

✓ Big Data Analytics Homework 04

Complete this assignment in Google Colab. Prior to submitting a copy of this notebook (.ipynb format), run every cell and ensure you have corrected all runtime errors. Be sure to fill in your Name and SUID in the following cell. As always, you must do your own work. This means you may not use answers to the following questions generated by any other person or a generative AI tool such as ChatGPT. You may, however, discuss this assignment with others in a general way and seek help when you need it, but, again, you must do your own work.

Name:

SUID:

✓ Setup

```
1 ! pip install pyspark -q

1 from pyspark.sql import SparkSession
2 from pyspark import SparkContext
3
4 sc = SparkContext.getOrCreate()
5
6 spark = SparkSession\
7     .builder\
8     .appName('Homework 04')\
9     .getOrCreate()
```

This assignment uses a data set containing information about data science programs at universities worldwide.

The dataset contains many columns that we can use to understand how these data science programs differ from one another.

```
1 # download the data science programs data set
2 %%bash
3 if [[ ! -f colleges-data-science-programs.csv ]]; then
4     wget https://syr-bda.s3.us-east-2.amazonaws.com/colleges-data-science-programs.csv -q
5 fi
```

✓ Q1

Read colleges-data-science-programs.csv into a Spark data frame named raw_ds_programs_text.

```
1 # your code here
2 raw_ds_programs_text = spark.read.csv("colleges-data-science-programs.csv", header=True, inferSchema=True)
3
```

```
1 # do not modify
2 print('rows: ', raw_ds_programs_text.count(),
3       ', cols:', len(raw_ds_programs_text.columns))
4
5 raw_ds_programs_text\
6     .show(5)
```

```
rows: 222 , cols: 28
```

id	name	url	program	degree	country	state	online	oncampus	department
1	South Dakota Stat...	http://www.sdstat...	Data Science	Masters	US	SD	false	true	Mathematics and S...
2	Dakota State Univ...	http://www.dsu.ed...	Analytics	Masters	US	SD	true	true	Business and Info...
3	Lewis University	http://www.lewisu...	Data Science	Masters	US	IL	true	true	Computer Science
4	Saint Joseph's Un...	http://online.sju...	Business Intellig...	Masters	US	PA	true	true	Business
5	University Of Leeds	http://www.engine...	Advanced Computer...	Masters	GB	NULL	false	true	Computer Science

only showing top 5 rows

✓ Q2

Starting with `raw_ds_programs_text`, create a new data frame named `ds_programs_text` which simply adds a column named `text` to the original data frame.

The `text` column will be a concatenation of the following columns, separated by a space: `program`, `degree`, and `department`. You will find the appropriate function in `pyspark.sql.functions`

An example of the `ds_programs_text_df` should give you:

```
ds_programs_text.orderBy('id').first().text
```

```
'Data Science Masters Mathematics and Statistics'
```

```
1 # your code here
2 from pyspark.sql.functions import concat_ws, col
3 ds_programs_text = raw_ds_programs_text.withColumn(
4     "text",
5     concat_ws(" ",
6         col("program"),
7         col("degree"),
8         col("department")
9     )
10 )
```

```
1 # do not modify
2 ds_programs_text.select('text')\
3     .show(5, truncate = False)
```

```

+-----+
|text|
+-----+
|Data Science Masters Mathematics and Statistics|
|Analytics Masters Business and Information Systems|
|Data Science Masters Computer Science|
|Business Intelligence & Analytics Masters Business|
|Advanced Computer Science(Data Analytics) Masters Computer Science|
+-----+
only showing top 5 rows

```

Q3

Create a pipeline named `pipe_features` that creates a new dataframe `ds_features`. The `pipe_features` pipeline should add a column, `features` to `ds_programs_text` that contains the `tfidf` of the `text` column.

