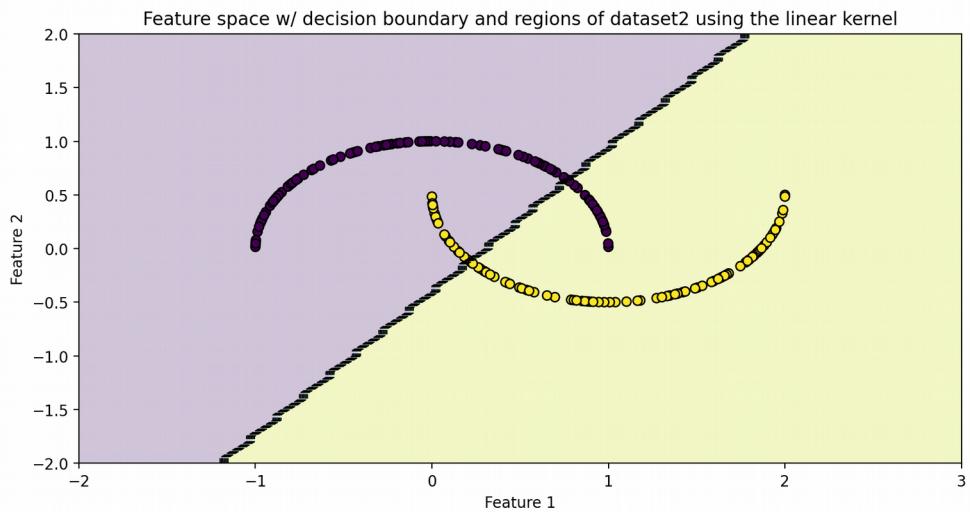
(g) The following figures illustrate the use of the **linear kernel**, for classification, on each dataset. Here, we plot the training data points, decision regions and boundary, in the feature space.

(The linear boundaries look a little rugged because of being stretched.)

For dataset 1:

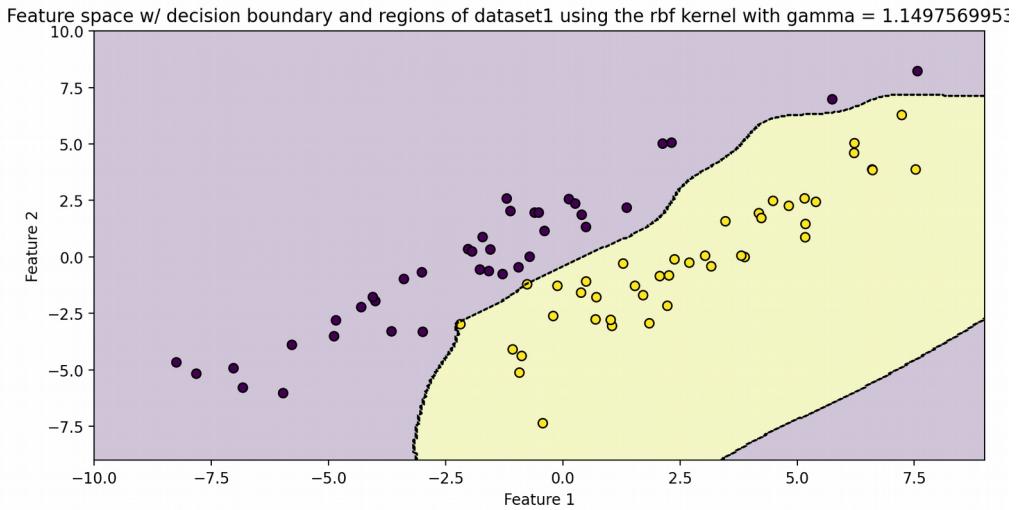


For dataset 2:

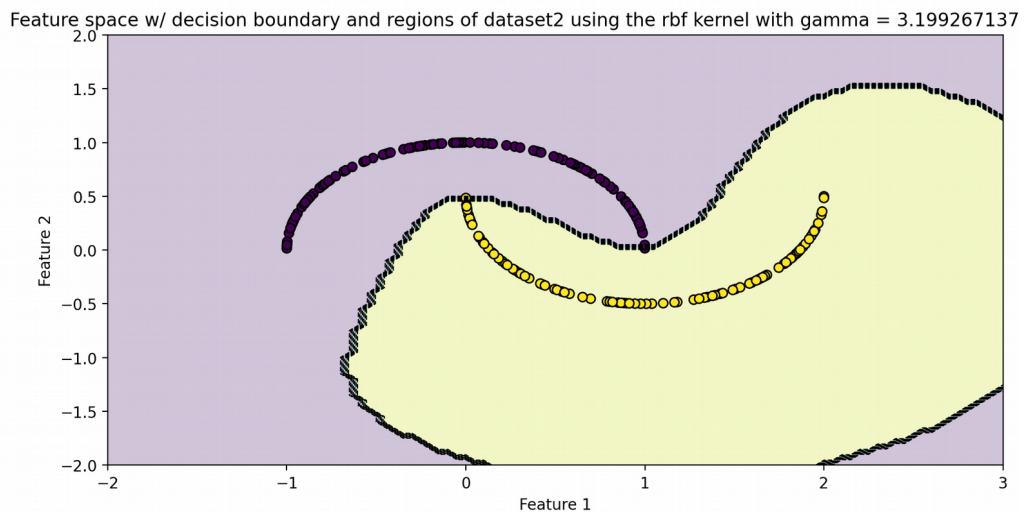


**(h)** The following figures illustrate the use of the **RBF kernel (w/ optimal gamma)**, for classification, on each dataset.

