

Perceptron's and SVM

**Rithvik Reddy
Ananth**

24th October 2020

Machine Learning

Cowan Charles

Perceptron's

As usual, we want to generate a data set to fit a perceptron onto. Recall that in the analysis of perceptrons in the notes, we assumed that all the data was contained within the unit sphere of whatever the underlying dimension was. Consider generating data points *on* the unit sphere in the following way: let \underline{Z} be a k -dimensional vector, where each entry is an i.i.d. standard normal, mean 0 and variance 1. Then define $\underline{X} = \underline{Z}/\|\underline{Z}\|$ (taking the norm as the 2-norm or Euclidean norm), so that \underline{X} is a random vector of length 1, lying exactly on the k -dimensional unit sphere. Note that because of the spherical symmetry of the \underline{Z} , the \underline{X} will be uniformly distributed over the surface of the sphere.

For a given value of $1 > \epsilon > 0$, discard any points where $|X_k| < \epsilon$ - this will create two hemispheres of data points separated by a gap along the equator (in this high dimensional space). Let us classify any remaining point where $X_k \geq \epsilon$ with $Y = +1$, and any remaining point with $X_k \leq -\epsilon$ with $Y = -1$. Note that for a data set defined in this way, we have a very natural perceptron that can be applied:

$$\text{classify}(\underline{x}) = \begin{cases} +1 & \text{if } 0x_1 + 0x_2 + \dots + 0x_{k-1} + 1x_k > 0 \\ -1 & \text{if } 0x_1 + 0x_2 + \dots + 0x_{k-1} + 1x_k < 0, \end{cases} \quad (1)$$

with a linear separator given by $x_k = 0$. This is by no means the only feasible perceptron, but it does show that the data is linearly separable.

Additionally, note that the margin of separation between the positive and negative data classes is *at least* ϵ .

We proved in class that the convergence of the Perceptron Learning Algorithm is bound independently of the dimension of the data and the number of data points, and can be bound entirely in terms of the maximum margin between the two data classes. I'd like you to try to verify this experimentally. For a given value of k and ϵ , and a desired value of m , consider repeatedly generating points \underline{x} but discarding them if they do not satisfy $|x_k| \geq \epsilon$ - repeat this process until m data points have been generated.

We generated the dataset for perceptron as per the given conditions so that we would surely have a linear separator for the generated dataset.

```
import numpy as np
import pandas as pd
from pprint import pprint
import math

def generate_data(k, m, eps):
    count = 0
    data = []
    # mean = 0 standard deviation = 1 of a normal distribution
    while(count != m):
        z = np.random.normal(0, 1, k)
        a = np.sum(np.square(z))
        magnitude = math.sqrt(a)
        X = z / magnitude
        if abs(X[k-1]) < eps:
            continue
        count += 1
        d = X.tolist()
        data.append(d)
    # generate the Y values based on the last feature
```

```

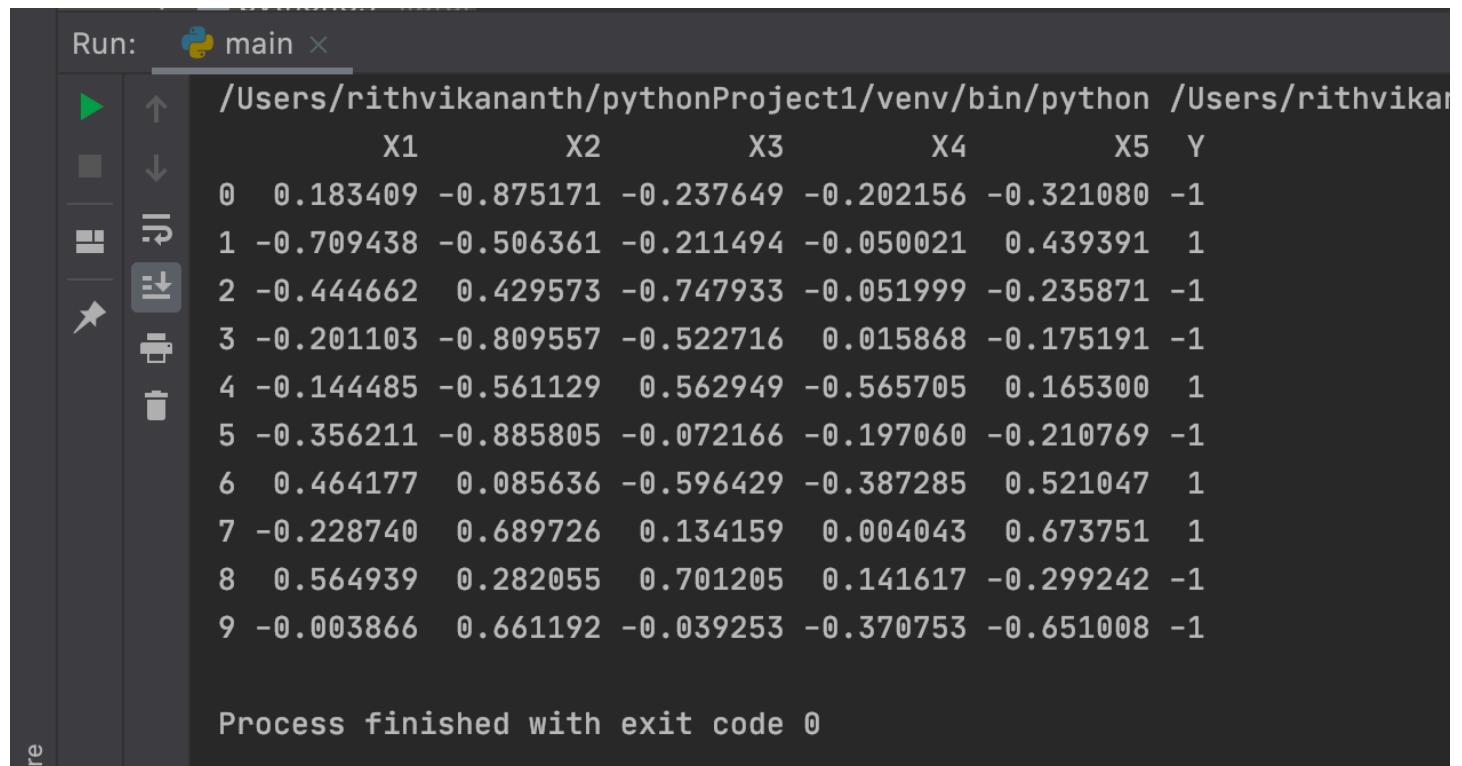
y = []
for row in range(m):
    if data[row][-1] > 0:
        y.append(1)
    else:
        y.append(-1)
# changing the column names
col_names = []
for i in range(1, k + 1):
    col_names.append("X" + str(i))
final_data = pd.DataFrame(data[0:], columns=col_names)
final_data["Y"] = y
return final_data

def sign(row):
    if row[-1] >= 1:
        return 1
    else:
        return -1

training_data = generate_data(5, 10, 0.3)