Make sure to create your pipeline using methodology similar to what was demonstrated in class. Aside from removing stop words and setting a minimum token length of 2, no further restrictions should be imposed on the resulting vocabulary.

```
1 # your code here
2 from pyspark.ml import Pipeline
3 from pyspark.ml.feature import Tokenizer, StopWordsRemover, HashingTF, IDF
4
5 tokenizer = Tokenizer(inputCol="text", outputCol="words")
6 stop_words_remover = StopWordsRemover(inputCol="words", outputCol="filtered_words")
7 hashing_tf = HashingTF(inputCol="filtered_words", outputCol="raw_features", numFeatures=20)
8 idf = IDF(inputCol="raw_features", outputCol="features")
9
10 pipe_features = Pipeline(stages=[tokenizer, stop_words_remover, hashing_tf, idf])
11 ds_features = pipe_features.fit(ds_programs_text).transform(ds_programs_text)
12
13 ds_features_fitted = pipe_features.fit(ds_programs_text)
14 ds_features = ds_features_fitted.transform(ds_programs_text)
```

```
1 # do not modify
2 ds_features.select('features')\
3     .show(5,
4         truncate = False)
```

```

+-----+
|features|
+-----+

```

```
| (20, [3, 6, 13, 15, 16], [1.1305056524440635, 0.2481164722455897, 2.76811444184486, 0.5868902058550819, 2.188295946591918])
| (20, [6, 8, 9, 10, 14], [0.2481164722455897, 0.7158238892309751, 2.188295946591918, 1.5785303749710238, 0.853294879859578])
| (20, [3, 6, 11, 15], [2.261011304888127, 0.2481164722455897, 1.6229821375418576, 0.5868902058550819])
| (20, [6, 8, 12, 14], [0.2481164722455897, 1.431647784619501, 1.5359707605522277, 1.706589759719156])
| (20, [3, 6, 11, 14, 18, 19], [1.1305056524440635, 0.2481164722455897, 3.245964275083715, 0.853294879859578, 4.020877410340228, 2.699121570357909])
+-----+
only showing top 5 rows
```

Q4

Create a pipeline model called `pipe_pca` that computes the first two principle components of the `features` column created by `pipe_features`, and creates a new column named `scores`.

Use `pipe_pca` to create a data frame, `ds_features_1` with the columns `id`, `name`, `url`, and `scores`.

Note: Prior to computing PCA scores, you will want to scale the TF-IDF outputs. Refer to lecture notes regarding the appropriate parameters to use during this step.

```
1 # your code here
2 from pyspark.ml import Pipeline
3 from pyspark.ml.feature import StandardScaler, PCA
4
5 scaler = StandardScaler(inputCol="features", outputCol="scaled_features", withMean=True, withStd=True)
6 pca = PCA(k=2, inputCol="scaled_features", outputCol="scores")
7
8 pipe_pca = Pipeline(stages=[scaler, pca])
9 pipe_pca_model = pipe_pca.fit(ds_features)
10 ds_features_1 = pipe_pca_model.transform(ds_features).select("id", "name", "url", "scores")

1 # do not modify
2 ds_features_1\
3   .select('scores')\
4   .show(5,
5       truncate = False)
```

```
+-----+
|scores|
+-----+
| [-0.8273713018426965, -0.09636389661089356]|
| [1.541109169337008, 1.0772516193312067]|
| [-1.6352404416200834, -1.316722397188407]|
| [2.3459831485125564, -0.7393473954578522]|
| [-1.4596210033256862, 1.4993072818402193]|
+-----+
only showing top 5 rows
```

Q5

In this question you will write code that makes recommendations for programs closest to a program of interest.

Create a function named `get_nearest_programs` that returns the 3 closest programs to a program of interest.

The `get_nearest_programs` function should take 1 argument: `program_of_interest`. Write the function so that it returns the 3 programs (as defined by the `name` column) closest to the program argument as defined by Euclidian (L2) distance. Do not return the program of interest as one of the names.