For dataset 1:



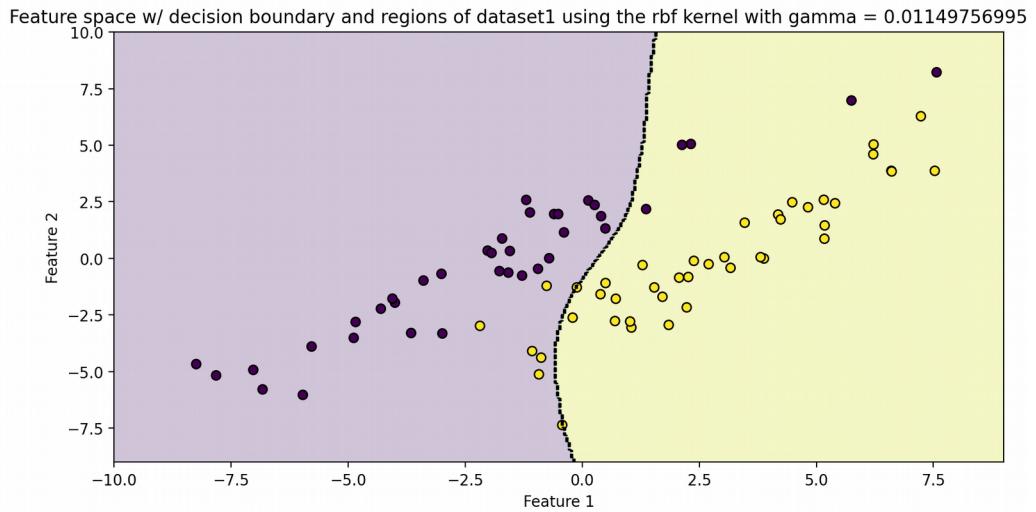
For dataset 2:



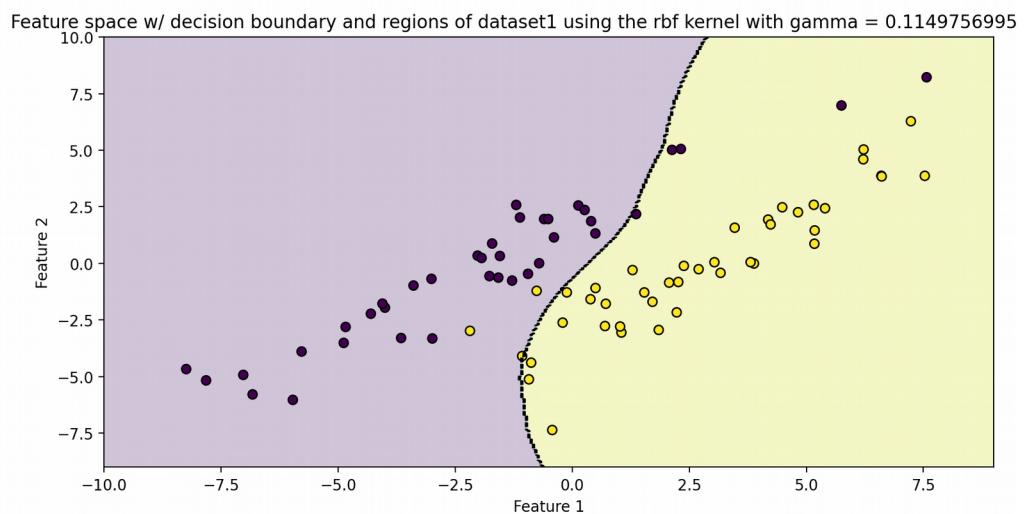
**(i)** In this subproblem, we repeat (h) but with different values of the multiples of optimal gamma i.e the multiples include 0.01, 0.1, 0.3, 3, 10, 100 for each dataset.

For dataset 1:

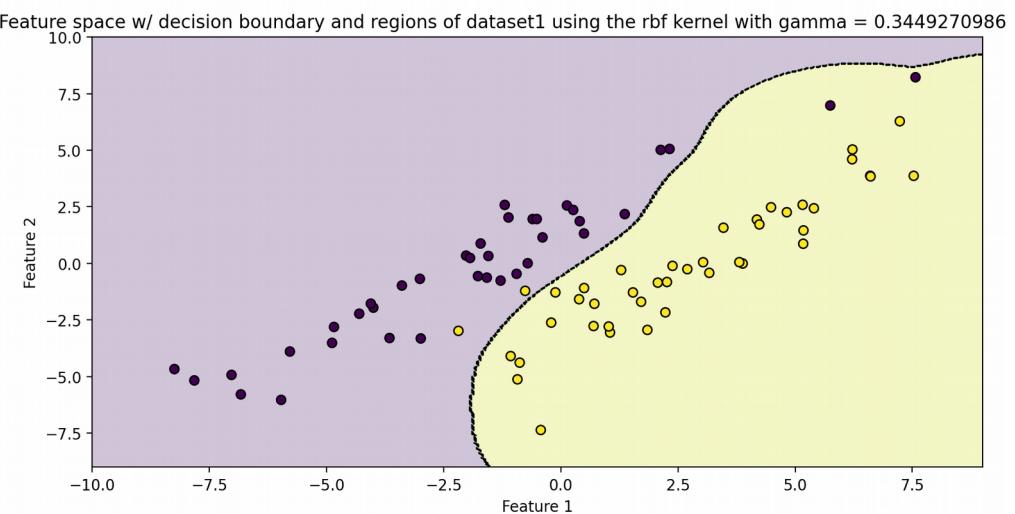
$$\text{gamma} = 0.01 * 1.1497 = 0.011497$$



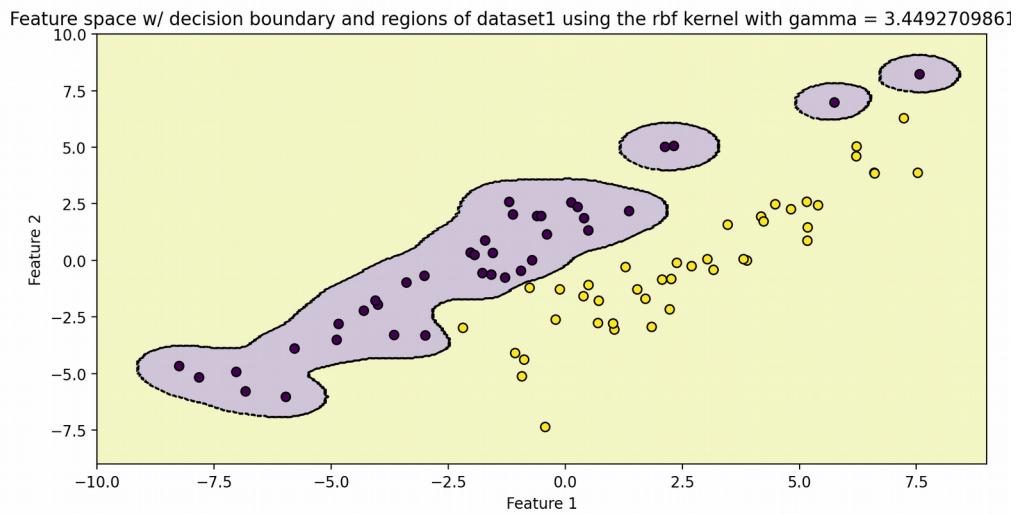
$$\text{gamma} = 0.1 * 1.14967 = 0.11497$$



