

UC San Diego Extension Cloud Services for Machine Learning

Summer 2020

Homework#3

Date Given: July 13, 2020

Due Date: July 19, 2020

=====
Classification Using GCP: **There are 2 problems in this assignment.**
=====

The file “HW03 Jobs Data.csv” contains tabular student’s job data. There are 4 predictor variables.

- Height in inches
- Major
- Points (on a scale from 1 to 100)
- Weight in lbs

There is one response variable.

- Job (Categorical): Yes or No

This data was artificially generated using a script which ran on Google’s Spreadsheet application “Sheets”. The source code of the Sheets script is displayed at the end (Appendix) of this assignment.

The relationship between ‘Job’ and ‘Major’ + ‘Points’ is as follows.

If (Major = CompSci or Medicine or Finance) AND Points > 50
Job = Yes

This relationship between ‘Job’ and “Major’ + ‘Points’ is non-linear. The ‘Height’ and the ‘Weight’ data are purely noise. There are 1,500 observations. The following table displays the first 20 observations.

	A	B	C	D	E	F
1	Random Number	Height	Major	Points	Weight	Job = Yes/No
2	0.545174171	67	Finance	68	199	Yes
3	0.123945224	52	ElecEngg	38	107	No
4	0.089229406	51	ElecEngg	36	99	No
5	0.821472747	77	MachLearn	87	260	No
6	0.513495121	66	Finance	65	192	Yes
7	0.324542285	59	CompSci	52	151	Yes
8	0.869788654	79	ChemEngg	90	271	No
9	0.817844911	77	MachLearn	87	259	No
10	0.059904576	50	ElecEngg	34	93	No
11	0.775520336	75	MachLearn	84	250	No
12	0.137275209	52	ElecEngg	39	110	No
13	0.064974899	50	ElecEngg	34	94	No
14	0.223490243	56	CompSci	45	129	No
15	0.771401754	75	MachLearn	83	249	No
16	0.702930964	73	MachLearn	79	234	No
17	0.93211503	81	ChemEngg	95	285	No
18	0.794962	76	MachLearn	85	254	No
19	0.842757815	78	ChemEngg	88	265	No
20	0.242065207	56	CompSci	46	133	No

Problem#1

Build a **classification** Machine Learning model (Neural Network) using Google Cloud Platform (GCP) with the data in the 'HW03 Jobs Data.csv' file. Ignore the first column 'Random Number' while building the model.

This model does not know the TRUE relationship between 'Job' and 'Major' + 'Points'. It must learn this relationship only by analyzing the data. To challenge the ML model, noise data of 'Height' and 'Weight' columns have been added to the dataset.

The procedure to build a **classification** model on GCP is as follows (same as regression model).

1. GCP/Storage
 - a. Create a Bucket in GCP
 - b. Region: us-central1(Iowa)
 - c. Upload Data file in that bucket
2. GCP/Table/Dataset
 - a. Import data in a GCP Dataset from the bucket: Takes time
3. GCP/Table/Model
 - a. Train the Model
 - i. Select Target Variable + Budget: Takes time
 - b. Evaluate the Model
 - c. Test & Use: Deploy the Model: Takes time
 - i. Prediction

Evaluate the model after the model is trained on GCP. Which variable is the most important for prediction? Copy the 'Feature Importance' plot generated by GCP in your answer document. Does your model built on GCP able to capture the TRUE relationship between 'Job' and 'Major' + 'Points' variables?

Using the **classification** model, you have built on GCP, **classify** the job status of the following 4 students. Compute the probability of the response variable 'Job'. Assume the cut-off probability as 0.5. If the probability is greater than 0.5, classification of the 'Job' variable is 'Yes'.

	Height in inches	Major	Points	Weight in lbs	Logical Value of Job variable	Probability computed by GCP Job = Yes	Probability computed by GCP Job = No	Classification Yes/No
1	75	CompSci	85	220	Yes	?	?	?
2	82	CompSci	49	200	50/50	?	?	?
3	62	MachLearn	81	151	No	?	?	?
4	67	Finance	51	95	50/50	?	?	?

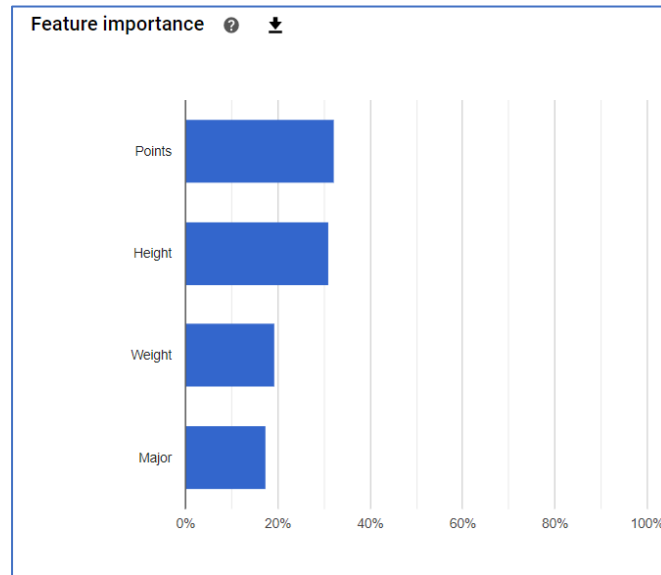
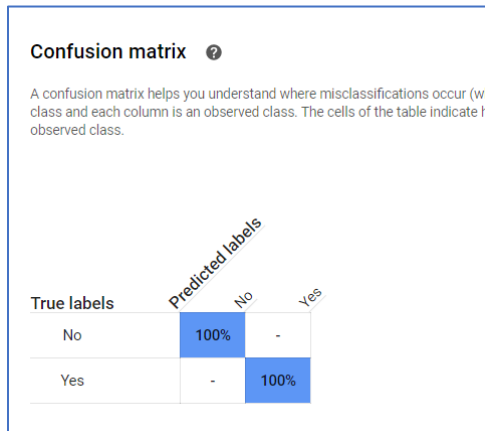
Problem#2: Can we build a kNN (k Nearest Neighbor) model for this dataset used in Problem#1? If no, why not?