Your function should **not** consider **Bachelors** programs.

```
1 # your code here
2 from pyspark.sql.functions import *
3 from pyspark.ml.linalg import Vectors
4 from pyspark.sql.functions import col
5 from pyspark.sql import DataFrame
6
7
8 def euclidean_distance(v1, v2):
9     return float(Vectors.squared_distance(v1, v2))**0.5
10
11 def get_nearest_programs(name: str, df=ds_features_1) -> list:
12     scores = df.filter(col('name') == name).select('scores').first()['scores']
```

```

13     return (df
14             .filter((col('degree') != 'Bachelors') & (col('name') != name))
15             .select('name', 'scores')
16             .rdd.map(lambda r: (r['name'], euclidean_distance(r['scores'], scores)))
17             .toDF(['name', 'distance'])
18             .orderBy('distance')
19             .limit(3)
20             .select('name')
21             .toPandas()['name'].tolist())

```

```

1 # do not modify
2 get_nearest_programs('Syracuse University')

```

```

↩ ['Columbia University',
   'New Jersey Institute of Technology',
   'Coventry University']

```

Q6

Create two Pandas dataframes `pc1` and `pc2` with the columns `word` and `absolute_loading` that contain the top 5 absolute values (descending order) of loadings.

```

1 # your code here
2 import pandas as pd
3 import numpy as np
4
5 def extract_principal_components(pipe_pca_model):
6     pca_model = pipe_pca_model.stages[-1]
7     pc_loadings = pca_model.pc.toArray()
8
9     result_df = pd.DataFrame({
10         'feature': range(len(pc_loadings)),
11         'load_pc1': pc_loadings[:, 0],
12         'load_pc2': pc_loadings[:, 1]
13     })
14     result_df['abs_pc1'] = np.abs(result_df['load_pc1'])
15     result_df['abs_pc2'] = np.abs(result_df['load_pc2'])
16     return result_df
17
18 def get_top_components(df, n=5):
19     pc1 = (df.sort_values('abs_pc1', ascending=False)
20           .head(n)[['feature', 'abs_pc1']]
21           .rename(columns={'abs_pc1': 'loading'}))
22     pc2 = (df.sort_values('abs_pc2', ascending=False)
23           .head(n)[['feature', 'abs_pc2']]
24           .rename(columns={'abs_pc2': 'loading'}))
25     return pc1, pc2
26
27 full_df = extract_principal_components(pipe_pca_model)
28 top_pc1, top_pc2 = get_top_components(full_df)
29

```

```

1 # do not modify
2 display(pc1.head())
3 display(pc2.head())

```



	word	absolute_loading	
0	test	0.408248	
1	another	0.408248	
2	important.	0.408248	
5	words	0.408248	
3	sentence	0.408248	

	word	absolute_loading	
6	sentence.	0.632456	
1	another	0.316228	
0	test	0.316228	
2	important.	0.316228	
3	sentence	0.316228	

Q7

Create a new pipeline called `pipe_pca_1` where you fit the maximum possible number of principal components for this dataset.