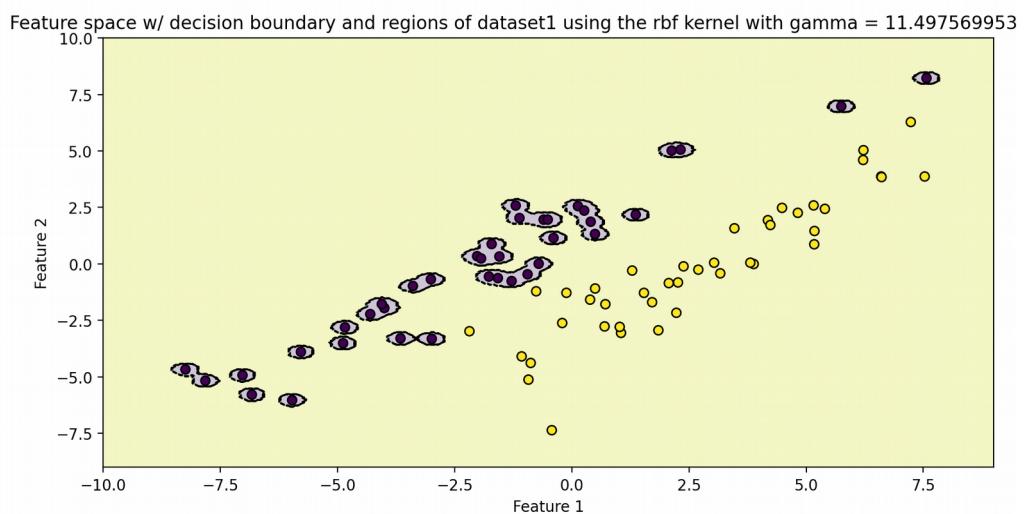
$$\text{gamma} = 0.3 * 1.14967 = 0.3449$$



