

CS 536 : PERCEPTRONS AND SVMs-ASSIGNMENT Nº 3

Aravinda Reddy Dandu, Rutgers University

10/24/2020

Contents

1 Perceptrons	1
1.1 Data generation	1
1.2 Perceptron algorithm	3
1.2.1 Information being returned	3
1.3 Question 1.1	5
1.4 Question 1.2	9
1.5 Question 1.3	11
1.6 Bonus Question	13
1.6.1 Varying m	14
1.6.2 Varying k	14
1.6.3 Varying ϵ	15
2 Support Vector Machines	15

1 Perceptrons

1.1 Data generation

Given method to generate data:

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.

Data is generated using the below function. This function serves as a general purpose data generator for all further operations in perceptrons

```

1 # Function to generate perceptron data. Takes epsilon, size of data and dimension ←
  as inputs
2 def gen_perceptron_data(e, m, k):
3     dataset = {}
4     # Initializing dataset storage with header
5     for x in range(1, k + 1):
6         dataset['x' + str(x)] = []
7     dataset['y'] = []
8     # Looping until dataset reaches required size which is m
9     while len(dataset['x1']) < m:
10         # Drawing k samples out of normal distribution using mean 0 and variance 1
11         z = np.random.normal(0, 1, k)
12         # Taking euclidian norm of Z
13         z_norm = np.linalg.norm(z, 2)
14         # Dividing Z vector with norm of Z to get X
15         x = z / z_norm
16         y = None
17         # checking if absolute value is greater than epsilon. Continuing otherwise
18         if abs(x[len(x) - 1]) > e:
19             # Assigning positive if xk is positive(will obviously be greater than ←
              e). Negative else
20             if x[len(x) - 1] > e:

```

```

21         y = 1
22     else:
23         y = -1
24     else:
25         continue
26     # Including X vector in dataset
27     for index, xk in enumerate(x):
28         dataset['x' + str(index + 1)].append(xk)
29     # Appending Y values to last column
30     dataset['y'].append(y)
31 # Returning a pandas dataframe
32 dataset = pd.DataFrame.from_dict(dataset)
33 return dataset

```

A sample data generated using k as 5, m as 100 and ϵ as 0.2 looks like below

	x1	x2	x3	x4	x5	y
0	-0.637690	-0.254533	-0.469553	-0.382885	0.401851	1
1	-0.105606	-0.212329	0.316824	-0.087033	-0.914227	-1
2	-0.125693	0.171695	0.715603	0.206205	-0.632545	-1
3	-0.272733	-0.556376	0.396314	0.612019	-0.290568	-1
4	0.778924	-0.269206	0.089614	-0.439899	-0.345346	-1
..
95	-0.149478	-0.248669	0.672143	0.054015	-0.679062	-1
96	-0.690397	0.511306	-0.000118	-0.203999	-0.469364	-1
97	-0.398089	0.392455	-0.060761	-0.438334	0.701196	1
98	0.422985	-0.606463	-0.248129	-0.521433	-0.346158	-1
99	0.309706	0.262561	0.355346	-0.035801	0.841184	1

[100 rows x 6 columns]

1.2 Perceptron algorithm

A class is written for Perceptron which has a method to generate a perceptron for any given data. Here we assume that data is linearly separable and try to find weights and bias for a corresponding dataset. Class variables will be weight vector and bias value because these are the only values which define a perceptron.

Initializing with weight vector as 0 and bias as 0, we keep moving the perceptron until all the points are correctly classified.

1.2.1 Information being returned

Number of steps taken to converge : A step is defined as one move where weight and bias value is changed. Every time we encounter a misclassified point, we change the weight and bias and it counts as a *step*

Weight vector and bias : After all the points are correctly classified by a given perceptron, respective weights and bias values will be stored in the class variables.

1 # Main class to fit a perceptron on given data

```

2 class Perceptron():
3     def __init__(self):
4         pass
5
6     # Function to fit a perceptron on given data
7     def fit_perceptron(self, data):
8         # Splitting the dataset for features and classifiers
9         X_vectors = np.asarray(data.iloc[:, :-1])
10        y_values = np.asarray(data.iloc[:, -1:])
11        # Getting size of sample and dimension from first split
12        m, k = X_vectors.shape
13
14        # Declaring to class variables weights and bias. These will be updated for every mis classified point
15        self.weights = np.zeros(shape=(k, 1))
16        self.bias_val = 0
17
18        # Initializing number of steps count
19        number_of_steps = 0
20        # Vector for storing y values calculated by perceptron at that instance
21        Calculated_val = [0] * m
22        i = 0
23        # Running an indefinite loop until all points are correctly classified. Here we know that perceptron will converge because of the dataset we generate
24        while True:
25            # Looping over the dataset looking for misclassified points
26            for i in range(m):
27                # Getting the dot product of weights and current X vector. This is w1*x1+ w2*x2 + ..... +wk*xk
28                if (np.dot(X_vectors[i].reshape((k, 1)).T, self.weights) + self.bias_val)[0][0] > 0:
29                    Calculated_val[i] = 1
30                else:
31                    Calculated_val[i] = -1
32                # If calculated value agrees with y value of vector, perceptron correctly classified the point
33                if float(Calculated_val[i]) != float(y_values[i]):
34                    # Updating weights and bias if wrongly classified
35                    self.weights += np.dot(X_vectors[i].reshape((k, 1)), y_values[i].reshape((1, 1)))
36                    self.bias_val += y_values[i]
37                    number_of_steps += 1
38                # Storing all calculated values in the dataset
39                Calculated_val = np.asarray(Calculated_val).reshape((m, 1))
40                i += 1
41                # If all points in dataset correctly classified, break the loop and return the number of steps taken. The
42                # class will have weights and bias values if needed
43                if (np.array_equal(y_values, Calculated_val)) or (number_of_steps >= 10000):
44

```

```

45         break
46     return number_of_steps

```

1.3 Question 1.1

- 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?

Here we generate 20 datasets for each value of m and vary m over a range of 10 to 5000. For each value of m , average number of steps taken is calculated.

Code for running iterations and calculating values below. This code also includes logic for bonus question(will be explained later).

```

1 # Varying M and getting number of steps taken and distance from ideal perceptron ←
2     in each case
3 def test_vary_m():
4     storage = []
5     # Range of m values from 10, 5000 with increments of 100
6     for m in tqdm(range(10, 5000, 100)):
7         steps = []
8         weights = []
9         biases = []
10        # Looping 20 times and taking average number of steps to get a better ←
11            estimate
12        for i in range(0, 20):
13            # Get perceptron data with e = 0.1, m as a variable and k as 5
14            data = gen.gen_perceptron_data(0.1, m, 5)
15            pt = Perceptron()
16            steps.append(pt.fit_perceptron(data))
17            # Normalizing all the weights to make them uniform
18            den = np.sqrt(np.square(np.linalg.norm(pt.weights)) + np.square(pt.←
19                bias_val))
20            pt.weights = pt.weights / den
21            pt.bias_val = pt.bias_val / den
22            # Get weights of perceptron and store them
23            weights.append(pt.weights)
24            # Get biases of perceptron and store them
25            biases.append(pt.bias_val)
26            # Calculating average weights using np.mean and axis as 0
27            avg_w = np.mean(weights, axis=0)
28            # Calculating bias the same way
29            avg_b = np.mean(biases, axis=0)
30            # Getting the distance from "ideal" perceptron
31            dist = get_dist_from_ideal(weights=avg_w, bias=avg_b)
32            # Returning all values calculated
33            storage.append(str(m) + "," + str(np.average(steps)) + "," + str(dist))

