

CS 484: Introduction to Machine Learning

Spring 2021 Assignment 4

Question 1 (40 Points)

In 2014, Allstate provided the data on Kaggle.com for the Allstate Purchase Prediction Challenge. The data contain transaction history for customers that ended up purchasing a policy. For each Customer ID, we know the quote history and the purchased coverage options.

The data is available on the Blackboard as **Purchase_Likelihood.csv**.

1. It contains 665,249 observations on 97,009 unique Customer ID.
2. The nominal target variable is **insurance** that has these categories 0, 1, and 2
3. The nominal features are (categories are inside the parentheses):
 - a. **group_size**. *How many people covered under the policy (1, 2, 3, or 4)?*
 - b. **homeowner**. *Whether the customer owns a home or not (0 = No, 1 = Yes)?*
 - c. **married_couple**. *Does the customer group contain a married couple (0 = No, 1 = Yes)?*

You will train a Naïve Bayes model without any smoothing using all the observations in the **Purchase_Likelihood.csv**. In other words, the Laplace/Lidstone alpha is zero. Please answer the following questions based on your model.

- a) (5 points) Show in a table the frequency counts and the Class Probabilities of the target variable.

insurance	0	1	2
Frequency Count	143691	426067	95491
Class Probability	0.2159958	0.6404624	0.1435417

- b) (5 points) Show the crosstabulation table of the target variable by the feature **group_size**. The table contains the frequency counts.

group_size	insurance		
	0	1	2
1	115460	329552	74293
2	25728	91065	19600
3	2282	5069	1505
4	221	381	93

- c) (5 points) Show the crosstabulation table of the target variable by the feature **homeowner**. The table contains the frequency counts.

homeowner	insurance		
	0	1	2
0	78659	183130	46734
1	65032	242937	48757

- d) (5 points) Show the crosstabulation table of the target variable by the feature **married_couple**. The table contains the frequency counts.

Married_couple	insurance		
	0	1	2
0	117110	333272	75310
1	26581	92795	20181

- e) (5 points) Calculate the Cramer's V statistics for the above three crosstabulations tables. Based on the Cramer's V statistics, which feature has the strongest association with the target insurance?

Feature	Cramer's V
group_size	0.02710201
homeowner	0.09708642
married_couple	0.03242165

Homeowner has the strongest association with the target insurance.

- f) (10 points) For each of the sixteen possible value combinations of the three features, calculate the predicted probabilities for insurance = 0, 1, 2 based on the Naïve Bayes model that includes features group_size, homeowner, and married_couple. List your answers in a table with proper labeling.

group_size	homeowner	married_couple	Prob(insurance = 0)	Prob(insurance = 1)	Prob(insurance = 2)
1	0	0	0.227037	0.627593	0.145370
1	0	1	0.214391	0.637467	0.148142
1	1	0	0.205588	0.654128	0.140284
1	1	1	0.193842	0.663414	0.142744
2	0	0	0.238441	0.614462	0.147097
2	0	1	0.225342	0.624635	0.150024
2	1	0	0.216281	0.641528	0.142192
2	1	1	0.204079	0.651128	0.144794
3	0	0	0.250201	0.601084	0.148715
3	0	1	0.236653	0.611546	0.151801
3	1	0	0.227342	0.628652	0.144006
3	1	1	0.214684	0.638559	0.146756
4	0	0	0.262308	0.587475	0.150218
4	0	1	0.248318	0.598215	0.153467
4	1	0	0.238767	0.615513	0.145720
4	1	1	0.225656	0.625720	0.148624

- g) (5 points) Based on your model, determine the value combination of group_size, homeowner, and married_couple that will yield the maximum value for this odds $\text{Prob}(\text{insurance} = 1) / \text{Prob}(\text{insurance} = 2)$? What is that maximum odds value?

group_size	homeowner	married_couple	odd value(p_in_1/p_in_0)
1	0	0	4.317213
1	0	1	4.303069
1	1	0	4.662868
1	1	1	4.647591
2	0	0	4.177260
2	0	1	4.163574
2	1	0	4.511710
2	1	1	4.496928
3	0	0	4.041845
3	0	1	4.028602
3	1	0	4.365452
3	1	1	4.351149
4	0	0	3.910818
4	0	1	3.898005
4	1	0	4.223935
4	1	1	4.210096

group_size	1.0000000
homeowner	1.0000000
married_couple	0.0000000
p_in_0	0.205588
p_in_1	0.654128
p_in_2	0.140284
odd value(p_in_1/p_in_2)	4.662868

The maximize odds value of $\text{Prob}(\text{insurance} = 1) / \text{Prob}(\text{insurance} = 2)$ is 4.662868 at value combination group_size=1, homeowner=1 & married_couple=0.