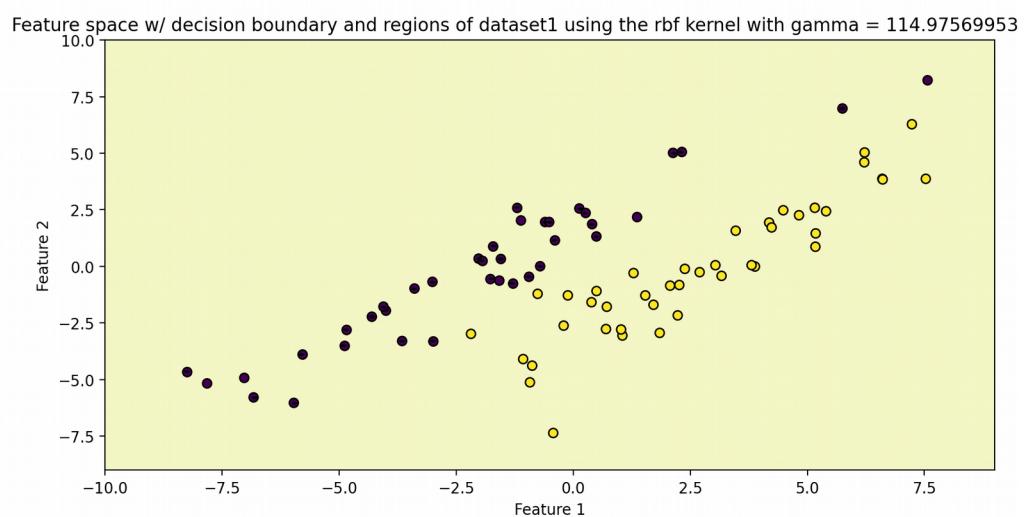
**gamma = 3 \* 1.14967 = 3.4492**



**gamma = 10 \* 1.14967 = 11.4967**

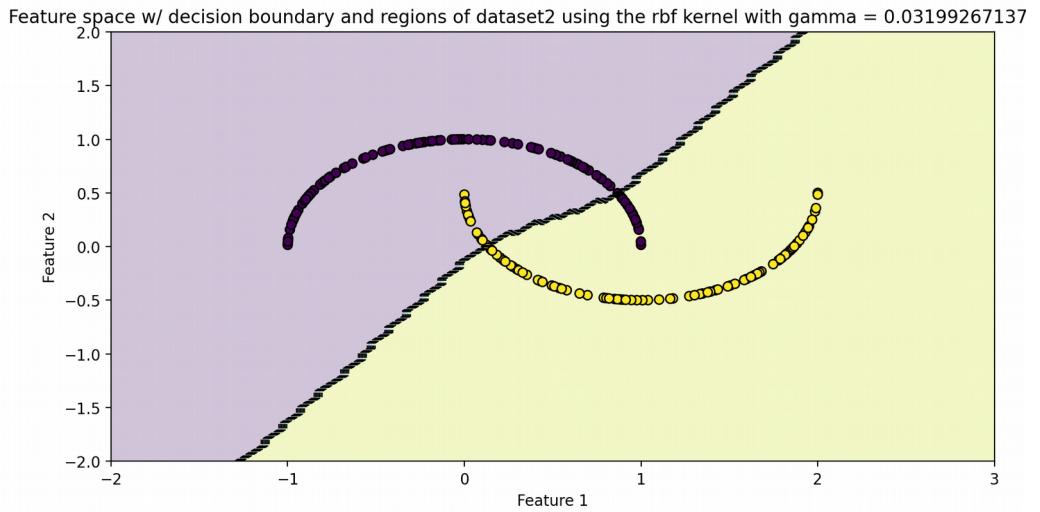


**gamma = 100 \* 1.14967 = 114.967**

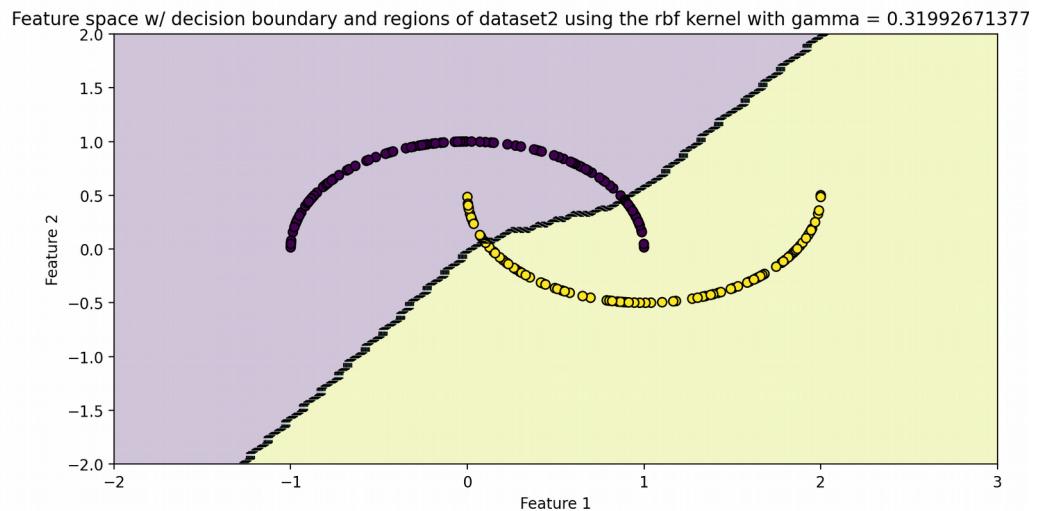


For dataset 2:

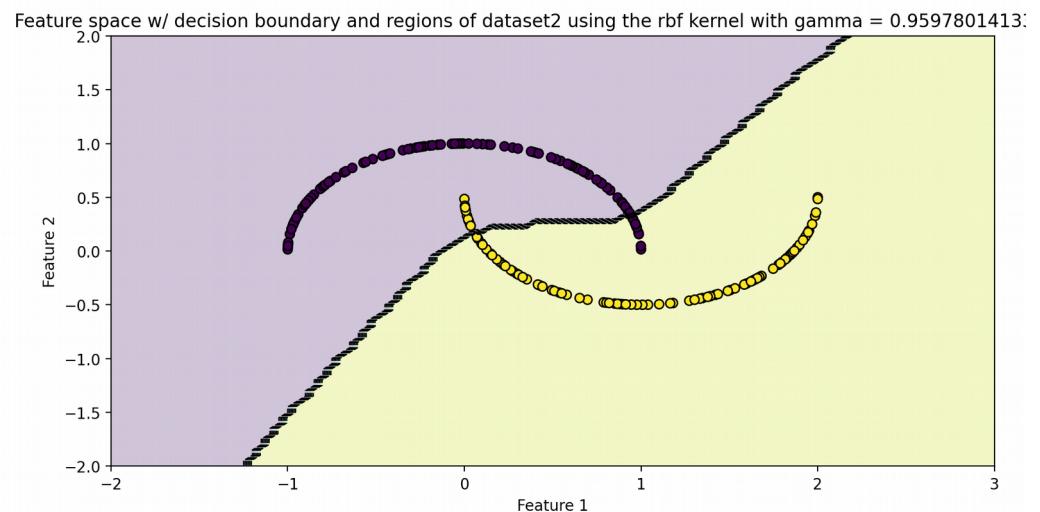
$$\text{gamma} = 0.01 * 3.1992 = 0.031992$$



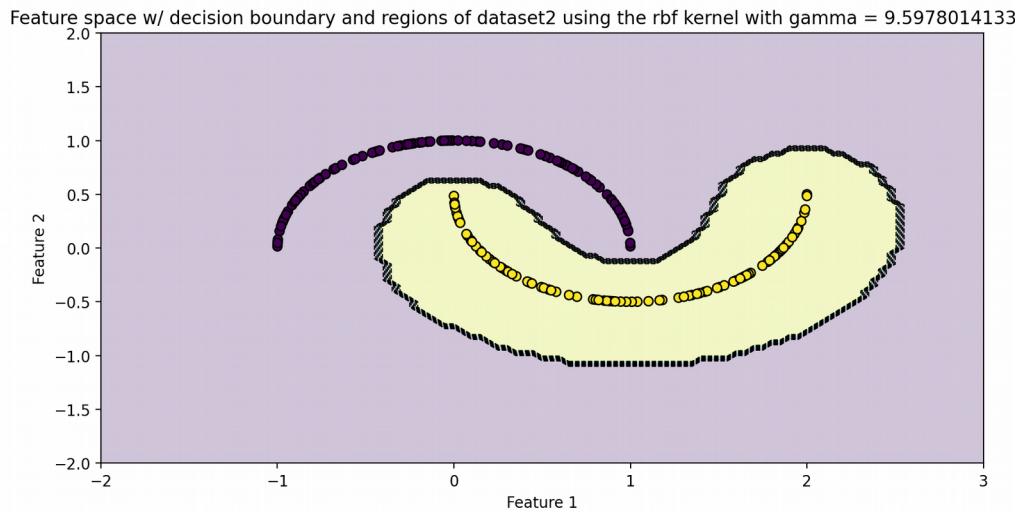
$$\text{gamma} = 0.1 * 3.1992 = 0.31992$$



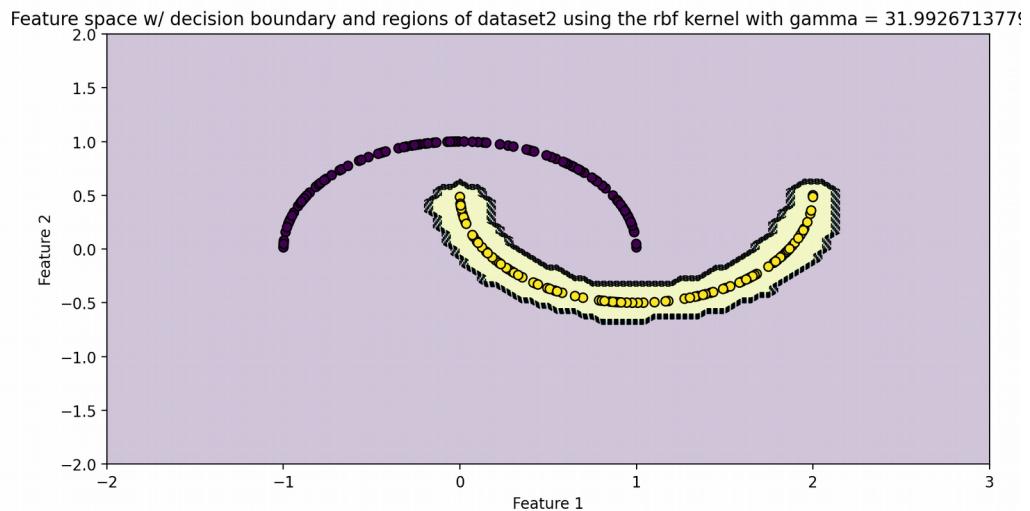
$$\text{gamma} = 0.3 * 3.1992 = 0.9597$$