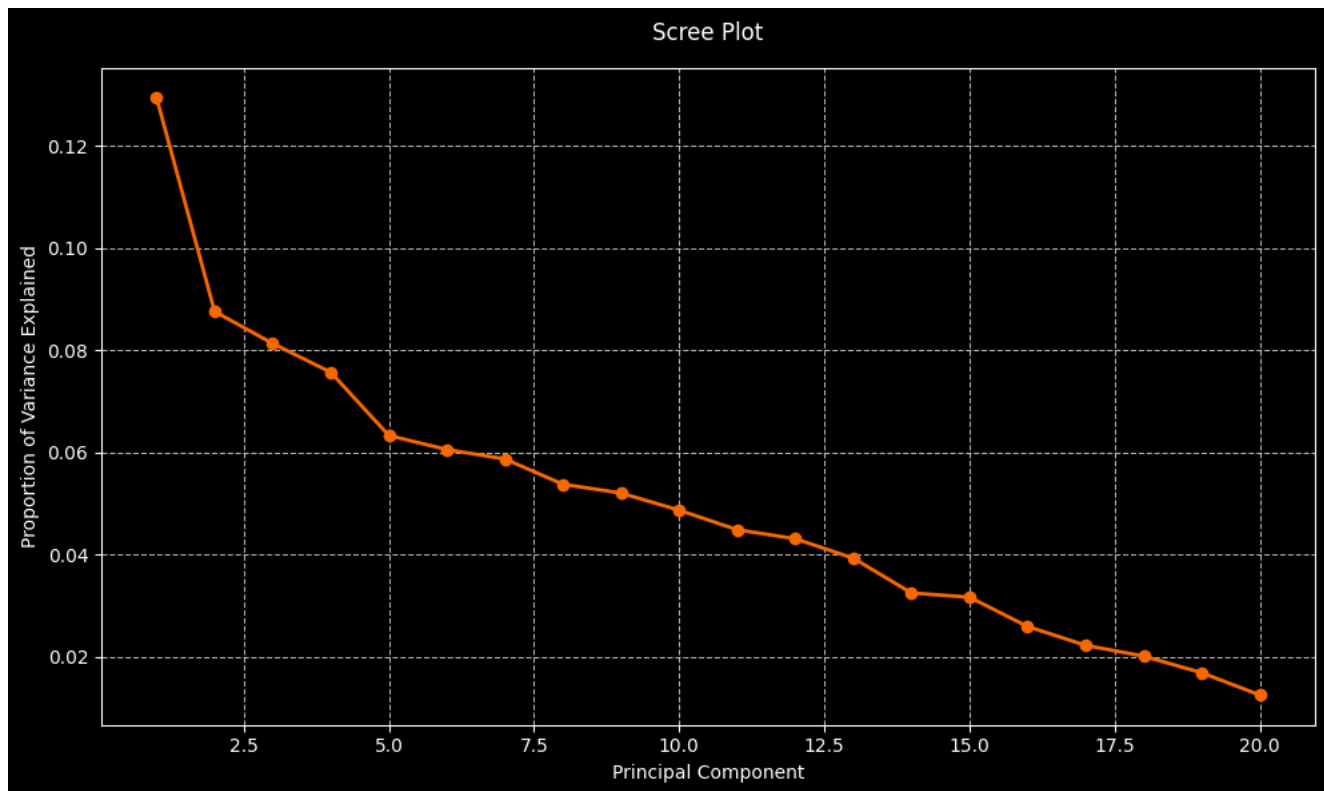
Create a scree plot and a plot of cumulative variance explained (exactly 2 plots).

Answer the following:

1. How many principal components were able to create (the maximum number)?
2. Based on either the scree or cumulative variance explained plot, how many principal components would you use if you were building a supervised machine learning model, and why?

```
1 # your code for new pipeline here
2 import matplotlib.pyplot as plt
3 import numpy as np
4 from pyspark.ml.feature import PCA
5
6 # Create a new pipeline model with the maximum number of principal components
7 max_pca = PCA(k=20, inputCol='scaled_features', outputCol='scores')
8 pipe_pca_1 = Pipeline(stages=[scaler, max_pca])
9 pipe_pca_model_1 = pipe_pca_1.fit(ds_features)
10 ds_features_max_pca = pipe_pca_model_1.transform(ds_features)
```