```

31 return storage

This gives below output

```
1 m,steps,dist
2 10,5.5,0.019489341654281102
3 110,18.0,0.003741245883055618
4 210,20.1,0.0008855192074725559
5 310,23.0,0.0005059914666722962
6 410,23.2,0.00025653320008841673
7 510,22.7,0.00029609120136727437
8 610,21.8,0.0007503066622136762
9 710,22.3,0.0003418595277035938
10 810,24.6,0.00014717323003123932
11 910,23.1,0.00023137543389946573
12 1010,26.6,0.0009667901065258704
13 1110,25.2,0.00012727281759664525
14 1210,26.7,0.0004519138170351651
15 1310,28.6,0.000380831454051946
16 1410,27.7,0.0003024906779735875
17 1510,29.0,0.0009825533419468453
18 1610,27.1,0.0019581086302622006
19 1710,26.7,0.00012581869760750503
20 1810,27.3,0.0005850925989935943
21 1910,25.9,0.00019724733504510606
22 2010,29.6,0.00022531452223603434
23 2110,25.9,0.0005475133482714749
24 2210,30.0,0.00034486506066963763
25 2310,27.5,0.0006426021847137075
26 2410,27.2,0.00040621419264086177
27 2510,31.1,0.00043452634846506466
28 2610,30.5,0.0005569729728738002
29 2710,29.7,0.0006764970999968959
30 2810,25.5,0.00033635171409906995
31 2910,28.0,0.0004318544179080017
32 3010,28.5,0.0005534896992075317
33 3110,27.8,0.0003803624637624665
34 3210,28.6,0.00038719553227329473
35 3310,29.8,0.00011749664108536564
36 3410,27.2,0.00015566645416848981
37 3510,27.4,0.0002740062616987382
38 3610,27.9,0.00027755673650489975
39 3710,27.1,0.0005468490446676837
40 3810,28.4,0.0002611695784311267
41 3910,28.4,0.00020052087857995798
42 4010,31.8,0.0007349533386939369
43 4110,27.6,0.00036789631161839004
44 4210,30.3,8.138274817110553e-05
45 4310,28.6,0.0005782565382362639
46 4410,31.7,0.00035775107189336107
```

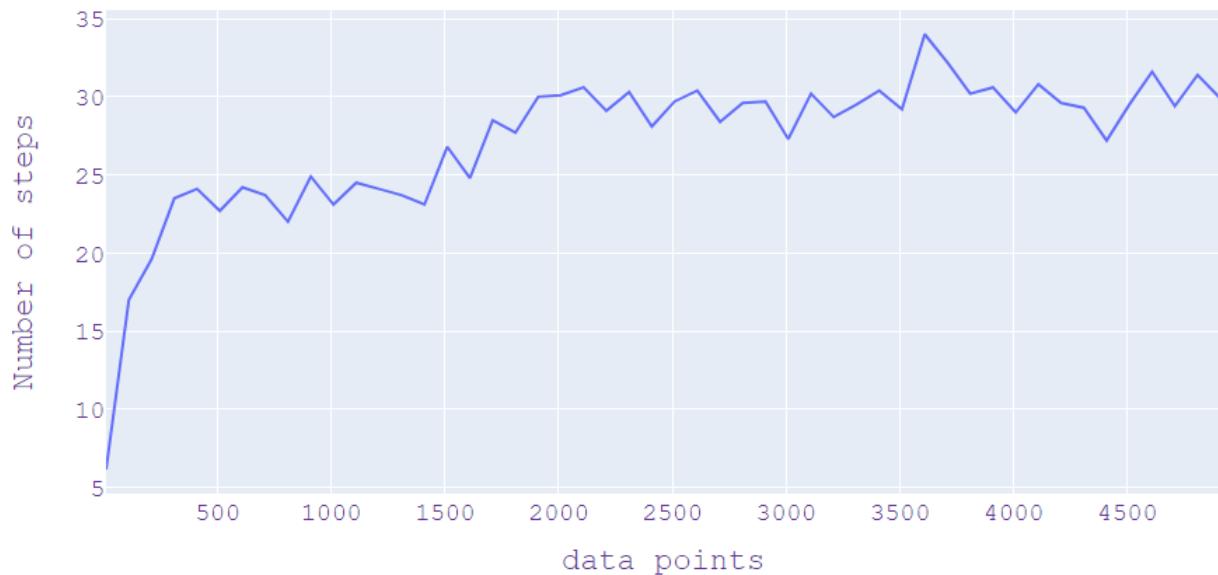
```
47 4510,33.0,0.000587761788711138
48 4610,29.3,0.00019728731303329835
49 4710,31.8,4.509583525873702e-05
50 4810,27.3,0.0002902779624813425
51 4910,33.1,0.00011546399748538787
```

Plotting using plotly

```
1 import plotly.graph_objects as go
2 import pandas as pd
3 cols=['m', 'steps', 'dist']
4 data= pd.read_csv("D:/Study/ML/Perceptrons_SVMs/data_m_var.csv",usecols=cols)
5 fig = go.Figure()
6
7 fig.add_trace(
8     go.Scatter(
9         x=data['m'],
10        y=data['steps'],
11    ))
12
13
14 fig.update_layout(
15     title="Data size-number of steps plot",
16     xaxis_title="data points",
17     yaxis_title="Number of steps",
18     font=dict(
19         family="Courier New, monospace",
20         size=18,
21         color="RebeccaPurple"
22     )
23 )
```

Number of steps taken versus data size will give below graph

Data size-number of steps plot



We can see that number of steps will converge at somewhere **between 25 and 35**. Even after increasing m drastically, number of steps will remain constant. This proves that perceptron complexity is independent of data size

Coming to the first few data points where it starts rising and converges, these data points are between 10 and 350. Reason why these m values have lesser number of steps is *Probability that at least one point falling at the separation line(Gap that we kept while generating dataset) is very less with less data points*. As points are far from separator, it gets easier to classify. Once a data point is located at the separation line, number of steps needed will remain regardless of number of data points.

This can be proven programmatically by repeatedly generating datasets of size 100 and checking number of times it is less than 25(converging value).

```
1 def test_m_100():
2     storage = []
3     for m in [100]:
4         steps = []
5         for i in range(0, 1000):
6             data = gen.gen_perceptron_data(0.1, m, 5)
7             pt = Perceptron()
8             nu = pt.fit_perceptron(data)
9             steps.append(nu)
10        storage.append(steps)
11    # storage.append(str(m) + "," + str(np.max(steps)))
12 steps = storage[0]
13 lesscnt = 0
14 for n in steps:
15     if n < 25:
16         lesscnt += 1
17 print('Number of times steps is less than 20 is ' + str(lesscnt) + ' out of ←
18 1000')
19 print('maximum steps is ' + str(max(steps)))
```

```
19
20 Number of times steps is less than 20 is 909 out of 1000
21 maximum steps is 40
```

Here we can see 909/1000 times, number of steps is less than 25. This shows, probability if points falling on separation line is low. This is why many steps is not required to converge as gap is large.

1.4 Question 1.2

- 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?