```

Generated Dataset:



The screenshot shows a Jupyter Notebook terminal window. The title bar says "Run: main". The terminal output displays a table of data with 10 rows and 6 columns (X1 to X5) and a Y column. The first five columns are numerical values, and the last column is labeled Y with values -1 or 1. The data is as follows:

	X1	X2	X3	X4	X5	Y
0	0.183409	-0.875171	-0.237649	-0.202156	-0.321080	-1
1	-0.709438	-0.506361	-0.211494	-0.050021	0.439391	1
2	-0.444662	0.429573	-0.747933	-0.051999	-0.235871	-1
3	-0.201103	-0.809557	-0.522716	0.015868	-0.175191	-1
4	-0.144485	-0.561129	0.562949	-0.565705	0.165300	1
5	-0.356211	-0.885805	-0.072166	-0.197060	-0.210769	-1
6	0.464177	0.085636	-0.596429	-0.387285	0.521047	1
7	-0.228740	0.689726	0.134159	0.004043	0.673751	1
8	0.564939	0.282055	0.701205	0.141617	-0.299242	-1
9	-0.003866	0.661192	-0.039253	-0.370753	-0.651008	-1

At the bottom of the terminal, it says "Process finished with exit code 0".

Perceptron Algorithm:

```
def perceptron(data):
    x = np.asarray(data.iloc[:, :-1])
    b = 0
    w = []
    d = (np.asarray(data.iloc[:, :-1])).tolist()
    d_for_y = (np.asarray(data.iloc[:, :])).tolist()
    for i in range(0, len(d[0])):
        w.append(0)
    steps = 0

    while 1:
        mismatch = 0
        # finding the misclassified point
        for i in range(0, len(d)):
            s = 0
            y = sign(d_for_y[i])
            for j in range(0, len(d[i])):
                s += d[i][j] * w[j]
            s += b
            if s > 0:
                f_y = 1
            else:
                f_y = -1
            if y == f_y:
                continue
            #continue if it is correctly classified
            else:
                # Update w and b and increase step number
                mismatch += 1
                steps += 1
                b += y
                for l in range(0, len(w)):
                    w[l] += d[i][l] * y
    # Terminate if there is no mismatch of Y when iterated through all the points in the
    # dataset
    if mismatch == 0:
        break
    # weights = pd.DataFrame(w[0:])
    # print("Steps: ", steps )
    return steps
```