The maximum odd value is 4.662868

Question 2 (60 points)

The **SpiralWithCluster.csv** contains four variables.

Name	Description	Measurement Level	Role
Id	Case Identifier	Nominal	Identifier
X	x-coordinate	Interval	Feature
Y	y-coordinate	Interval	Feature
SpectralCluster	Cluster Identifier	Binary	Target

Please use the Support Vector Machine (SVM) algorithm to classify SpectralCluster. You will use the `sklearn.svm.SVC` function with the following specifications.

1. The linear kernel
2. The decision function shape is One Over Rest (OVR)
3. No limit on the number of iterations
4. The random seed is 20210325

Please answer the following questions based on your model.

- a) (5 points) What is the equation of the separating hyperplane in the Slope-Intercept form? Please state the coefficients up to seven decimal places.

Intercept = [0.003345]

Coefficients = [[0.0533351 0.3286838]]

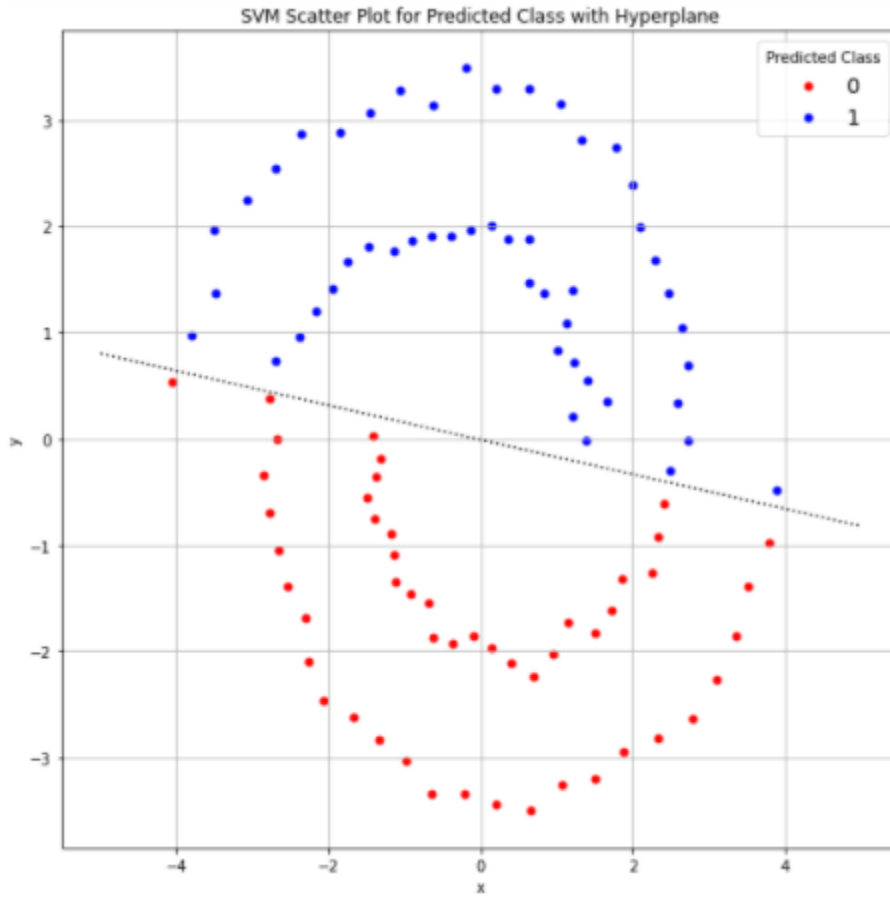
Equation of the separating hyperplane: $(0.003345) + (0.0533351 \cdot x) + (0.3286838 \cdot y) = 0$