Same as m , we vary k for range between 1 and 200

```
1 # Varying K and getting number of steps taken and distance from ideal perceptron ←
  in each case
2 def test_vary_k():
3     storage = []
4     # Range of k values from 1, 200 with increments of 5
5     for k in tqdm(range(1, 200, 5)):
6         steps = []
7         weights = []
8         biases = []
9         # Looping 50 times and taking average number of steps to get a better ←
           estimate
10        for i in range(0, 50):
11            # Get perceptron data with e = 0.05, m as 100 and k as a variable
12            data = gen.gen_perceptron_data(0.05, 100, k)
13            pt = Perceptron()
14            steps.append(pt.fit_perceptron(data))
15            den = np.sqrt(np.square(np.linalg.norm(pt.weights)) + np.square(pt.←
               bias_val))
16            pt.weights = pt.weights / den
17            pt.bias_val = pt.bias_val / den
18            # Get weights of perceptron and store them
19            weights.append(pt.weights)
20            # Get biases of perceptron and store them
21            biases.append(pt.bias_val)
22            # Calculating average weights and biases using np.mean and axis as 0
23            avg_w = np.mean(weights, axis=0)
24            avg_b = np.mean(biases, axis=0)
25            # Getting the distance from "ideal" perceptron
26            dist = get_dist_from_ideal(weights=avg_w, bias=avg_b)
27            # Returning all values calculated
28            storage.append(str(k) + "," + str(np.average(steps)) + "," + str(dist))
29    return storage
```

This gives below data

```
1 k,steps,dist
2 1,1.0,0.5857864376269051
3 6,35.12,0.000792159091595642
4 11,41.96,0.00031879176668000563
5 16,48.28,0.002254649180593259
6 21,51.08,0.003400524838979116
7 26,54.04,0.0038550581261809582
8 31,54.64,0.003711919347351471
9 36,55.0,0.008862846496991766
10 41,56.28,0.00614584153155933
11 46,58.12,0.007023159307302922
12 51,56.56,0.007688828000580725
13 56,55.84,0.011459824234848238
14 61,60.0,0.011155824076315938
15 66,58.6,0.00849497549660349
16 71,58.8,0.01325382244363926
17 76,60.8,0.01110629363120088
18 81,55.12,0.016952300943098672
19 86,58.72,0.015701964987362517
20 91,59.64,0.015569634807497248
21 96,58.8,0.019670416214498727
22 101,57.04,0.02201314269528506
23 106,60.72,0.01766211538212445
24 111,59.68,0.017378414236108863
25 116,60.32,0.017127343120082607
26 121,58.52,0.021073563037323682
27 126,57.0,0.019483854280533577
28 131,60.88,0.024178362458176155
29 136,56.84,0.02630801908628984
30 141,59.56,0.02467376302048073
31 146,59.64,0.026099787316625253
32 151,58.36,0.022567438467676185
33 156,58.52,0.024212178692972354
34 161,60.36,0.02593934400425663
35 166,57.76,0.025307979318785707
36 171,59.16,0.028588054401244144
37 176,58.8,0.026819771515990912
38 181,59.36,0.033574121955407696
39 186,56.84,0.028033870434750377
40 191,58.68,0.028245189496670018
41 196,57.72,0.030840883821933104
```

Plotting the output

Data dimension-number of steps plot



Here we can see that the line neatly converges to $\tilde{60}$ and irrespective of increase of dimension, stays there. So, we can infer that complexity is independent of dimension as well.

1.5 Question 1.3

- 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?

Same way, we generate data by varying ϵ between 0 and 1. Step size is taken as 0.02

```

1 # Varying M and getting number of steps taken and distance from ideal perceptron ←
  in each case
2 def test_vary_e():
3     storage = []
4     # Range of e values from 0.01, 1 with increments of 0.02
5     for e in tqdm(range(1, 100, 2)):
6         steps = []
7         weights = []
8         biases = []
9         e = e / 100
10        # Looping 30 times and taking average number of steps to get a better ←
11          estimate
12        for i in range(0, 30):
13            # Get perceptron data with e as a variable, m as 100 and k = 5
14            data = gen.gen_perceptron_data(e, 100, k=5)
15            pt = Perceptron()
16            steps.append(pt.fit_perceptron(data))
17            den = np.sqrt(np.square(np.linalg.norm(pt.weights)) + np.square(pt.←
18                          bias_val))
18            pt.weights = pt.weights / den

```

```

18         pt.bias_val = pt.bias_val / den
19         # Get weights and biases of perceptron and store them
20         weights.append(pt.weights)
21         biases.append(pt.bias_val)
22         # Calculating average weights and biases using np.mean and axis as 0
23         avg_w = np.mean(weights, axis=0)
24         avg_b = np.mean(biases, axis=0)
25         # Getting the distance from "ideal" perceptron
26         dist = get_dist_from_ideal(weights=avg_w, bias=avg_b)
27         # Returning all values calculated
28         storage.append(str(e) + "," + str(np.average(steps)) + "," + str(dist))
29     return storage

```

This gives below data

```

1 e,steps,dist
2 0.01,79.3,0.000251573281694657
3 0.03,44.166666666666664,0.00038497288786862836
4 0.05,30.466666666666665,0.0008442735446659691
5 0.07,23.4,0.0006796210666929297
6 0.09,18.0,0.00045169585805523825
7 0.11,15.866666666666667,0.0008552694318976155
8 0.13,13.33333333333334,0.0006728406190542547
9 0.15,11.466666666666667,0.0008774147777884127
10 0.17,9.93333333333334,0.006073310901259595
11 0.19,9.266666666666667,0.0024628739411135855
12 0.21,8.4,0.0017496265566388533
13 0.23,7.4,0.002741258975014222
14 0.25,6.0,0.0018680011915344326
15 0.27,5.73333333333333,0.003606969573007876
16 0.29,5.066666666666666,0.002600249430509449
17 0.31,5.266666666666667,0.006418991432606572
18 0.33,5.03333333333333,0.010954987395197702
19 0.35,4.73333333333333,0.003567297494209923
20 0.37,4.2,0.0034202039673193866
21 0.39,4.0,0.005934374635189128
22 0.41,3.733333333333334,0.0014837974666425764
23 0.43,3.666666666666665,0.0021682842515879323
24 0.45,3.333333333333335,0.009082218527788934
25 0.47,3.4,0.0014736669621873565
26 0.49,3.0,0.004366047165908399
27 0.51,3.066666666666667,0.009538339063128794
28 0.53,2.666666666666665,0.004030965159470492
29 0.55,2.4,0.008314647429324758
30 0.57,2.466666666666667,0.011931785718988914
31 0.59,2.13333333333333,0.002715749944516967
32 0.61,2.2,0.005419412874700856
33 0.63,2.1,0.014677085603064775
34 0.65,2.0,0.0028529555732480626
35 0.67,2.0,0.008392828839299626