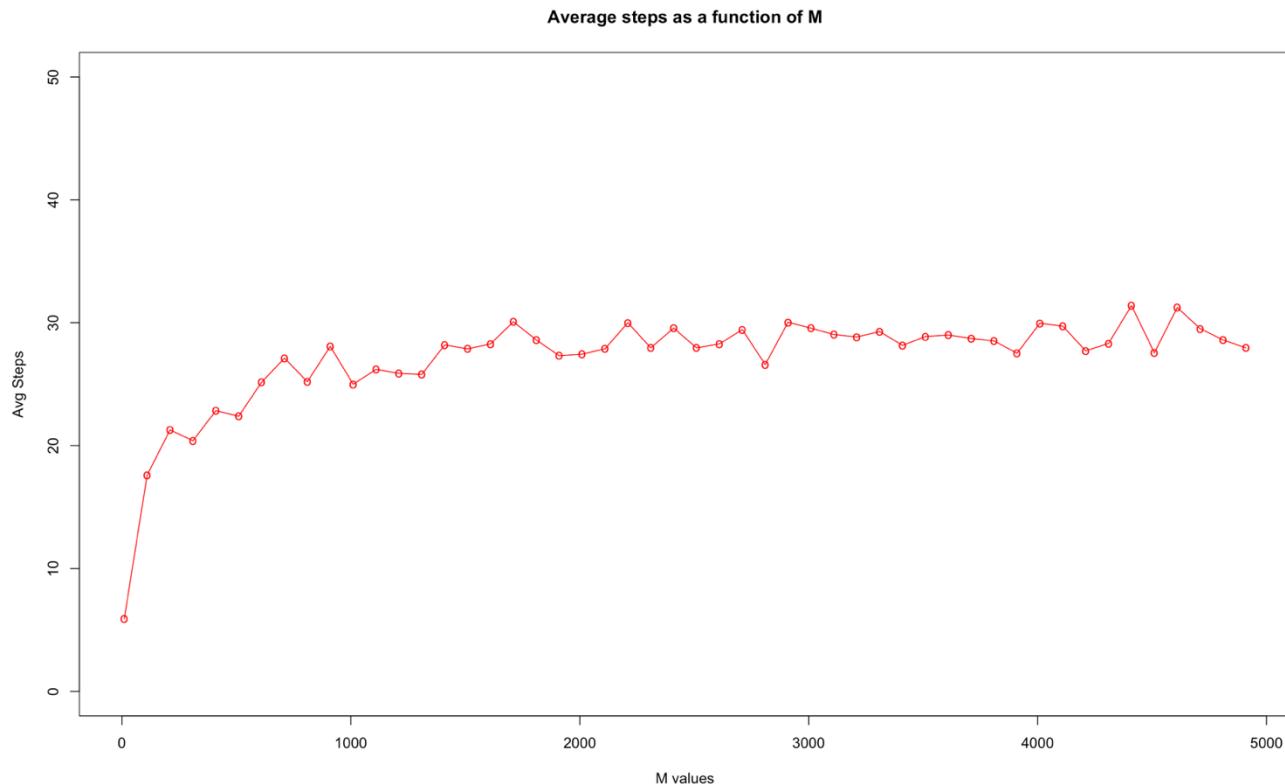
- 1) For $k = 5, \epsilon = 0.1$, for a range of possible m values, repeatedly generate data sets of size m and fit a perceptron to them. Plot, as a function of m , the average number of steps needed for the Perceptron Learning Algorithm to converge. Do your results make sense? Do you think looking at larger and larger m values outside the range you plot will produce anything different?

Finding the average step value for a range of given m values.

```
def function_of_steps():
    #for k = 5 eps = 0.1
    mv = []
    step = []
    for m in range(10, 5000, 100):
        s = []
        mv.append(m)
        for i in range(50):
            data = generate_data(5, m, 0.1)
            steps = perceptron(data)
            s.append(steps)
        step.append((sum(s)/len(s)))
    print(mv)
    print(step)
```

```
103 | function_of_steps()
main x
/Users/rithvikanth/pythonProject1/venv/bin/python /Users/rithvikanth/pythonProject1/main.py
[10, 110, 210, 310, 410, 510, 610, 710, 810, 910, 1010, 1110, 1210, 1310, 1410, 1510, 1610, 1710, 1810, 1910, 2010, 2110, 2210, 2310, 2410, 2510, 2610, 2710, 2810, 2910, [6.3, 16.68, 21.32, 22.36, 21.16, 23.76, 24.24, 25.04, 27.0, 23.84, 24.64, 26.92, 25.84, 27.36, 26.32, 29.24, 27.44, 28.48, 27.04, 28.76, 28.88, 28.48, 29.36, 28.08, 28]
```

We observe an increase in the average step size with the increase in M values initially, but after m-value greater than 2000 we observe that average steps are almost the same. The larger and larger m values outside the range of plot would also maintain the same average steps observed according to my understanding.



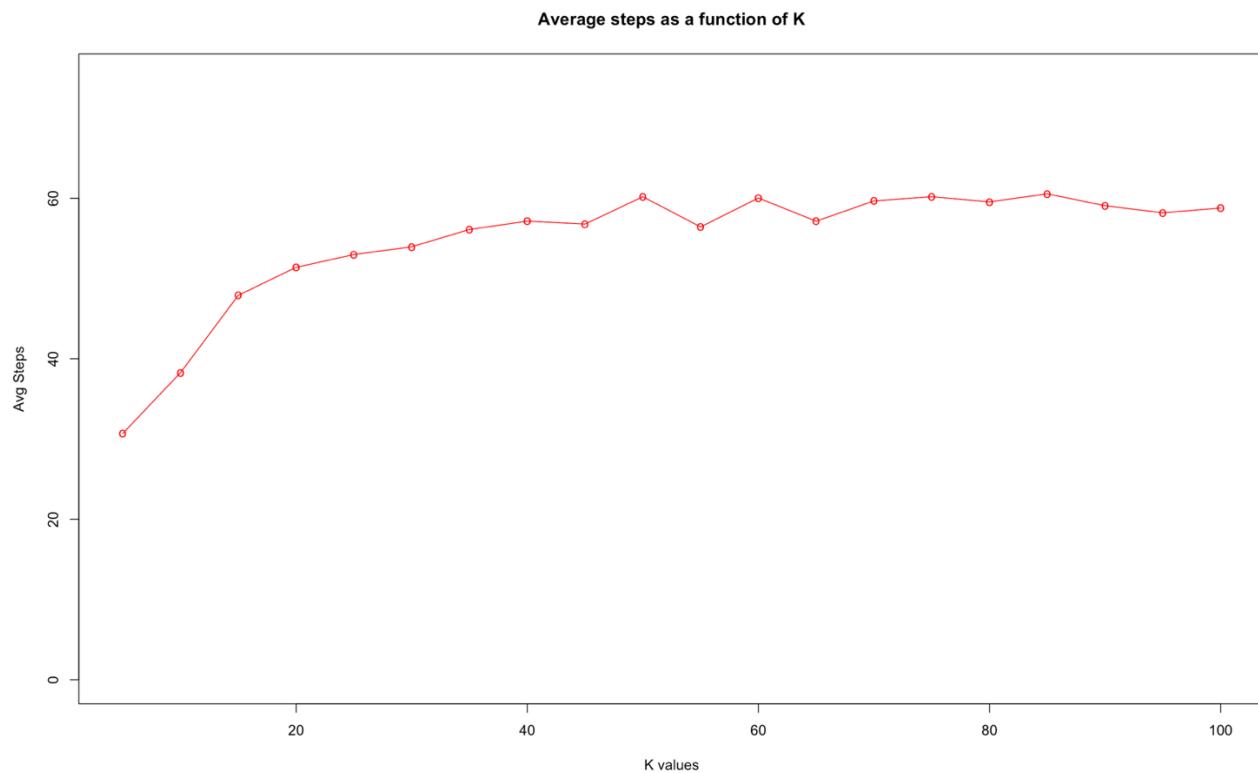
- 2) For $m = 100$, $\epsilon = 0.05$, for a range of possible k values, repeatedly generate data sets of dimension k , and fit a perceptron to them. Plot, as a function of k , the average number of steps needed for the Perceptron Learning Algorithm to converge. Do your results make sense? Do you think looking at larger and larger k values outside the range you plot will produce anything different?

```
def function_of_features():
    kv = []
    step = []
    for k in range(5, 101, 5):
        s = []
        kv.append(k)
        for i in range(50):
            data = generate_data(k, 100, 0.05)
            steps = perceptron(data)
            s.append(steps)
        step.append((sum(s) / len(s)))
    print(kv)
    print(step)
```

```
114 |     function_of_features()

main ×
/Users/rithvikananth/pythonProject1/venv/bin/python /Users/rithvikananth/pythonProject1/main.py
[5, 10, 15, 20, 25, 30, 35, 40, 45, 50, 55, 60, 65, 70, 75, 80, 85, 90, 95, 100]
[30.7, 38.28, 47.92, 51.4, 53.0, 53.96, 56.12, 57.16, 56.8, 60.2, 56.44, 60.04, 57.16, 59.68, 60.2, 59.56, 60.56, 59.08, 58.2, 58.8]
```