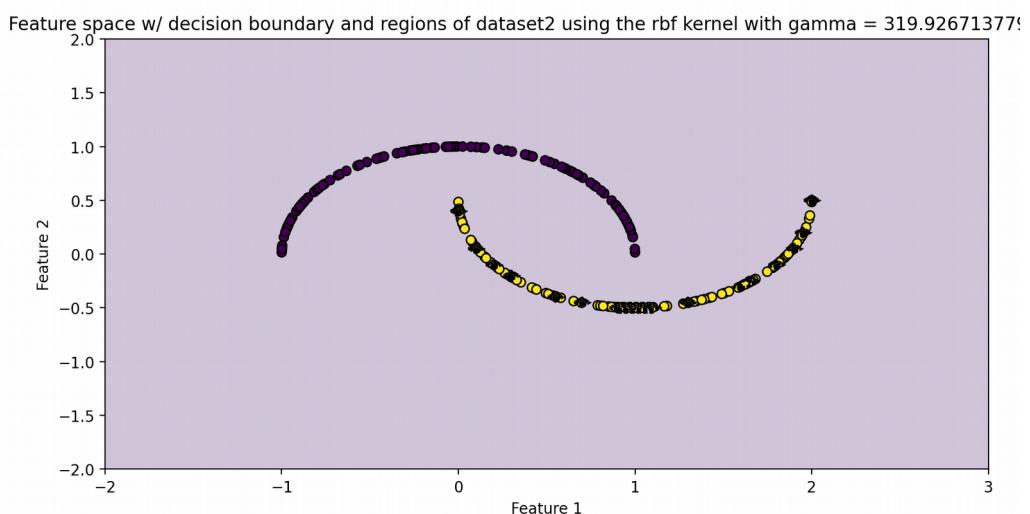
**gamma = 3 \* 3.1992 = 9.597**



**gamma = 10 \* 3.1992 = 31.992**



**gamma = 100 \* 3.1992 = 319.92**



We observe that as the value of *gamma* is increased beyond the optimal gamma (for each dataset), the performance of non-linear classification begins to get better. As said earlier, gamma parameterizes and models the similarity between two data points. And so, a higher value *tries* to bring a higher similarity between the training and unknown (validation/test) data points. In a standard RBF kernel, gamma is defined as the inverse of the variance of the Gaussian/normal distribution (that the RBF's probability distribution underlies). And so, larger values of gamma ensure that the variance between the two data points is as low as possible, resulting in high similarity and better model performance on unknowns. At higher values of the optimal gamma, the decision boundary performs extremely well in classifying between the two classes of data, and so the error rate reduces massively (expected).

2. Given criterion function of a 2-class problem (using augmented notation) :

$$J(\underline{w}) = - \sum_{n=1}^N \{ [g(z_n \underline{x}_n) \leq 0] \} \frac{\underline{w}^T z_n \underline{x}_n}{\|\underline{w}\|_2}$$

(a) Based on the indicator function,  $J(\underline{w})$  can be rewritten as :

$$J(\underline{w}) = \begin{cases} - \sum_{n=1}^N \frac{\underline{w}^T z_n \underline{x}_n}{\|\underline{w}\|_2} & g(z_n \underline{x}_n) \leq 0 \rightarrow \text{misclassified} \\ 0 & \text{else} \rightarrow \text{correct classification} \end{cases}$$

Now,  $\frac{\underline{w}^T z_n \underline{x}_n}{\|\underline{w}\|_2}$  can be interpreted as the distance of the point  $\underline{x}_n$  from the hyperplane  $H$ .

or  $J(\underline{w}) = - (\text{Sum of all distances of the data points from the hyperplane})$  ; given the current weight vector.

And so, when there are correct classifications,  $J(\underline{w}) = 0$ .

If the misclassified points are close to each other, the magnitude of  $J(\underline{w})$  can be brought down.

$$\begin{aligned} J(\underline{w}) &= - \sum_{n=1}^N \{ [g(z_n \underline{x}_n) \leq 0] \} \frac{\underline{w}^T z_n \underline{x}_n}{\|\underline{w}\|_2} \\ &= \begin{cases} - \sum_{n=1}^N \frac{\underline{w}^T z_n \underline{x}_n}{\|\underline{w}\|_2} & \text{if } g(z_n \underline{x}_n) \leq 0 \\ 0 & g(z_n \underline{x}_n) > 0 \end{cases} \end{aligned}$$

Taking the gradient of  $J(\underline{w})$  wrt  $\underline{w}$ :

$$\nabla_{\underline{w}} J(\underline{w}) = - \sum_{n=1}^N \frac{\|\underline{w}\|_2 z_n \underline{x}_n - (\underline{w}^T z_n \underline{x}_n) \frac{1}{2} (\underline{w}^T \underline{w}) 2 \underline{w}}{\|\underline{w}\|_2^2}$$

$$= - \sum_{n=1}^N \frac{\|\underline{w}\|_2 z_n \underline{x}_n - \underline{w}^T z_n \underline{x}_n \cdot \frac{\underline{w}}{\|\underline{w}\|_2}}{\|\underline{w}\|_2^2}$$

$$= - \sum_{n=1}^N \frac{z_n \underline{x}_n}{\|\underline{w}\|_2} - \frac{\underline{w}^T z_n \underline{x}_n \underline{w}}{\|\underline{w}\|_2^3}$$

$$\text{or } \nabla_{\underline{w}} J(\underline{w}) = \begin{cases} - \sum_{n=1}^N \frac{z_n \underline{x}_n}{\|\underline{w}\|_2} - \frac{\underline{w}^T z_n \underline{x}_n \underline{w}}{\|\underline{w}\|_2^3} & \text{if } g(z_n \underline{x}_n) \leq 0 \\ 0 & \text{else.} \end{cases}$$

### Batch gradient descent

The weight update rule is as follows:

$$\underline{w}(i+1) = \underline{w}(i) - \eta(i) \nabla_{\underline{w}} J(\underline{w}) \quad \eta(i) \rightarrow \text{learning rate}$$

$$= \underline{w}(i) + \eta(i) \sum_{n=1}^N \left( \frac{z_n \underline{x}_n}{\|\underline{w}\|_2} - \frac{\underline{w}^T z_n \underline{x}_n \underline{w}}{\|\underline{w}\|_2^3} \right)$$

$$\text{Stop when } \sum_{n=1}^N \frac{z_n \underline{x}_n}{\|\underline{w}\|_2} - \frac{\underline{w}^T z_n \underline{x}_n \underline{w}}{\|\underline{w}\|_2^3} = 0 \checkmark$$

$$\underline{w}(i+1) = \underline{w}(i) + \eta(i) \sum_{n=1}^N [g(z_n \underline{x}_n \leq 0)] \left( \frac{z_n \underline{x}_n}{\|\underline{w}\|_2} - \frac{\underline{w}^T z_n \underline{x}_n \underline{w}}{\|\underline{w}\|_2^3} \right) \checkmark$$