Slope-Intercept Form: $y = -0.0101770 - 0.1622687x$

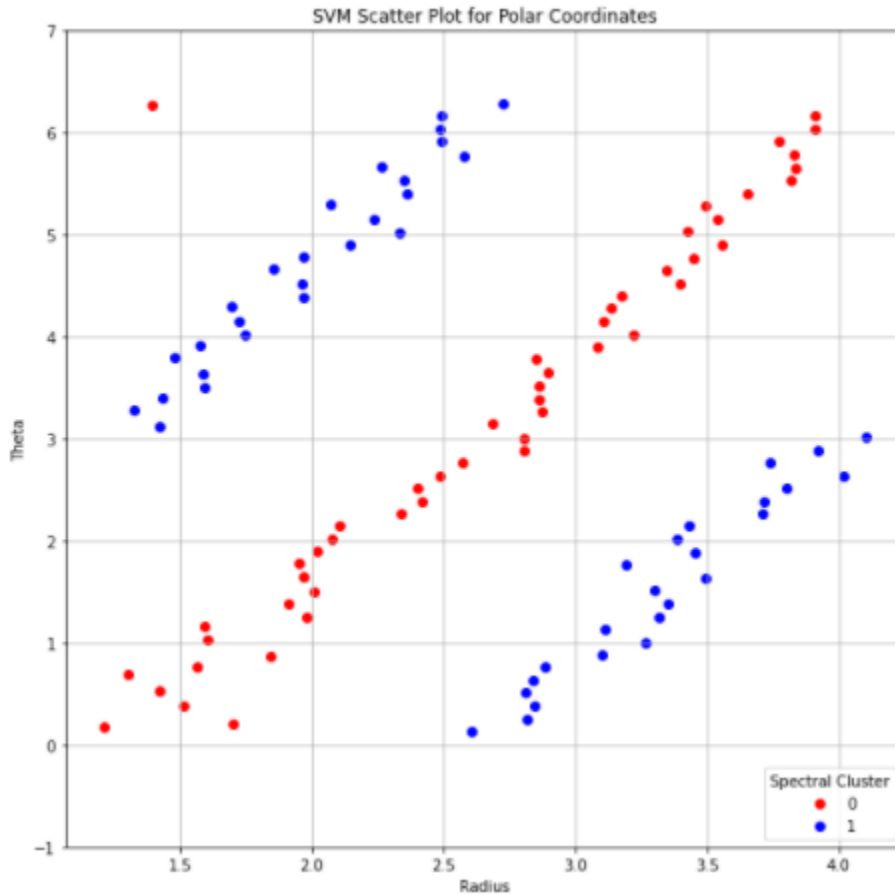
- b) (5 points) What is the misclassification rate?