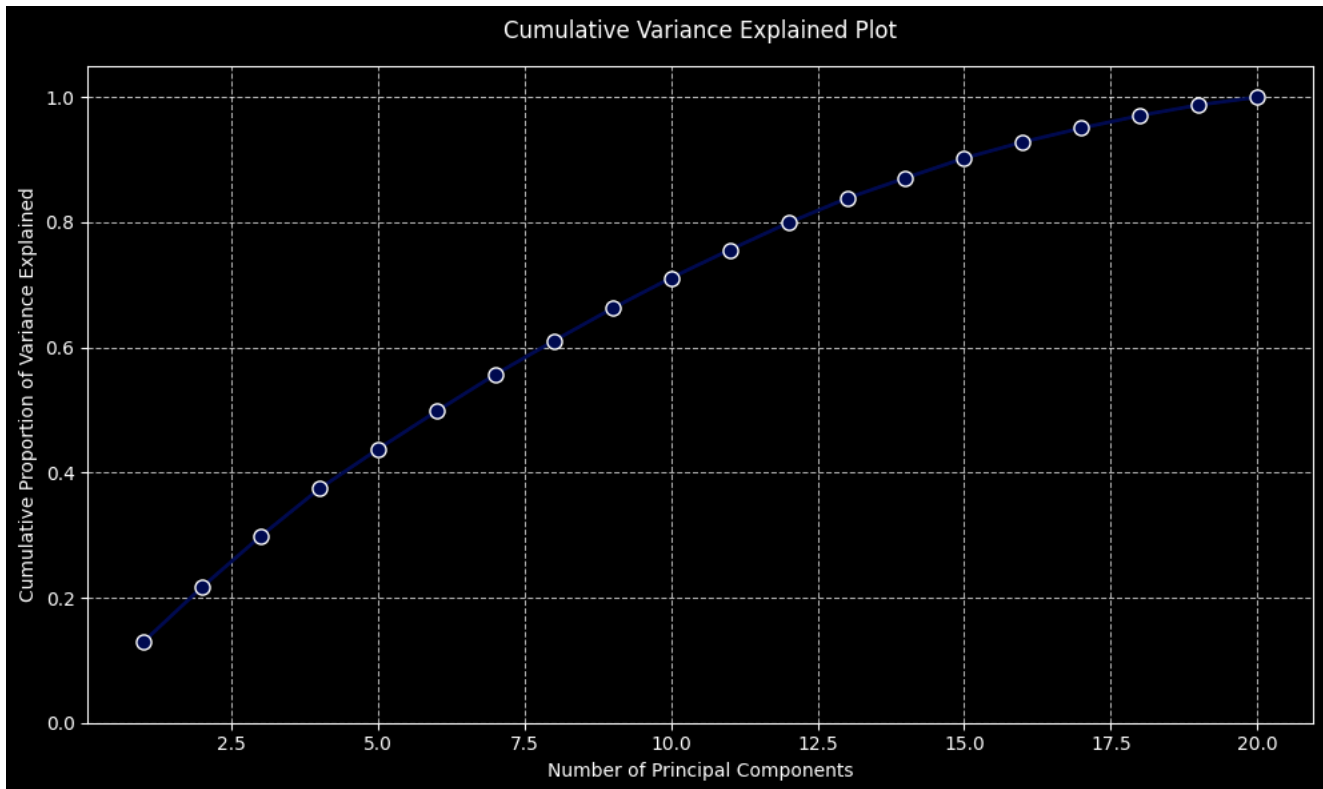
```
1 # your code for scree plot here
2
3 import matplotlib.pyplot as plt
4 import numpy as np
5 plt.style.use('dark_background')
6 SYRACUSE_ORANGE = '#F76900'
7 pca_model = pipe_pca_model_1.stages[-1]
8 n_components = len(pca_model.explainedVariance)
9
10 plt.figure(figsize=(10, 6))
11 plt.plot(
12     np.arange(1, n_components + 1),
13     pca_model.explainedVariance,
14     marker='o',
15     color=SYRACUSE_ORANGE,
16     linewidth=2
17 )
18
19 plt.title('Scree Plot', pad=15)
20 plt.xlabel('Principal Component')
21 plt.ylabel('Proportion of Variance Explained')
22 plt.grid(True, linestyle='--', alpha=0.7)
23
24 plt.tight_layout()
25 plt.show()
```



```

1 # your code for cumulative variance explained plot here
2
3 plt.style.use('dark_background')
4 SYRACUSE_BLUE = '#000E54'
5
6 pca_model = pipe_pca_model_1.stages[-1]
7 n_components = len(pca_model.explainedVariance)
8
9 cumulative_variance = np.cumsum(pca_model.explainedVariance)
10
11 plt.figure(figsize=(10, 6))
12 plt.plot(
13     np.arange(1, n_components + 1),
14     cumulative_variance,
15     marker='o',
16     color=SYRACUSE_BLUE,
17     linewidth=2,
18     markersize=8,
19     markeredgewidth=1,
20     markeredgewidth=1,
21     markeredgewidth=1,
22 )
23 plt.title('Cumulative Variance Explained Plot', pad=15)
24 plt.xlabel('Number of Principal Components')
25 plt.ylabel('Cumulative Proportion of Variance Explained')
26 plt.grid(True, linestyle='--', alpha=0.7)
27
28 plt.ylim(bottom=0, top=1.05)
29
30 plt.tight_layout()
31 plt.show()
32
33 plt.show()