## Stochastic gradient descent - variant 1 :

(Randomly shuffle before each epoch)

Iterate over  $m = 1, 2, 3, \dots$  until no change in  $w$  over 1 full epoch.

i) Randomly shuffle dataset

ii), ~~or~~ Define the sequential epoch variable  $m$  as:

for  $m$  in range( $\text{len}(N)$ ) :  $N \rightarrow \text{data points}$

$$i = (m-1)N + (n+1)$$

$$\begin{cases} \underline{w}(i+1) = \underline{w}(i) + \eta(i) \nabla_{\underline{w}} J(\underline{w}) & \text{if } \underline{w}^T \underline{x}_n \underline{x}_n \leq 0 \\ \underline{w}(i+1) = \underline{w}(i) & \text{if } \underline{w}^T \underline{x}_n \underline{x}_n > 0 \end{cases}$$

$\nabla_{\underline{w}} J(\underline{w}) \rightarrow$  gradient of  $J(w)$  defined earlier

$$= - \sum_{n=1}^N \frac{\underline{x}_n \underline{x}_n}{\|\underline{w}\|_2} - \frac{\underline{w}^T \underline{x}_n \underline{x}_n \underline{w}}{\|\underline{w}\|_2^3} \quad \text{if } g(\underline{x}_n \underline{x}_n) \leq 0$$

(i) Total features =  $D = 2$

$$\eta(i) = \frac{1}{1+i}$$

For epoch 1 :  $(\underline{x}_1, \underline{x}_2, \underline{x}_3, \underline{x}_4)$

(i)  $S_1 : \underline{x}_1 = [-4, 1]^T, \underline{x}_2 = [-3, -2]^T$

$S_2 : \underline{x}_3 = [-1, 3]^T, \underline{x}_4 = [3, -2]^T$

Since we're dealing w/ a non-augmented feature space

iii),  $\underline{w}(1) = [-1, 4]^T$

$$w_1 = -1; w_2 = 4. \quad (\text{since non-aug})$$

iv), As found above,

$$\underline{w}(i+1) = \underline{w}(i) + \eta(i) \nabla_{\underline{w}} J(\underline{w})$$

$$\text{or } \underline{w}(2) = \underline{w}(1) + \eta(1) \left[ \frac{\underline{z}_2 \underline{z}_2^T}{\cancel{\|\underline{w}(1)\|_2}} - \frac{\underline{w}(1)^T \underline{z}_2 \underline{z}_2 \underline{w}(1)}{\|\underline{w}(1)\|_2^3} \right]$$

$$\eta(1) = \frac{1}{2}$$

$$\underline{w}(2) = [-1, 4]^T + \frac{1}{2} \left[ \frac{\begin{bmatrix} -3 \\ -2 \end{bmatrix}}{\|[-1, 4]^T\|_2} - \frac{[-1, 4]^T [-3, -2] [-1, 4]^T}{\|[-1, 4]\|_2^3} \right]$$

$$= [-1, 4]^T + \frac{1}{2} \left[ \frac{\begin{bmatrix} 3 \\ -2 \end{bmatrix}}{\sqrt{1^2 + 4^2}} - \frac{\begin{bmatrix} 5 \\ -20 \end{bmatrix}}{(\sqrt{1^2 + 4^2})^{3/2}} \right]$$

$$= [-1, 4]^T + \frac{1}{2} \left[ \frac{\begin{bmatrix} 3 \\ -2 \end{bmatrix}}{\sqrt{17}} - \frac{\begin{bmatrix} 5 \\ -20 \end{bmatrix}}{(\sqrt{17})^3} \right]$$

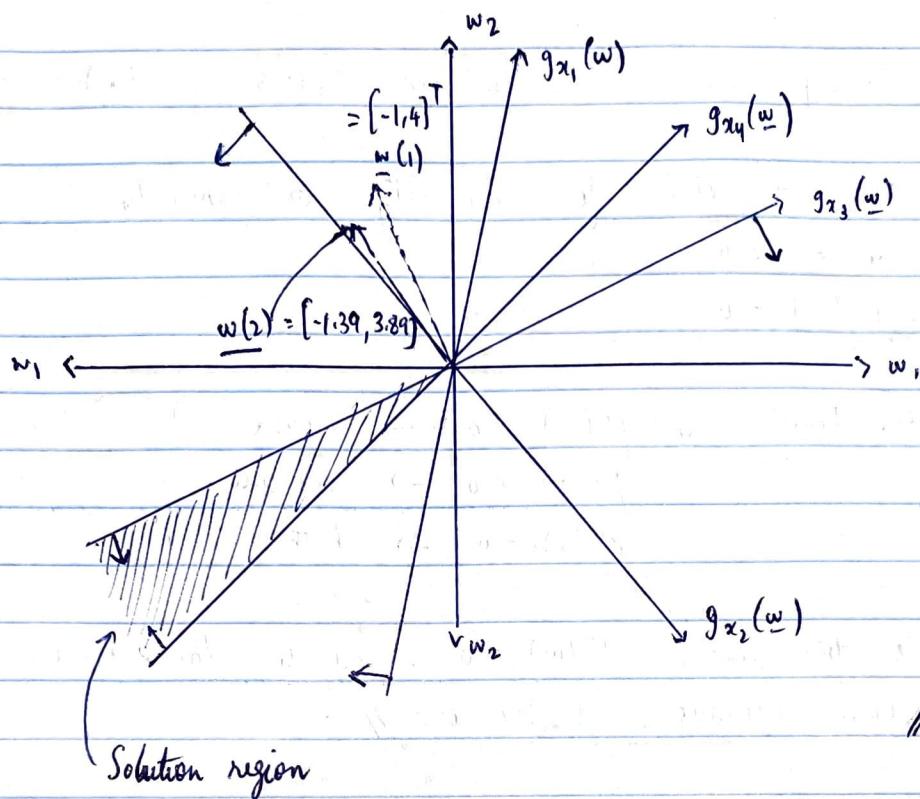
$$= \begin{bmatrix} -1 \\ 4 \end{bmatrix} + \frac{1}{2} \left[ \begin{bmatrix} -0.1727 \\ -0.485 \end{bmatrix} - \begin{bmatrix} -0.071 \\ -0.98 \end{bmatrix} \right]$$

$$\text{or } \underline{w}(2) = \begin{bmatrix} -1.399 \\ 3.897 \end{bmatrix}$$

$$\text{Or, at iteration 2, } \left. \begin{array}{l} w_1 = -1.399 \\ w_2 = 3.897 \end{array} \right\} \in \underline{w}(2)$$