Misclassification rate is 0.5

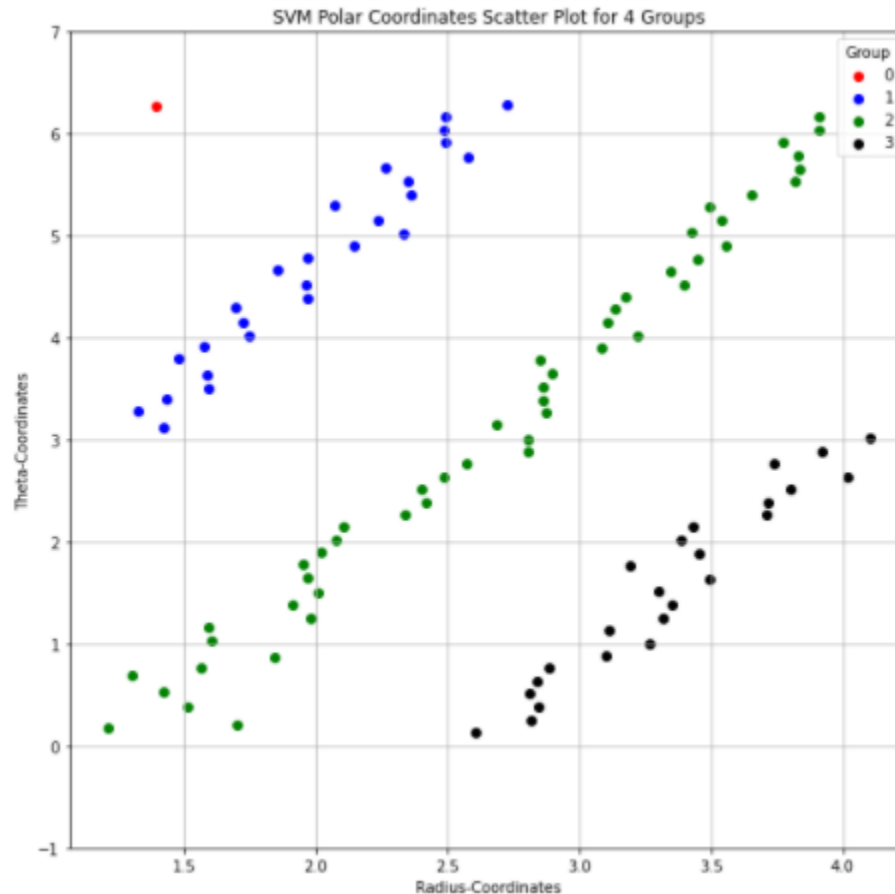
- c) (10 points) Please plot the y-coordinate against the x-coordinate in a scatterplot. Please color-code the points using the predicted SpectralCluster (0 = Red and 1 = Blue). Besides, plot the hyperplane as a dotted line to the graph. To obtain the full credits, you should properly label the axes, the legend, and the chart title. Also, please add grid lines to the axes.



- d) (10 points) Please express the data in polar coordinates. Please plot the theta-coordinate against the radius-coordinate in a scatterplot. Please color-code the points using the SpectralCluster variable (0 = Red and 1 = Blue). To obtain the full credits, you should properly label the axes, the legend, and the chart title. Also, please add grid lines to the axes.



the full credits, you should properly label the axes, the legend, and the chart title. Also, grid lines should be added to the axes.



- f) (10 points) Since the graph in (e) has four separate but neighboring segments, we will apply the Support Vector Machine algorithm differently. Instead of applying SVM once on a multi-class target variable, you will SVM three times, each on a pair of groups.

SVM 0: Group 0 versus Group 1

SVM 1: Group 1 versus Group 2

SVM 2: Group 2 versus Group 3

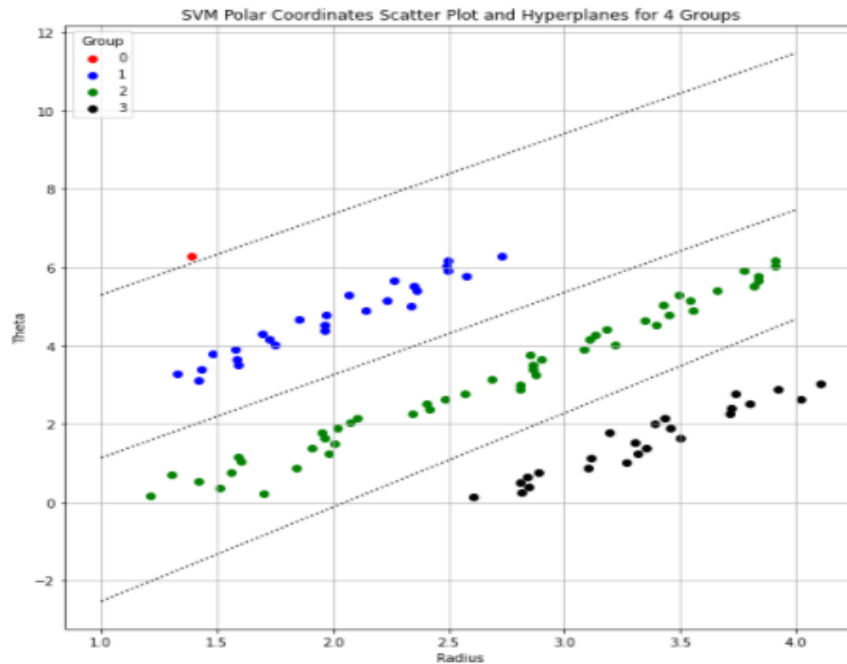
Please give the equations of the three hyperplanes.