```

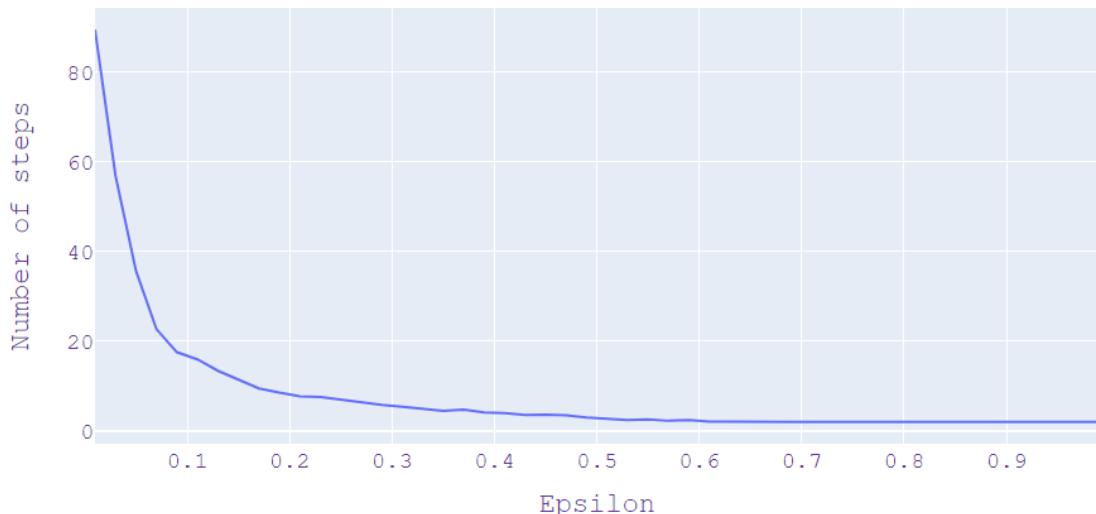
```

36 0.69,2.0,0.005282135462813741
37 0.71,2.0,0.009546609580059049
38 0.73,2.0,0.010786362991078912
39 0.75,2.0,0.004168047709009224
40 0.77,2.0,0.015064499366818236
41 0.79,2.0,0.018148851550467797
42 0.81,2.0,0.01448580727980522
43 0.83,2.0,0.002171327448825599
44 0.85,2.0,0.0005352089681034504
45 0.87,2.0,0.002831933180653132
46 0.89,2.0,0.0038489144130229116
47 0.91,2.0,0.0020277901924322084
48 0.93,2.0,0.0019955801434043555
49 0.95,2.0,0.0013001818080805889
50 0.97,2.0,0.0006986272544607081
51 0.99,2.0,0.00010029932603725317

```

Plotting the above data gives this graph. Using plotly

Separator(epsilon)-number of steps plot



We can clearly say that *larger the gap, easier to find a perceptron*. Here for smaller ϵ values, it takes larger number of steps as gap is small. For higher values, gap increases and it gets easier to solve.

1.6 Bonus Question

A method to calculate distance of any perceptron from the ideal perceptron is written. We normalize the input vector to calculate distance.

```

1 # Method to get distance of any perceptron my the ideal perceptron
2 def get_dist_from_ideal(weights, bias):
3     # getting dimension from input
4     k = len(weights)
5     # Calculating the denominator for normalizing the input perceptron to 1
6     den = np.sqrt(np.square(np.linalg.norm(weights)) + np.square(bias))

```

```

7     weights = weights / den
8     bias = bias / den
9     # Defining ideal perceptron
10    ideal_weights = np.zeros(shape=(k, 1))
11    ideal_weights[k - 1] = 1
12    ideal_bias_val = 0
13    # Calculating distance
14    distance = np.square(np.linalg.norm(weights - ideal_weights)) + np.square(np.←
15      linalg.norm(bias - ideal_bias_val))
15    return distance

```

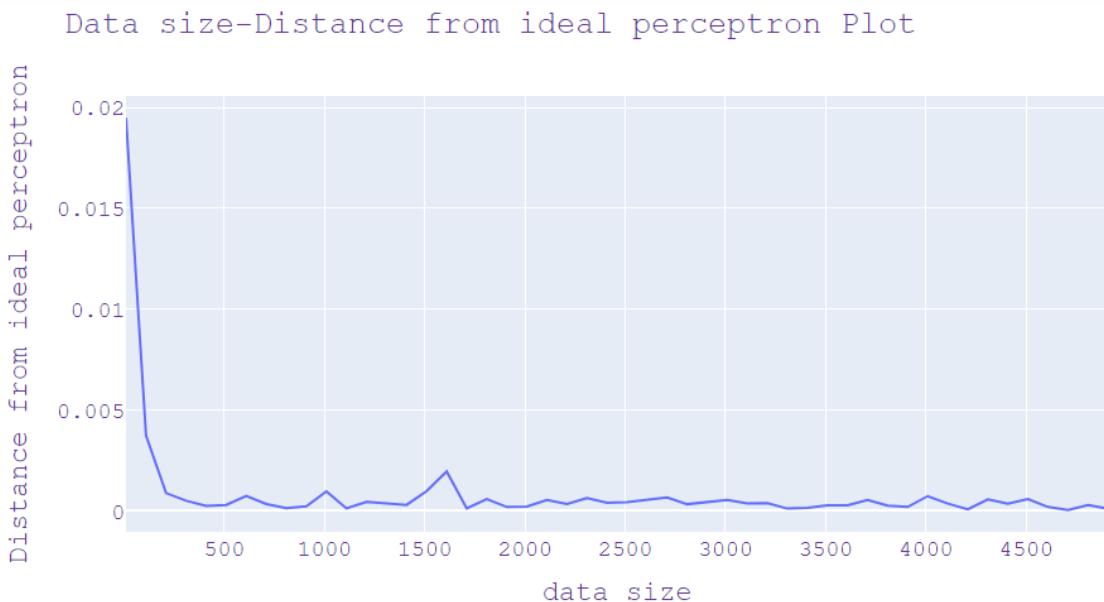
In all the methods to calculate number of steps while varying different parameters, we also store the weights and biases.

An important point to note is that, we have to normalize each weight vector before finding average of that specific parameter value. Will be explained below

1.6.1 Varying m

While varying, m , distance from ideal vector converges very fast at 200 data size to a small value of 0.0008855192074725559. For smaller values of m , error is greater because gap γ is larger. Because of this, perceptron algorithm stops after correctly classifying points and gap still remains.

Plotting values gives below graph



1.6.2 Varying k

While varying, k error rate slowly increases. My intuition is that- For a fixed size of $m=100$, increasing the dimension, makes larger gap in higher dimension and perceptron easily converges at some farther value from the ideal one which is at the centre. For higher values of k , it stays at $\tilde{0.02}$ for many values

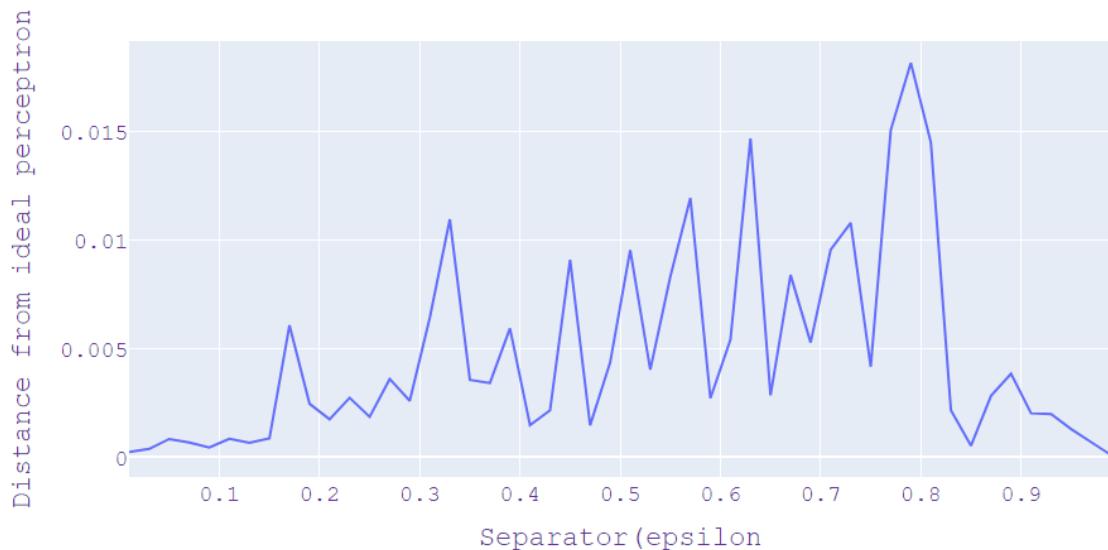
Data dimension-Distance from ideal perceptron Plot



1.6.3 Varying ϵ

While varying ϵ , distance from ideal vector doesn't make any sense. It varies wildly from 0 to 1.