## Plot for 2

(v)



(d) We're given,  $w = \underline{w}$ ,  $\underline{x}_n = x$  and  $w_0 = 0$  } scalars .

$$J(w) = \begin{cases} - \sum_{n=1}^N [g(z_n x_n) \leq 0] \frac{w z_n x_n}{\|w\|_2} & g(z_n x_n) \leq 0 \\ 0 & \text{else} \end{cases}$$

The L2 norm of a scalar is its absolute value .

$$= \begin{cases} - \sum_{n=1}^N \mathbb{I}_{w z_n x_n \leq 0} \frac{w z_n x_n}{\|w\|} & w z_n x_n \leq 0 \\ 0 & \text{else} \end{cases}$$

$$= \begin{cases} \sum_{n=1}^N z_n x_n - \sum_{n=1}^N z_n x_n & w \leq 0 \\ 0 & w > 0 \end{cases} \quad \left. \begin{array}{l} w \leq 0 \\ w > 0 \end{array} \right\} w z_n x_n \leq 0 .$$

8-3-19

Clearly, since  $z_n$  &  $x_n$  are scalars,  $J(w) = \text{constant} = \lambda$

Now, to check for convexity or concavity,

$$J'(w) = 0$$

$$J''(w) = 0$$

We know, if  $f''(x) > 0 \rightarrow \text{convex}$   
 $f''(x) \leq 0 \rightarrow \text{concave}$   
 $f''(x) = 0 \rightarrow \text{both}$

In this case,  $J''(w) = 0$  and so the criterion function is both concave and convex. //

This clearly states that the weight update will happen for an infinite amount of time and won't reach convergence.

3. Total no. of training data points =  $N$ .

No. of features =  $D = 3$ .

(a) In approach A1, we do a degree 3 or cubic polynomial mapping from  $\underline{x}$  to  $\underline{u}$  (expanded feature space) and then use a linear regressor in  $\underline{u}$  space. Let the features be  $(x_1, x_2, x_3)$ .

$\phi(\underline{x}) = \underline{u} =$	$1$	$w' =$	$w_0$
$x_1$			$w_{01}$
$x_2$			$w_{02}$
$x_3$			$w_{03}$
$x_1^2$			$w_{01}^2$
$x_2^2$			$w_{02}^2$
$x_3^2$			$w_{03}^2$
$x_1 x_2$			$w_{01} w_{02}$
$x_1 x_3$			$w_{01} w_{03}$
$x_2 x_3$			$w_{02} w_{03}$
$x_1^3$			$w_{01} w_{01} w_{01}$
$x_2^3$			$w_{02}^3$
$x_3^3$			$w_{03}^3$
$x_1 x_2^2$			$w_{01} w_{02}^2$
$x_1^2 x_2$			$w_{01}^2 w_{02}$
$x_1^2 x_3$			$w_{01}^2 w_{03}$
$x_2^2 x_3$			$w_{02}^2 w_{03}$
$x_1 x_2 x_3$			$w_{01} w_{02} w_{03}$

$y_N \rightarrow$  target values

(augmented notation)

$20 \times 1$        $20 \times 1$

The expanded feature space & the weight vector have a total dimension of 20. If the labels are denoted by  $y_N$  in the new space,  $J_{A1}(w') = \frac{1}{N} \|\phi(\underline{x}) w' - y_N\|_2^2$

$$= \frac{1}{N} \|\underline{u} w' - y_N\|_2^2 //$$

(b) degrees of freedom is defined as the total number of parameters that have to be fit to the data. Clearly, the total parameters = 20.

$$\text{And so, } \text{dof} = 20$$

(c) In approach 1,  $J(\underline{w}') = \frac{1}{N} \|\underline{w}' - \underline{y}_N\|_2^2$ .

To obtain the optimal weight vector, we use the Moore-Penrose method and so,

$$\hat{\underline{w}}' = \underline{U}' \underline{y}_N \quad \text{where } \underline{U}' = (\underline{U}^T \underline{U})^{-1} \underline{U}^T$$

$\underline{U} \rightarrow$  New feature space of dimension = 20.

$\underline{y}_N \rightarrow$  Labels in the new space.

(d) As per approach 2, in the original feature space,  $\hat{f}(\underline{x}) = \underline{w}^T \underline{x}$ .

In the non-linear space,  $\hat{y}(\underline{x}) = w_1' (\hat{f}(\underline{x})) + w_2' (\hat{f}(\underline{x}))^2 + w_3' (\hat{f}(\underline{x}))^3$ .

Let  $\underline{w}' = [w_1', w_2', w_3']$ .

$$\therefore J(\underline{w}') = \frac{1}{N} \|\hat{y}(\underline{x}) - \underline{y}_N\|_2^2 \quad N \rightarrow \text{total data points}$$

$$= \frac{1}{N} \|w_1' \hat{f}(\underline{x}) + w_2' \hat{f}(\underline{x})^2 + w_3' \hat{f}(\underline{x})^3 - \underline{y}_N\|_2^2$$

$$J(\underline{w}, \underline{w}') = \frac{1}{N} \left\| \begin{bmatrix} \hat{f}(\underline{x}) & \hat{f}(\underline{x})^2 & \hat{f}(\underline{x})^3 \end{bmatrix} \underline{w}' - \underline{y}_N \right\|_2^2$$

$$= \frac{1}{N} \left\| \begin{bmatrix} \underline{w}^T \underline{x} & (\underline{w}^T \underline{x})^2 & (\underline{w}^T \underline{x})^3 \end{bmatrix} \underline{w}' - \underline{y}_N \right\|_2^2$$

$\underline{y}_N \rightarrow$  training labels  
or target values

(e) Since, it's a 3 feature training matrix

$\underline{w} \rightarrow 4$  dimensional (since augmented)

In the expanded feature space,  $\underline{w}' = [w_1' \ w_2' \ w_3'] = 3$ .