```



your answers here

Q8

Starting with `pipe_pca_1` from the previous question, transform the pipeline and save the resulting dataframe to a variable named `pca_fun`.

Extract the output from the standard scaler column from the first row of `pca_fun` and store in a variable named `row1_centered`.

Manually compute 5 PCA scores by projecting `row1_centered` onto the first 5 loading vectors which were computed in your PCA object. Save the 5 projected pca scores in a variable called `proj_scores`.

Extract the first 5 PCA scores from the first row of the `pca_fun` scores column and save them in a variable named `pca_fun_scores`.

```
1 # your answer here
2 pca_fun = pipe_pca_model_1.transform(ds_features)
3 row1_centered = pca_fun.select('scaled_features').first()['scaled_features']
4 loading_vectors = pipe_pca_model_1.stages[1].pc.toArray()[:, :5]
5 proj_scores = np.dot(row1_centered.toArray(), loading_vectors)
6 pca_fun_scores = pca_fun.select('scores').first()['scores'][:5]
```

```
1 # do not modify
2 print(proj_scores)
3 print(pca_fun_scores)
```

```
[-0.8273713 -0.0963639  1.94471186 -0.09190493  2.02351546]
[-0.8273713 -0.0963639  1.94471186 -0.09190493  2.02351546]
```

Q9

Perform an **inverse transform** on the `proj_scores` variable and store the result in a variable named `inverse`.

The grading cell below prints `inverse` and the original `row1_centered` data such that they are right next to each other.

If `inverse` is different than `row1_centered`, explain why. How you could modify the forward and reverse transformation process such that the resulting `inverse` data almost exactly matches `row1_centered`.

```
1 # your code inverse_pca = np.dot(proj_scores, loading_vectors.T)
2
3 inverse_pca = np.dot(proj_scores, loading_vectors.T)
```

```

4 scaler_model = pipe_pca_model_1.stages[0]
5 mean = scaler_model.mean
6 std = scaler_model.std
7 inverse = inverse_pca * std + mean

```

your answer here

```

1 # do not modify
2 print(row1_centered[0:5])
3 print(inverse[0:5])

```

```

[-0.33920093 -0.31538444 -0.11677724  0.80235477 -0.23850672]
[ 0.2392728  -0.39532054 -0.05361332  0.47784174 -0.08726838]

```

Q10

Implement your modification so that `row1_centered` and `inverse` match almost exactly.

```

1 # your code here
2 pca_fun = pipe_pca_model_1.transform(ds_features)
3 row1_centered = pca_fun.select('scaled_features').first()['scaled_features']
4 loading_vectors = pipe_pca_model_1.stages[1].pc.toArray() # All loading vectors
5 proj_scores = np.dot(row1_centered.toArray(), loading_vectors.T)
6 inverse_pca = np.dot(proj_scores, loading_vectors)
7
8 scaler_model = pipe_pca_model_1.stages[0]

```