A marginal increase in the average step size is observed with the increase in K values but. For K-value greater than 45 the step size is constant. To my understanding, larger and larger K values outside the range of plot would also maintain the same average steps observed.



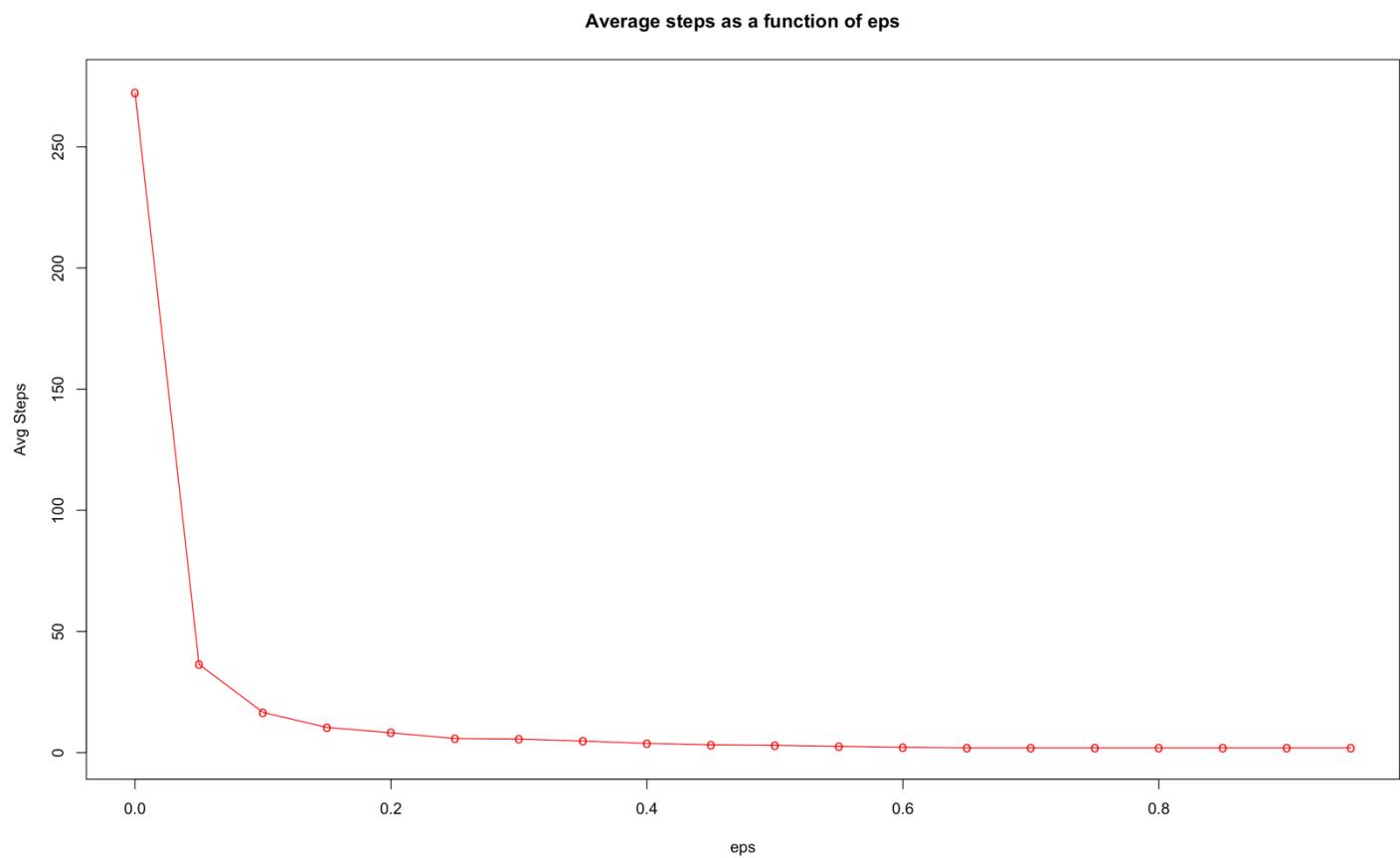
- 3) For $k = 5, m = 100$, for a range of possible ϵ values in $[0, 1]$, repeatedly generate data sets with an ϵ -threshold cutoff and fit a perceptron to them. Plot, as a function of ϵ , the average number of steps needed for the Perceptron Learning Algorithm to converge. Do your results make sense? Do you think looking at different k, m values will produce anything different?

```
def function_of_eps():
    eps_val = []
    step = []
    eps = 0
    for i in range(21):
        s = []
        eps_val.append(eps)
        for j in range(50): # Taking the average
            data = generate_data(5, 100, eps)
            steps = perceptron(data)
            s.append(steps)
        step.append((sum(s) / len(s)))
        eps += 0.05
    print(eps_val)
    print(step)
```