Separator(epsilon)-Distance from ideal perceptron Plot



2 Support Vector Machines

SVMs

Saturday, October 24, 2020 1:15 PM

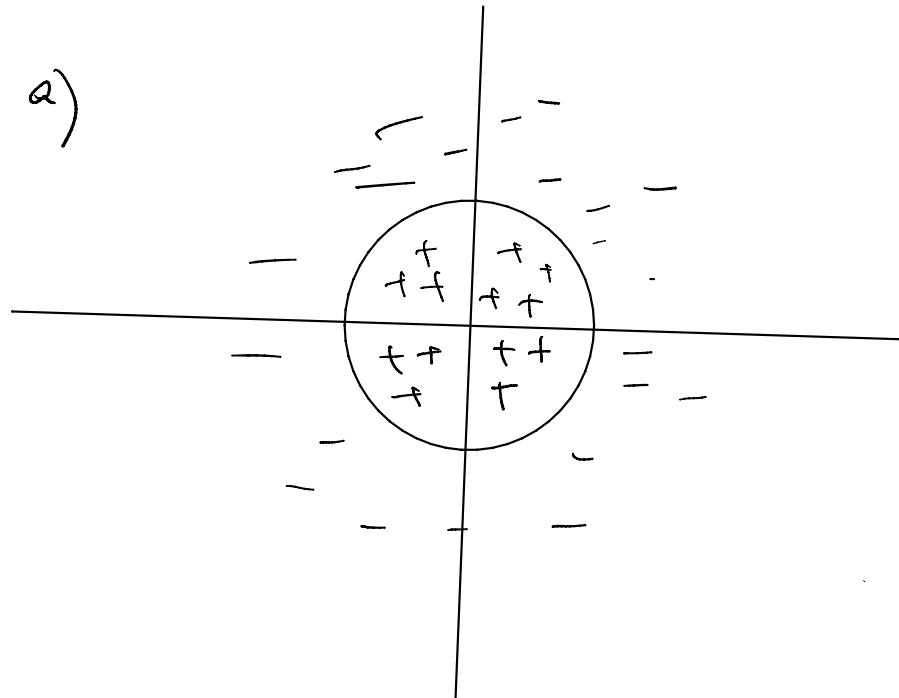


HW-
Perceptro...

SVMs

- 1) Suppose you had a data set in two dimensions that satisfied the following: the positive class all lay within a certain radius of a point, the negative class all lay outside that radius.
 - Show that under the feature map $\phi(x_1, x_2) = (1, x_1, x_2, x_1x_2, x_1^2, x_2^2)$ (or equivalently, with the kernel $K(\underline{x}, \underline{y}) = (1 + \underline{x} \cdot \underline{y})^2$), a linear separator can always be found in this embedded space, *regardless of radius and where the data is centered*.
 - In fact show that if there is an ellipsoidal separator, regardless of center, width, orientation (and dimension!), a separator can be found in the quadratic feature space using this kernel.
- 2) As an extension of the previous problem, suppose that the two dimensional data set satisfied the following: the positive class lay within one of two (disjoint) ellipsoidal regions, and the negative class was everywhere else. Argue that the kernel $K(\underline{x}, \underline{y}) = (1 + \underline{x} \cdot \underline{y})^4$ will recover a separator.
- 3) Suppose that the two dimensional data set is distributed like the following: the positive class lays in a circle centered at some point, the negative class lies in a circular band surrounding it of some radius, and then additional positive points lie outside that radius. Argue that the kernel $K(\underline{x}, \underline{y}) = (1 + \underline{x} \cdot \underline{y})^4$ will recover a separator.
- 4) Consider the XOR data (located at $(\pm 1, \pm 1)$). Express the dual SVM problem and find a separator using
 - $K(\underline{x}, \underline{y}) = (1 + \underline{x} \cdot \underline{y})^2$
 - $K(\underline{x}, \underline{y}) = \exp(-\|\underline{x} - \underline{y}\|^2)$.For each, determine the regions of (x_1, x_2) space where points will be classified as positive or negative. Given that each produces a distinct separator, how might you decide which of the two was preferred?

1) a)



Above diagram represents data being separated by the non-linear circular separator.

equation of a circle is

$$x^2 + y^2 = r^2$$

Changing x to u_1 and y to u_2 to
avoid confusion with sign and y

$$u_1^2 + u_2^2 = r^2$$

$$x_1^2 + x_2^2 = \gamma^2$$

Let

$$f(u) = x_1^2 + x_2^2 - \gamma^2$$

if this is < 0 , point lies inside the circle. Otherwise point lies outside the circle.

So, we can define our separator as

$$x_1^2 + x_2^2 - \gamma^2 < 0 \quad +1$$

$$x_1^2 + x_2^2 - \gamma^2 > 0 \quad -1$$

Our sign function is the above function which is circle inequation.

Here we needed x_1^2 and x_2^2 order features which will recover a separator

Both of these can be found in the feature set $(1, u_1, x_2, x_1, x_2, x_1^2, x_2^2)$.

So, a feature map in the above format will recover a circular separator

will recover a circular separator.

b)

An ellipse is represented by following equation

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} = 1$$

$$b^2 x^2 + a^2 y^2 = a^2 b^2$$

switching variables again

$$b^2 x_1^2 + a^2 y_2^2 = a^2 b^2$$

if we have to shift the origin, equation becomes

$$b^2 (x_1 - z_1)^2 + a^2 (x_2 - z_2)^2 = a^2 b^2$$

Here we can write the sign function as

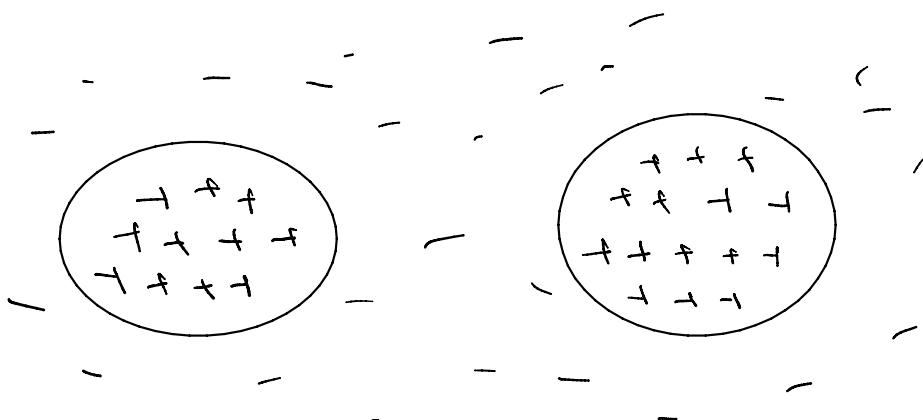
$$b^2 (x_1 - z_1)^2 + a^2 (x_2 - z_2)^2 - a^2 b^2 < 0 \rightarrow +1$$

$$b^2(x_1 - z_1)^2 + a^2(x_2 - z_2)^2 - a^2 b^2 > 0 \rightarrow -1$$

This separator also has the highest order of x_1 and x_2 as 2 and no other terms other than its first order.