Total dof  $\rightarrow 4 + 3/2 = 7$

$$(f) J(\underline{w}') = \frac{1}{N} \|\hat{y}(x) - y_N\|_2^2$$

$$\text{but } \hat{y}(x) = w_1' (f(x)) + w_2' (f(x))^2 + w_3' (f(x))^3.$$

$$f(x) = \underline{w}' \underline{x}$$

Clearly,  $\hat{y}(x)$  is non-linear wrt to  $x$  (since it's a polynomial)

$\hat{y}(x)$  is non-linear wrt to  $w$  and linear wrt  $w'$ .

$\hat{y}(x)$  is directly parameterised as a linear function of  $w' = [w_1', w_2', w_3']$ .

$$J(w, w')$$

(g) Criterion function of A2  $\rightarrow \underline{J(w)} = \frac{1}{N} \|\hat{y}(x) - y_N\|_2^2$

$$= \frac{1}{N} \| [f(x) \ f(x)^2 \ f(x)^3]^T w' - y_N \|_2^2.$$

Every convex function is continuous and differentiable. Clearly,  $\underline{J(w')}$  is convex since it satisfies the Hessian double differential check. The Hessian is positive semi-definite everywhere.

$$J(w, w')$$

So,  $\underline{J(w)}$  is continuous and differentiable in terms of  $w'$  and  $w$ .

$$(h) \quad \text{Given, } f(\underline{x}) = \underline{w}^T \underline{x} .$$

$$J(\underline{w}) = \frac{1}{N} \| \underline{w}^T \underline{x} - y_N \|_2^2 .$$

$$J(\underline{w}, \underline{w}') = \frac{1}{N} \| [f(\underline{x}), f(\underline{x})^2, f(\underline{x})^3]^T \underline{w}' - y_N \|_2^2 .$$

Taking the gradient with respect to  $\underline{w}'$ ,

$$\nabla_{\underline{w}'} J(\underline{w}, \underline{w}') = \frac{2}{N} \left( [f(\underline{x}), f(\underline{x})^2, f(\underline{x})^3]^T \underline{w}' - y_N \right) \cdot [f(\underline{x}), f(\underline{x})^2, f(\underline{x})^3]^T$$

The weight update formula is as follows; at iteration  $i$ :

$$\underline{w}'^{(i+1)} = \underline{w}'^{(i)} - \rho(i) \nabla_{\underline{w}'} J(\underline{w}, \underline{w}')$$

$$= \underline{w}'^{(i)} - \rho(i) \cdot \frac{2}{N} \left( [f(\underline{x}), f(\underline{x})^2, f(\underline{x})^3]^T \underline{w}' - y_N \right) \cdot [f(\underline{x}), f(\underline{x})^2, f(\underline{x})^3]^T$$

where  $\rho(i) \rightarrow$  learning rate update of  $\underline{w}'$ .

Taking the gradient wrt  $\underline{w}$ ,

~~$$J(\underline{w}, \underline{w}') = \frac{1}{N} \| w_1' \underline{w}^T \underline{x} + w_2' (\underline{w}^T \underline{x})^2 + w_3' (\underline{w}^T \underline{x})^3 - y_N \|_2^2 .$$~~

$$\nabla_{\underline{w}} J(\underline{w}, \underline{w}') = \frac{2}{N} \left( w_1' \underline{w}^T \underline{x} + w_2' (\underline{w}^T \underline{x})^2 + w_3' (\underline{w}^T \underline{x})^3 - y_N \right) \times \left( \cancel{\underline{w}^T \underline{x}} + \frac{2w_2' (\underline{w}^T \underline{x}) \underline{x}}{w_1' \underline{w}^T \underline{x}} + 3w_3' (\underline{w}^T \underline{x})^2 \underline{x} \right)$$

$$= \frac{2}{N} \left( w_1' \hat{f}(\underline{x}) + w_2' (\hat{f}(\underline{x}))^2 + w_3' (\hat{f}(\underline{x}))^3 - y_N \right) \left( \cancel{\underline{w}^T \underline{x}} + 2w_2' \hat{f}(\underline{x}) \underline{x} + 3w_3' \hat{f}(\underline{x})^2 \underline{x} \right) \|$$

The weight update formula is as follows at iteration  $i$ :

$$\underline{w}(i+1) = \underline{w}(i) - \eta(i) \nabla_{\underline{w}} J(\underline{w})$$

$$= \underline{w}(i) - \eta(i) \frac{1}{N} \left( w_1' f(x) + w_2' (f(x))^2 + w_3' (f(x))^3 - y_N \right) \left( w_1' x + 2w_2' f(x)x + 3w_3' f(x)^2 x \right) //$$

(i) For approach A1,  $J(\underline{w}) = \frac{1}{N} \| \underline{w}' - y_N \|_2^2$

$$= \| \phi(x) \underline{w}' - y_N \|_2^2$$

With L2 regularization or by ridge regression, regularization:

$$J(\underline{w}) = \frac{1}{N} \| \underline{w}' - y_N \|_2^2 + \lambda \| \underline{w} \|_2^2 \quad \text{where } \lambda \text{ is the regularization parameter.}$$

$\lambda$  is a penalty factor that controls the range of values of the weights. In this way, it prevents overfitting on the training set.

(j) For approach A2,  $J(\underline{w}, \underline{w}') = \frac{1}{N} \| \hat{f}(x) - y_N \|_2^2$

With L2 regularization or ridge regularization:

$$J(\underline{w}, \underline{w}') = \frac{1}{N} \| \hat{f}(x) - y_N \|_2^2 + \lambda \| \underline{w} \|_2^2 + \lambda' \| \underline{w}' \|_2^2$$

$\rightarrow \lambda$  &  $\lambda'$  are the regularization parameters, to prevent overfitting.

(k) When there's plenty of data i.e.  $N$  is more, approach A1 will be more suitable. The degrees of freedom in this case is roughly  $D+1 = 4$ . And when  $N > 100$ , we obtain an over-determined case and so, the error will be less on the unseen data. We will obtain a better fitted polynomial for the entire data.

(d)

When the training data is limited, A<sub>2</sub> is preferable. Here, the amount of data i.e N is less. The degrees of freedom of such a model is 7. Clearly, during training of the regression model, the no. of parameters to parameterize and tune are more and as gradient descent happens, the regressor begins to fit well to the data. This will result in the model to perform fairly well on the unseen data set. So, A<sub>2</sub> is preferable even though it's an underdetermined case.