137 | function_of_eps()

main ×

```
/Users/rithvikanth/pythonProject1/venv/bin/python /Users/rithvikanth/pythonProject1/main.py
[0, 0.05, 0.1, 0.15000000000000002, 0.2, 0.25, 0.3, 0.35, 0.3999999999999997, 0.4499999999999996, 0.4999999999999994, 0.5499999999999999, 0.6, 0.65, 0.7001272.3, 36.4, 16.6, 10.4, 8.2, 5.8, 4.8, 3.8, 3.2, 3.0, 2.6, 2.2, 2.0, 2.0, 2.0, 2.0, 2.0, 2.0, 2.0]
```



As the margin increases, we can clearly observe that the average number of steps taken to converge reduces. At margin 0.3 or greater we observe the average step size is 2. Irrespective of what the K and m values are, the graph for epsilon vs average steps would remain the same.

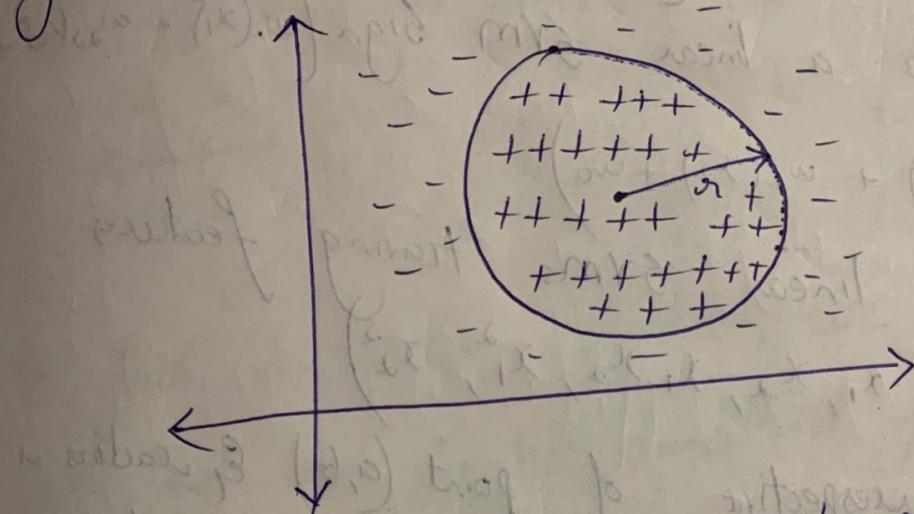
R- code:

```
plot(m, steps, type="o", col="red", pch="o", lty=1, ylim=c(0,50), ylab="Avg Steps", xlab = "M values")  
+ title("Average steps as a function of M")
```

```
plot(K, steps, type="o", col="red", pch="o", lty=1, ylim=c(0,75), ylab="Avg Steps", xlab = "K values") +  
title("Average steps as a function of K")
```

```
plot(eps, steps, type="o", col="red", pch="o", lty=1, ylim=c(0,22), ylab="Avg Steps", xlab = "eps") +  
title("Average steps as a function of eps")
```

Given that positive class lay within a certain radius of a point, negative ones lie outside that radius.



As told in the class of a similar example (unit circle)
the features of the data are x_1, x_2, x_1^2, x_2^2
 $\phi(x_1, x_2) = \begin{pmatrix} 1 & x_1 & x_2 & x_1^2 & x_2^2 \end{pmatrix}$

Equivalently with kernel

$$k(x, y) = (1 + x \cdot y)^2$$

The positive points are defined as

$$x_1^2 + x_2^2 \leq 1 \quad (\text{or}) \quad 0 \leq 1 - x_1^2 - x_2^2$$

In class (for unit circle) classifier was shown as $\Rightarrow \text{Sign}(1 - x_1^2 - x_2^2)$

for the given data at point (a, b) .

with radius a

$$\Rightarrow \text{Sign}(a^2 - (x_1 - a)^2 - (x_2 - b)^2)$$

$$\Rightarrow \text{sign} \left(\alpha^2 - (x_1^2 + \alpha^2 - 2\alpha x_1) - (x_2^2 + b^2 - 2bx_2) \right)$$

$$\Rightarrow \text{sign} \left(\alpha^2 - \alpha^2 - b^2 + 2x_1 \alpha + 2x_2 b - x_1^2 - x_2^2 \right)$$

Now to a linear form $\text{sign} \left(w_0(x_1^2) + w_1(x_1) + w_2(x_2)^2 + \dots \right)$

$$w_0(x_1^2) + w_1(x_1) + w_2(x_2)^2 + \dots$$

with linear terms having features

$$(1, x_1, x_2, x_1 x_2, x_1^2, x_2^2)$$

So irrespective of point (a, b) & radius r
the data will always form a linear separator.