So $\{x_1, x_2, x_1 x_2, x_1^2, x_2^2\}$ feature space will recover if there is a separator as defined above.

2)



From the discussion board, it is

From the discussion board, it is understood that data is separated in the above format.

We have two ellipses, both containing positive data and negatives everywhere else

Let first ellipse be

$$b_1^2 (x_1 - z_1)^2 + a_1^2 (x_2 - z_2)^2 = a_1^2 b_1^2$$

This can be defined as

$$K_1 \Rightarrow b_1^2 (x_1 - z_1)^2 + a_1^2 (x_2 - z_2)^2 - a_1^2 b_1^2$$

Second be

$$b_2^2 (x_1 - z_3)^2 + a_2^2 (x_2 - z_4)^2 = a_2^2 b_2^2$$

$$K_2 \Rightarrow b_2^2 (x_1 - z_3)^2 + a_2^2 (x_2 - z_4)^2 - a_2^2 b_2^2$$

Let a function K be defined as product of K_1 and K_2

Same way we represent two circles next to each other by product of their equations, we define

$$K = K_1 \times K_2$$

Here we have 3 classes of points

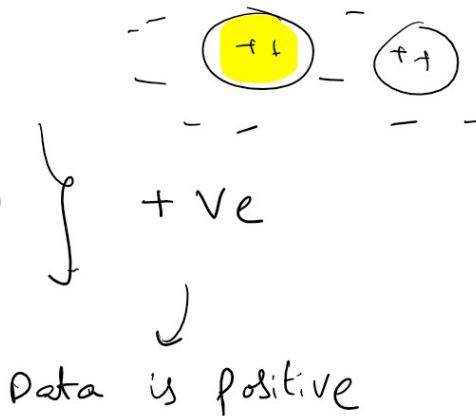
Class 1) Inside the first ellipse

Inequalities will be

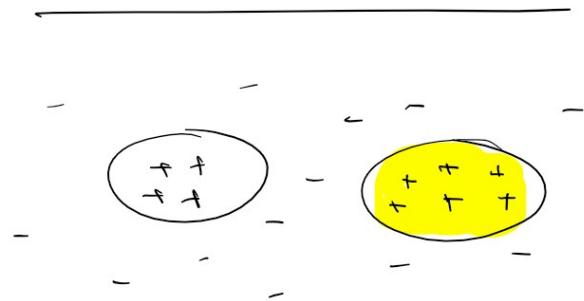
$$K_1 < 0$$

$$K_2 > 0$$

$$\Rightarrow K_1 K_2 < 0$$



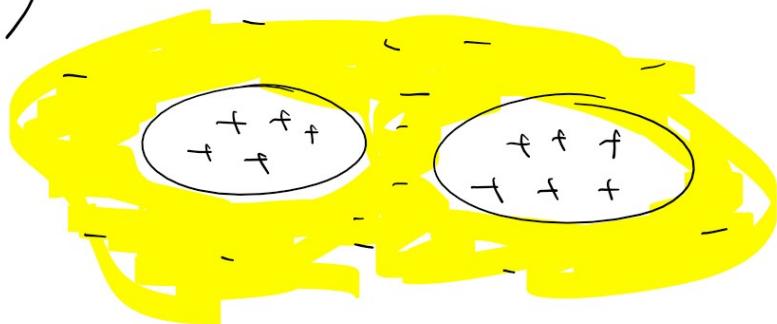
Class 2) Inside second ellipse



Same as inside first ellipse

$$\begin{aligned} K_1 > 0 \\ K_2 < 0 \end{aligned} \Rightarrow K_1, K_2 < 0 \quad \left. \begin{array}{l} \\ \end{array} \right\} +ve$$

Class 3)



Here

$$\begin{aligned} K_1 > 0 \\ K_2 > 0 \end{aligned} \Rightarrow K_1, K_2 > 0 \quad \left. \begin{array}{l} \\ \end{array} \right\} -ve$$

Points are -ve

So, our separator function can just be
the below

$$\begin{array}{ll} K_1, K_2 < 0 & +ve \\ K_1, K_2 > 0 & -ve \end{array}$$

K_1, K_2 will look like

$$\left(b_1^2 (x_1^2 - z_1^2) + a_1^2 (x_2^2 - z_2^2) - a_1^2 b_1^2 \right) \times \left(b_2^2 (x_1^2 - z_3^2) + a_2^2 (x_2^2 - z_4^2) - a_2^2 b_2^2 \right)$$

This raises to powers of 4 and variables
will be

$$\left\{ x_1^4, x_2^4, x_1^2 x_2^2 \right\}$$

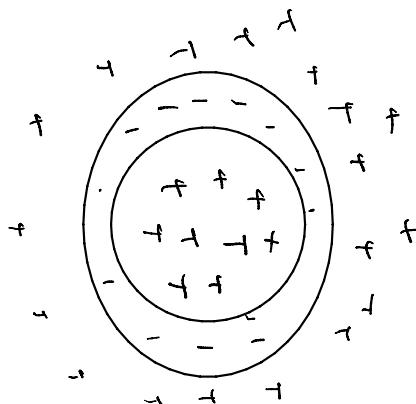
All these features can be recovered
from Kernel with dimension 4.

from Kernel with dimension 4^p .

So,

$(1 + \underline{x} \cdot \underline{y})^4$ will recover a separator in the given case.

3)



Again, two inequations of circle

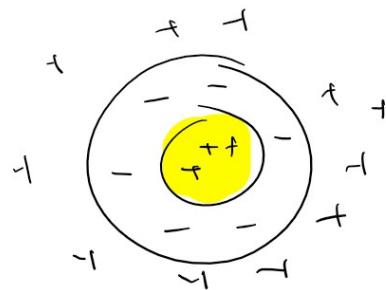
$$K_1 \Rightarrow (x_1 - a_1)^2 + (x_2 - b_1)^2 - r_1^2 \rightarrow \text{Inner circle}$$

$$K_2 \Rightarrow (x_1 - a_2)^2 + (x_2 - b_2)^2 - r_2^2 \rightarrow \text{Outer circle}$$