Equation of the separating hyperplane for SVM 0: $(1.4691251) + (0.9337841 * x) + (-0.4538025 * y) = 0$

Equation of the separating hyperplane for SVM 1: $(-0.8768943) + (1.8920953 * x) + (-0.8961325 * y) = 0$

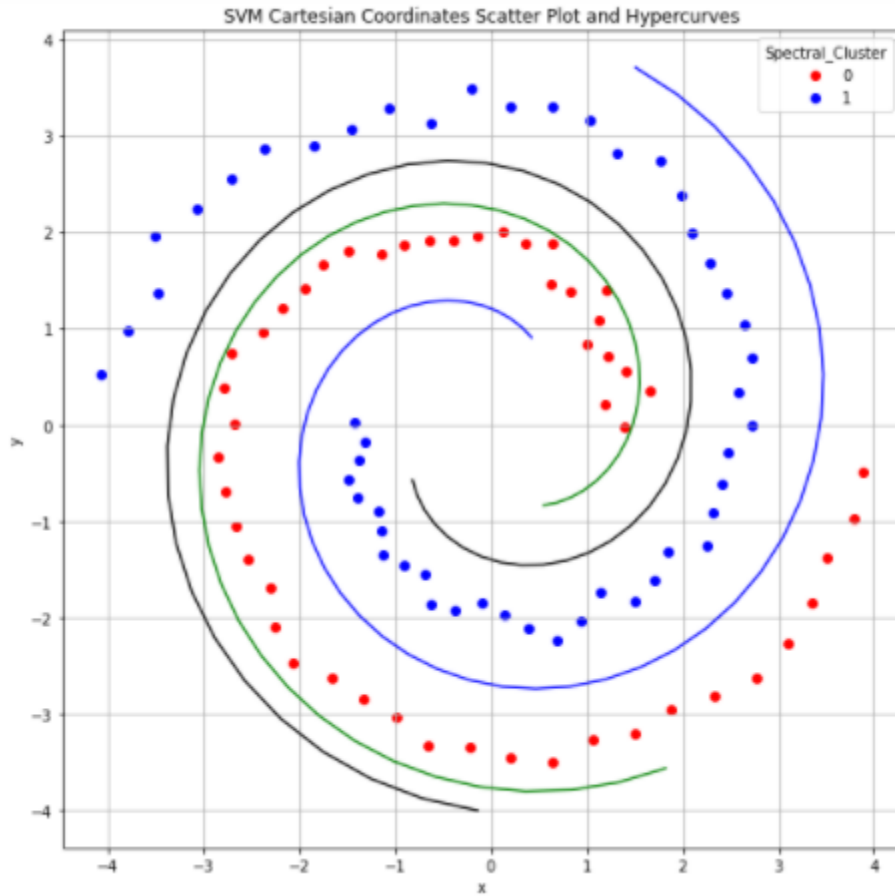
Equation of the separating hyperplane for SVM 2: $(-4.1328449) + (2.0125835 * x) + (-0.8375616 * y) = 0$

- g) (5 points) Please plot the theta-coordinate against the radius-coordinate in a scatterplot. Please color-code the points using the new Group target variable (0 = Red, 1 = Blue, 2 = Green, 3 = Black). Please add the hyperplanes to the graph. To obtain the full credits, you should properly label the axes, the legend, and the chart title. Also, grid lines should be added to the axes.



- h) (5 points) Convert the observations along with the hyperplanes from the polar coordinates back to the Cartesian coordinates. Please plot the y-coordinate against the x-coordinate in a scatterplot. Please color-code the points using the SpectralCluster (0 = Red and 1 = Blue). Besides, plot the hyper-curves as dotted lines to the graph. To obtain the full credits, you should properly label the axes, the legend, and the chart title. Also, grid lines should be added to the axes.

Based on your graph, which hypercurve do you think is not needed?



Based on the graph the green Curve is not needed. It doesn't classify the classes correctly.