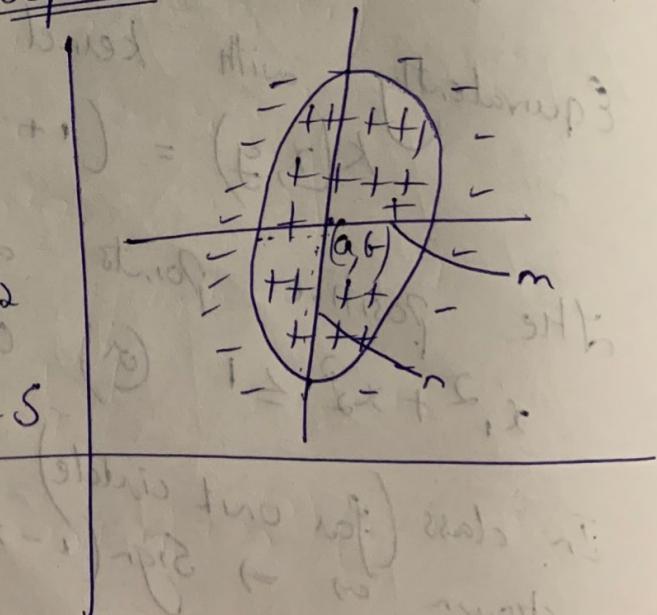
for ellipsoidal Separator

we classify (\pm)

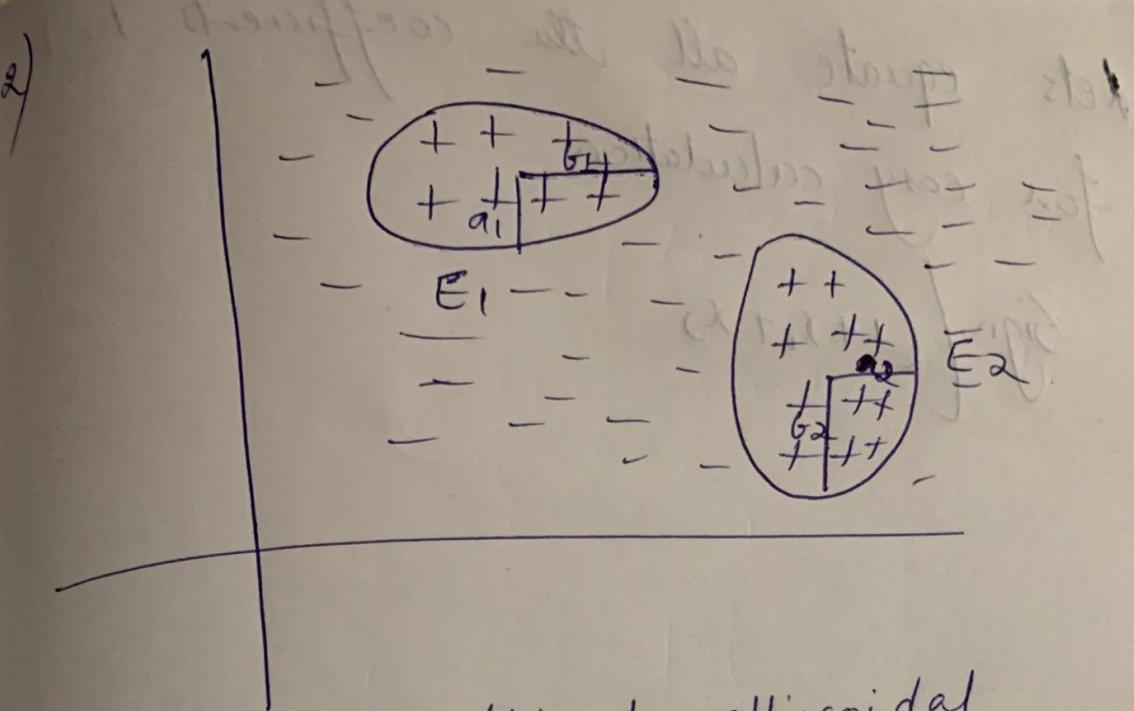
as

$$\text{sign} \left(\frac{x_1^2}{m^2}(x_1 - a)^2 + \frac{x_2^2}{r^2}(x_2 - b)^2 - m^2 r^2 \right) = S$$

$$\Rightarrow \begin{cases} +1 & S \leq 0 \\ -1 & S > 0 \end{cases}$$



We see that the parameters for this equation can also be achieved by the same mapping function.