We have 3 classes of points again

We have 3 classes of points again

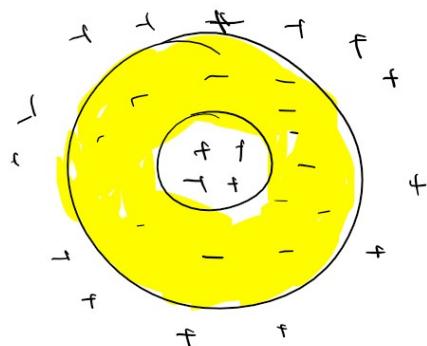
1) Inside inner circle



Point is inside both circles, so

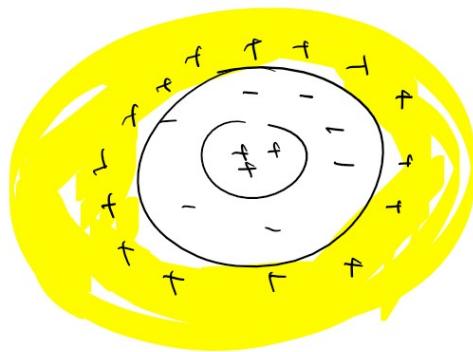
$$\begin{aligned} \kappa_1 &< 0 \\ \kappa_2 &< 0 \end{aligned} \quad \left. \begin{array}{l} \kappa_1, \kappa_2 > 0 \Rightarrow +ve \\ \text{Data is +ve} \end{array} \right\}$$

2) Inside outer circle, outside inner circle



$$\begin{array}{l} \kappa_1 > 0 \\ \kappa_2 < 0 \end{array} \quad \left. \begin{array}{l} \kappa_1 \kappa_2 < 0 \Rightarrow -ve \end{array} \right\}$$

3) outside both circles,



$$\begin{array}{l} \kappa_1 > 0 \\ \kappa_2 > 0 \end{array} \quad \left. \begin{array}{l} \kappa_1 \kappa_2 > 0 \end{array} \right\} +ve$$

So our classifier can be

$$\kappa = \kappa_1 \kappa_2 \quad \left. \begin{array}{l} > 0 \\ < 0 \end{array} \right\} \begin{array}{l} +ve \\ -ve \end{array}$$

κ can be expressed as

K can be expressed as

$K \Rightarrow$

$$\left((x_1 - a_1)^2 + (x_2 - b_1)^2 - \gamma_1^2 \right) \times \\ \left((x_1 - a_2)^2 + (x_2 - b_2)^2 - \gamma_2^2 \right)$$

This has

$x_1^4, x_2^4, x_1^2 x_2^2$ terms in it.

All of these can be retrieved using
Kernel

$$(1 + x \cdot y)^4$$

4)

- a) Written on paper. Will be attached along with this submission.

b) In a, we derived expanded form of the Kernel function as

$$\begin{aligned}
 & \text{Max}_{\alpha_1, \dots, \alpha_4} (\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4) - \frac{1}{2} \left[K_{11} y_1^2 \alpha_1^2 + K_{22} y_2^2 \alpha_2^2 \right. \\
 & \quad + K_{33} y_3^2 \alpha_3^2 + K_{44} y_4^2 \alpha_4^2 + 2 \left[y_1 y_2 \alpha_1 \alpha_2 K_{12} + \right. \\
 & \quad y_2 y_3 \alpha_2 \alpha_3 K_{23} + y_3 y_4 \alpha_3 \alpha_4 K_{34} + y_4 y_1 \alpha_4 \alpha_1 K_{41} \\
 & \quad \left. \left. + y_1 y_3 \alpha_1 \alpha_3 K_{13} + y_2 y_4 \alpha_2 \alpha_4 K_{24} \right] \right]
 \end{aligned}$$

Now we find Kernel values using

$$K(\underline{x}, \underline{y}) = \exp(-\|\underline{x} - \underline{y}\|^2)$$

$$\begin{aligned}
 K_{11} &= \exp\left(-\|(-1 - 1) - (-1 - 1)\|^2\right) \\
 &= \exp(0) = 1 \quad \text{Same for } K_{22}, K_{33}, K_{44}
 \end{aligned}$$

$$K_{12} = e^{-\rho} \left(- \left\| \begin{pmatrix} -1 & -1 \\ -1 & +1 \end{pmatrix} \right\|^2 \right)$$

$$= e^{-\rho} \left(- \left\| \begin{pmatrix} 0 & -2 \end{pmatrix} \right\|^2 \right)$$

$$\approx e^{-\rho} (-2^2) = e^{-4\rho}$$

$$K_{23} = e^{-\rho} \left(- \left\| \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix} \right\|^2 \right)$$

$$= e^{-\rho} \left(- \left\| \begin{pmatrix} -2 & 2 \end{pmatrix} \right\|^2 \right)$$

$$= e^{-8\rho}$$

$$K_{34} = e^{-\rho} \left(- \left\| \begin{pmatrix} +1 & -1 \\ 1 & 1 \end{pmatrix} \right\|^2 \right)$$

$$\approx e^{-4\rho}$$

$$K_{41} = e^{-\rho} \left(- \left\| \begin{pmatrix} 1 & 1 \\ -1 & -1 \end{pmatrix} \right\|^2 \right)$$

$$\approx e^{-8\rho}$$

$$K_{13} = e^{w\left(- \left\| \begin{pmatrix} 1 & -1 \\ 1 & -1 \end{pmatrix} - \begin{pmatrix} 1 & -1 \\ 1 & -1 \end{pmatrix} \right\|^2 \right)}$$

$$\approx e^{-4}$$

$$K_{24} = e^{w\left(- \left\| \begin{pmatrix} -1 & 1 \\ 1 & 1 \end{pmatrix} - \begin{pmatrix} 1 & -1 \\ 1 & -1 \end{pmatrix} \right\|^2 \right)}$$

$$\approx e^{-4}$$

Using these

$$\begin{aligned} \text{Max}_{\alpha_1, \dots, \alpha_4} & (\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4) - \frac{1}{2} \left[K_{11} y_1^2 \alpha_1^2 + K_{22} y_2^2 \alpha_2^2 + \right. \\ & \left. K_{33} y_3^2 \alpha_3^2 + K_{44} y_4^2 \alpha_4^2 + 2 \left[y_1 y_2 \alpha_1 \alpha_2 K_{12} + \right. \right. \\ & y_2 y_3 \alpha_2 \alpha_3 K_{23} + y_3 y_4 \alpha_3 \alpha_4 K_{34} + y_4 y_1 \alpha_4 \alpha_1 K_{41} \\ & \left. \left. + y_1 y_3 \alpha_1 \alpha_3 K_{13} + y_2 y_4 \alpha_2 \alpha_4 K_{24} \right] \right] \end{aligned}$$

$$\text{Min}_{\alpha_1, \dots, \alpha_4} \sum \alpha - \frac{1}{2} \left[1 (\sum \alpha^2) + 2 e^{-4} (-\alpha_1 \alpha_2 \right. \quad 77$$

$$\alpha_1 \sim \alpha_4 = -2 \left[\alpha_1^2 + \alpha_2^2 + \alpha_3^2 + \alpha_4^2 - (\alpha_3 \alpha_4 - \alpha_1 \alpha_2 - \alpha_2 \alpha_4) + e^{-8} (\alpha_2 \alpha_3 + \alpha_4 \alpha_1) \right]$$

d - 1/2((a^2 + b^2 + c^2 + d^2) + 2*(exp(-4))((-a*b) - (c*d) - (a*c) - (b*d)) + exp(-8)((b*c) + (d*a))), {a,b,c,d}]