=====

Building a classification model on GCP will cost a certain amount. Please check the GCP charges on your account before and after you complete this assignment. Make sure you do not deplete the \$300 credit you have on your account.

Problem#1

Model Evaluation



	A	B	C
1	Column name	Feature importance	
2	Points	0.3211139	
3	Height	0.3109972	
4	Weight	0.19298065	
5	Major	0.17490573	
6			

The TRUE relationship between 'Job' and 'Major' + 'Points' is as follows.

If (Major = CompSci or Medicine or Finance) AND Points > 50
Job = Yes

This relationship between 'Job' and "Major" + 'Points' is non-linear. The 'Height' and the 'Weight' data are purely noise.

The Neural Network model was NOT able to capture the TRUE relationship between 'Job' and 'Major' + 'Points'. The Feature Importance 'weight' associated with 'Weight' & 'Height' predictor variables should be close to zero.

The reason why this model failed to capture TRUE relationship is because we did not feed enough data to the model. We only have 1,500 observations. We need to repeat this experiment with more than 1,500 (close to 10,000) observations. Then only Neural Network would be able to predict correctly.

	Height in inches	Major	Points	Weight in lbs	Logical Value of Job variable	Probability computed by GCP Job = Yes	Probability computed by GCP Job = No	Classification Yes/No
1	75	CompSci	85	220	Yes	28.1	71.9	No
2	82	CompSci	49	200	50/50	32.3	67.7	No
3	62	MachLearn	81	151	No	63.2	36.8	Yes
4	67	Finance	51	95	50/50	85.1	14.9	Yes

Predict label
Job

Prediction result
No
Yes
Confidence score: 0.719
Confidence score: 0.281

Feature column name	Column ID	Data type	Status ↓	Value	Local feature import
Height	4482481307535802368	Numeric	Required	75	
Major	1239889575829045248	Categorical	Required	CompSci	
Points	9094167325963190272	Numeric	Required	85	
Weight	5851575594256433152	Numeric	Required	220	

☐ Generate feature importance

PREDICT RESET

Predict label
Job

Prediction result
No
Yes
Confidence score: 0.677
Confidence score: 0.323

Feature column name	Column ID	Data type	Status ↓	Value	Local feature import
Height	4482481307535802368	Numeric	Required	82	
Major	1239889575829045248	Categorical	Required	CompSci	
Points	9094167325963190272	Numeric	Required	49	
Weight	5851575594256433152	Numeric	Required	200	

☐ Generate feature importance

PREDICT RESET

Predict label
Job

Prediction result
No
Confidence score: 0.368
Yes
Confidence score: 0.632

Feature column name	Column ID	Data type	Status ↓	Value
Height	4482481307535802368	Numeric	Required	<input type="text" value="62"/>
Major	1239889575829045248	Categorical	Required	<input type="text" value="MachLearn"/>
Points	9094167325963190272	Numeric	Required	<input type="text" value="81"/>
Weight	5851575594256433152	Numeric	Required	<input type="text" value="151"/>

☐ Generate feature importance

PREDICT **RESET**

Predict label
Job

Prediction result
No
Confidence score: 0.149
Yes
Confidence score: 0.851

Feature column name	Column ID	Data type	Status ↓	Value
Height	4482481307535802368	Numeric	Required	<input type="text" value="67"/>
Major	1239889575829045248	Categorical	Required	<input type="text" value="Finance"/>
Points	9094167325963190272	Numeric	Required	<input type="text" value="51"/>
Weight	5851575594256433152	Numeric	Required	<input type="text" value="95"/>

☐ Generate feature importance

PREDICT **RESET**

=====

Problem#2: Can we build a kNN (k Nearest Neighbor) model for this dataset used in Problem#1? If no, why not?

=====

The kNN modeling technique works only when the predictor variables are numerical, and the response variable is categorical.

Data given in this assignment contains categorical predictor variable. There are the 4 predictor variables in this dataset.

- Height in inches (numerical)
- **Major (categorical)**
- Points (on a scale from 1 to 100) (numerical)
- Weight in lbs (numerical)

Since the 'Major' variable is categorical, we cannot build kNN model on this data.

Only way to predict the class variable 'Job' (Categorical) is build a Neural Network with 3 numerical and 1 categorical variable.