Consider 2 disjoint ellipsoidal regions with centers (m_1, n_1) & (m_2, n_2) .
 major axis $\sqrt{a_1^2 + b_1^2}$ minor axis $\sqrt{a_2^2 + b_2^2}$

Given that the positive points lie
 inside E_1 & E_2 , for any positive
 class data points (x_1, x_2)

$$\left(1 - \frac{(x_1 - m_1)^2}{a_1^2} - \frac{(x_2 - n_1)^2}{b_1^2} \geq 0 \right) \quad (1)$$

$$\left(1 - \frac{(x_1 - m_2)^2}{a_2^2} - \frac{(x_2 - n_2)^2}{b_2^2} \geq 0 \right)$$

The classifier for above would be

$$\text{sign} \left[\left(a_1^2 b_1^2 - a_1^2 (x_1 - m_1)^2 - b_1^2 (x_2 - n_1)^2 \right) \times \right.$$

$$\left. \left(a_2^2 b_2^2 - a_2^2 (x_1 - m_2)^2 - b_2^2 (x_2 - n_2)^2 \right) \right]$$

Equating all coefficients to 1
for easy output

$$\Rightarrow \text{Sign} \left((1 + x_1 + x_2 + x_1^2 + x_2^2) (1 + x_1 + x_2 + x_1^2 + x_2^2 + x_1^3 + x_2^3 + x_1^4 + x_2^4) \right)$$
$$\text{Sign} \left(1 + x_1 + x_2 + x_1^2 + x_2^2 + x_1^3 + x_2^3 + x_1^4 + x_2^4 + x_1^2 x_2^2 + x_2^2 x_1^2 + x_1^3 x_2^2 + x_2^3 x_1^2 + x_1^4 x_2^2 + x_2^4 x_1^2 \right)$$

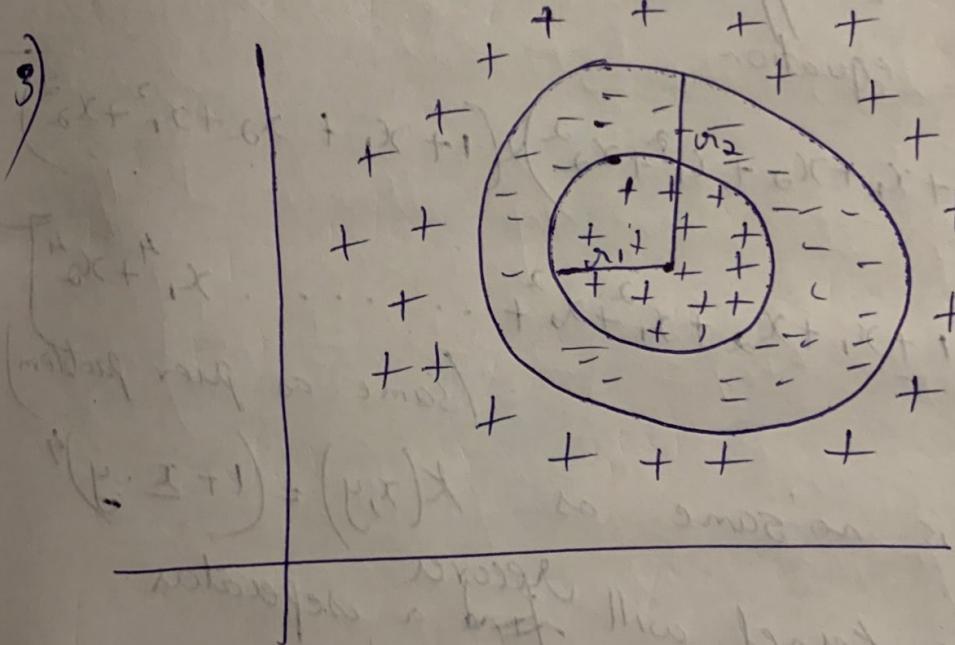
we get other parameters as

$$(1, x_1, x_2, x_1x_2, x_1^2, x_2^2, x_1^2x_2^2, x_1x_2^2 \dots \\ \dots, x_1^4, x_2^4)$$

We can hence use the kernel

function to represent this mapping function through

$$k(x, y) = (1 + x \cdot y)^4$$



Let r_1, r_2 be the radii of inner
outer circles respectively

Center is (a, b)

Positive points :- $(x_1 - a)^2 + (x_2 - b)^2 \leq r_1^2$

(a) $(x_1 - a)^2 + (x_2 - b)^2 > r_2^2$

The classifier for above would be

$$\text{Sign} \left[\begin{aligned} & [\alpha_1^2 - (x_1 - a)^2 - (x_2 - b)^2] \\ & [\alpha_2^2 - (x_1 - a)^2 - (x_2 - b)^2] \end{aligned} \right]$$

$$\text{Sign} \left[\begin{aligned} & [\alpha_1^2 - x_1^2 - a^2 + 2x_1 a - x_2^2 - b^2 + 2x_2 b] \\ & [\alpha_2^2 - x_1^2 - a^2 + 2x_1 a - x_2^2 - b^2 + 2x_2 b] \end{aligned} \right]$$

Equating co-efficients to 1 & calculating
the equation

$$\text{Sign} \left[(1 + x_1 + x_2 + x_1^2 + x_2^2) \times (1 + x_1 + x_2 + x_1^2 + x_2^2) \right]$$

$$\Rightarrow \text{Sign} \left[1 + x_1 + x_2 + x_1^2 + x_2^2 + \dots + x_1^4 + x_2^4 \right] \quad (\text{Same as previous problem})$$

$$\text{This is no same as } K(x, y) = (1 + \Sigma x_i y_i)^4$$

& this kernel will ~~not~~ ^{not record} a separator

4) Given XOR data $(\pm 1, \pm 1)$

x	x_1	x_2	$\frac{4}{-1}$	Result
	+1	-1		
	-1	+1		
	+1	-1		
	+1	+1		

$$k(x, y) = (1 + x \cdot y)^2$$

we can find the values of k as the following

$$K_{11} = \left[1 + \begin{bmatrix} -1 \\ -1 \end{bmatrix} \begin{bmatrix} 1 & -1 \end{bmatrix} \right]^2$$