Extended Keyboard Upload

Examples Random

Input interpretation:

maximize	$a + b + c + d - \frac{1}{2} ((a^2 + b^2 + c^2 + d^2) + 2 \exp(-4) (-a b - c d - a c - b d) + \exp(-8) (b c + d a))$	for	<table border="1"> <tr><td>a</td></tr> <tr><td>b</td></tr> <tr><td>c</td></tr> <tr><td>d</td></tr> </table>	a	b	c	d
a							
b							
c							
d							

Global maxima:

(no global maxima found)

Local maximum:

Exact form More digits

$$\max \left(a + b + c + d - \frac{1}{2} ((a^2 + b^2 + c^2 + d^2) + 2 \exp(-4) (-a b - c d - a c - b d) + \exp(-8) (b c + d a)) \right) \approx 2.07533 \text{ at } (a, b, c, d) \approx (1.03766, 1.03766, 1.03766, 1.03766)$$

Download Page

POWERED BY THE WOLFRAM LANGUAGE

Using Wolfram solver, we get α values

as

$$(\alpha_1, \alpha_2, \alpha_3, \alpha_4) \cong (1.037, 1.037, 1.037, 1.037)$$

Calculating bias

$$b = y^i - \underline{w} \cdot u$$

$$\begin{aligned}\underline{w} \cdot u_i &= \sum_{i=1}^n \alpha_i y^i K(u^i, u) \\ &= \alpha_1 y_1 K_{1u} + \alpha_2 y_2 K_{2u} + \alpha_3 y_3 K_{3u} + \alpha_4 y_4 K_{4u} \\ &= 1.037 \left[-K_{1u} + K_{2u} + K_{3u} - K_{4u} \right]\end{aligned}$$

Using $i = 1$

$$\approx 1.037 \left[-1 + e^{-4} + e^{-4} - e^{-8} \right]$$

$$\approx -0.9993 \approx 1$$

$$b = y_i - \underline{w}_i u_i = -1 - (-1)$$

$$= 0$$

So separator is $= 0$

$$1.038[-K_{1n} + K_{2n} + K_{3n} - K_{4n}] + 0$$

where

$$K_{in} = \exp\left(-\|x_i - x\|^2\right)$$



2.4) XOR data is given by

	①	②	y
x_1	-1	-1	-1
x_2	-1	+1	1
x_3	1	-1	1
x_4	1	1	-1
↓ data points			

Kernel function for SVM is given by

$$\text{Max } \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

$$\text{s.t. } \sum_{i=1}^m \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

Substituting,

for
first we find Kernel values

$$K \text{ is given by } (1 + x^i \cdot x^j)^2$$

$$\begin{aligned}
 & \text{Max } (\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4) - \frac{1}{2} \left[K_{11} y_1^2 \alpha_1^2 + K_{22} y_2^2 \alpha_2^2 + K_{33} y_3^2 \alpha_3^2 \right. \\
 & \quad \left. + K_{44} y_4^2 \alpha_4^2 + 2y_1 y_2 \alpha_1 \alpha_2 K_{12} + 2y_1 y_3 \alpha_1 \alpha_3 K_{13} + 2y_1 y_4 \alpha_1 \alpha_4 K_{14} \right. \\
 & \quad \left. + 2y_2 y_3 \alpha_2 \alpha_3 K_{23} + 2y_2 y_4 \alpha_2 \alpha_4 K_{24} \right. \\
 & \quad \left. + 2y_3 y_4 \alpha_3 \alpha_4 K_{34} \right]
 \end{aligned}$$

To find values,

$$\begin{aligned}
 K_{11} &= \left(1 + \underline{\alpha_1 \alpha_1} \right)^2 = \left(1 + \begin{bmatrix} -1 \\ -1 \end{bmatrix} \cdot \begin{bmatrix} -1 & -1 \end{bmatrix} \right)^2 \\
 &= (1+2)^2 = 9
 \end{aligned}$$

for all K_{22}, K_{33}, K_{44} ; dot product w.r.t. result in 2, thus $\rightarrow 9$

$$\begin{aligned}
 K_{12} &= 1 + \left(\begin{bmatrix} -1 \\ -1 \end{bmatrix} \cdot \begin{bmatrix} -1 & 1 \end{bmatrix} \right)^2 = 1 + 1 + (-1) = (1+0)^2 = 1 \\
 K_{23} &= 1 + \left(\begin{bmatrix} -1 \\ -1 \end{bmatrix} \cdot \begin{bmatrix} 1 & -1 \end{bmatrix} \right)^2 = (1-1-1)^2 = (-1)^2 = 1
 \end{aligned}$$

$$\begin{aligned}
 K_{14} &= \left(1 + \left(\begin{bmatrix} -1 \\ -1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 \end{bmatrix} \right) \right)^2 = (1+(-1)-1)^2 \\
 &= 1
 \end{aligned}$$

\therefore all will be 1

and $y_1 = -1$ $y_2 = 1$ $y_3 = 1$ $y_4 = -1$
 So, substituting in equation gives \downarrow
 all squares are 1

$$\max_{\alpha_1, \alpha_2, \alpha_3, \alpha_4} (\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4) - \frac{1}{2} [9\alpha_1^2 + 9\alpha_2^2 + 9\alpha_3^2 + 9\alpha_4^2 + 2(-\alpha_1\alpha_2 + \alpha_2\alpha_3 - \alpha_3\alpha_4 + \alpha_4\alpha_1 - \alpha_1\alpha_3 - \alpha_2\alpha_4)]$$

Now we have a maximization function with $\alpha_1, \alpha_2, \alpha_3, \alpha_4$

$$\max_{\alpha_1, \alpha_2, \alpha_3, \alpha_4} \left(\sum_{i=1}^4 \alpha_i \right) - \frac{9}{2} [\alpha_1^2 + \alpha_2^2 + \alpha_3^2 + \alpha_4^2 + 2(-\alpha_1\alpha_2 + \alpha_2\alpha_3 - \alpha_3\alpha_4 + \alpha_4\alpha_1 - \alpha_1\alpha_3 - \alpha_2\alpha_4)]$$

Substituting this in Wolfram-Solver, I got

Max is at $\frac{4}{17}, \frac{1}{17}$ and

$$(\alpha_1, \alpha_2, \alpha_3, \alpha_4) = \left(\frac{2}{17}, \frac{2}{17}, \frac{2}{17}, \frac{2}{17} \right)$$

Using this values, we make our sign function.

$\underline{w} \cdot \underline{x}$ is given by

$$\sum_{i=1}^m \alpha_i y^i K(x^i, x)$$

Substituting values

$$\Rightarrow \alpha_1 y_1 K_{1n} + \alpha_2 y_2 K_{2n} + \alpha_3 y_3 K_{3n} + \alpha_4 y_4 K_{4n}$$

$$= \frac{1}{8} \left[-K_{1n} + K_{2n} + K_{3n} - K_{4n} \right]$$

finding bias

$$b = y^i - \underline{w} \cdot \underline{x}$$

Using $i=1$

$$\underline{w} \cdot \underline{x}_1 = \frac{1}{8} [-K_{11} + K_{21} + K_{31} - K_{41}]$$

$$= \frac{1}{8} [-9 + 1 + 1 - 1]$$

$$= \frac{1}{8} (-8) = -1$$

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

bias is 0

So, separator is

$$\Rightarrow \frac{1}{8} [-K_{1n} + K_{2n} + K_{3n} - K_{4n}] + 0$$

$$\text{here } \underline{w} \text{ is } (\underline{1} + \underline{x_i} \cdot \underline{x})^2$$