$$(1+2)^2 = 9.$$

$$K_{12} = \left[1 + \begin{bmatrix} -1 \\ +1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} \right]^2$$

$$\Rightarrow [1+0]^2 \Rightarrow 01$$

K matrix \Rightarrow

$$\begin{bmatrix} 9 & -1 & -1 & -1 \\ -1 & 9 & 1 & -1 \\ -1 & 1 & 9 & -1 \\ -1 & -1 & -1 & 9 \end{bmatrix}$$

$$\begin{bmatrix} k \times \text{desire}(y) \\ \begin{bmatrix} 9 & -1 & -1 & -1 \\ -1 & 9 & 1 & -1 \\ -1 & 1 & 9 & -1 \\ -1 & -1 & -1 & 9 \end{bmatrix} \end{bmatrix} = \begin{bmatrix} 9 & -1 & -1 & -1 \\ -1 & 9 & 1 & -1 \\ -1 & 1 & 9 & -1 \\ -1 & -1 & -1 & 9 \end{bmatrix}$$

Kernel function for SVM

$$\max_{i=1}^m \alpha_i = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i y_j k_{ij} y_j \alpha_j$$

such that $\sum_{i=1}^m \alpha_i y_i = 0$ and $\alpha_i \geq 0$

So now

$$\max_{\lambda_1, \lambda_2, \lambda_3, \lambda_4} (\lambda_1 + \lambda_2 + \lambda_3 + \lambda_4) - \frac{1}{2} \left[9\lambda_1^2 + 9\lambda_2^2 + 9\lambda_3^2 + 9\lambda_4^2 + 2(-\lambda_1\lambda_2 + \lambda_2\lambda_3 + \lambda_3\lambda_4 - \lambda_3\lambda_4 - \lambda_2\lambda_4 - \lambda_1\lambda_3) \right]$$

We have max fn with $\lambda_1, \lambda_2, \lambda_3, \lambda_4$

$$\max_{\lambda_1, \lambda_2, \lambda_3, \lambda_4} \left(\sum_{i=1}^4 \lambda_i \right) - \frac{9}{2} \left[\lambda_1^2 + \lambda_2^2 + \lambda_3^2 + \lambda_4^2 \right] + 2(-\lambda_1\lambda_2 + \lambda_2\lambda_3 + \lambda_3\lambda_4 - \lambda_3\lambda_4 - \lambda_2\lambda_4 - \lambda_1\lambda_3)$$

Solving the above eqn Max is at $1/4$

$$(\lambda_1, \lambda_2, \lambda_3, \lambda_4) = \left(\frac{1}{8}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8} \right)$$

We find other sign function

$\underline{\omega} - \underline{x}$ is given by

$$\sum_{i=1}^m d_i y_i k(x_i, \underline{x})$$

$$\Rightarrow -\lambda_1 y_1 k_{1x} + \lambda_2 y_2 k_{2x} + \lambda_3 y_3 k_{3x} + \lambda_4 y_4 k_{4x}$$

$$\Rightarrow \frac{1}{8} [-k_{1x} + k_{2x} + k_{3x} - k_{4x}]$$

bias $b = y^i - \underline{\omega} x$

$$\underline{\omega} x_i = \frac{1}{8} [-k_{11} + k_{21} + k_{31} - k_{41}]$$

$$= \frac{1}{8} [-1 + 1 + 1 - 1]$$

$$= -1$$

$$b = -1 - (-1) = 0$$

the separator

$$\frac{1}{8} [-k_{1x} + k_{2x} + k_{3x} - k_{4x}] + 0$$

$$K = (1 + x_i \cdot x)^2$$

$$\Rightarrow K(x, y) = e^{-\frac{(x-y)^2}{2}}$$

calculating the kernel values

$$|x-y| = (x_1 - y_1)^2 + (x_2 - y_2)^2$$

$$= e^{-(\epsilon_1+1)^2 + (\epsilon_1-1)^2}$$

$$K_{11} = e^0 = 1$$

$$K_{12} = e^{-(\epsilon_1+1)^2 + (\epsilon_1-1)^2}$$

$$= e^{-4} = \frac{1}{e^4}$$

Doing so for XOR we get
the K matrix as

$$\begin{bmatrix} 1 & e^{-4} & e^{-4} & e^{-8} \\ e^{-4} & 1 & e^{-8} & e^{-4} \\ e^{-4} & e^{-8} & 1 & e^{-4} \\ e^{-8} & e^{-4} & e^{-4} & 1 \end{bmatrix}$$

Finding

$$\max_{x_i} \sum_{i=1}^m \lambda_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda_i y_i^T k_{ij} y_j$$

We get $\lambda = [e^{0.03}, e^{0.03}, e^{-0.03}, e^{-0.03}]$

by solving the eqn

$$\lambda \approx [1.03, 1.03, 1.03, 1.03]$$

$$b = y_i - \underline{\omega} x_i$$

$$\underline{\omega} x_i = \sum_{i=1}^m \lambda_i y_i^T k(x_i, x)$$

$$= 1.03 [-k_{1x} + k_{2x} + k_{3x} - k_{4x}]$$

Put $x_i = 1$

$$1.03 [-1 + e^{-4} + e^{-4} - e^{-8}]$$

$$\Rightarrow 1.03 - 0.999 \approx -0.001$$

$$b = y_i - \underline{\omega} x_i = -1 - (-1) = 0$$

Separator is

$$1.03 [-k_{1x} + k_{2x} + k_{3x} - k_{4x}] + 0$$

$$k_{1x} = e^{-|x_1 - x|^2}$$