

# Momento de Retroalimentación: Módulo 1 Utilización, procesamiento y visualización de grandes volúmenes de datos (Portafolio Análisis)

Andrés Alejandro Guzmán González - A01633819

Inteligencia artificial avanzada para la ciencia de datos II

Grupo 502

1 - Configura correctamente tu entorno de trabajo en Colab para utilizar PySpark (Preparación del ambiente de trabajo para Big Data en Colab)

1.1 - Instalar y preparar Google Colab para trabajar con PySpark

```
1 !pip install pyspark py4j

Collecting pyspark
  Downloading pyspark-3.5.0.tar.gz (316.9 MB)
    Preparing metadata (setup.py) ... done
Requirement already satisfied: py4j in /usr/local/lib/python3.10/dist-packages (0.10.9.7)
Building wheels for collected packages: pyspark
  Building wheel for pyspark (setup.py) ... done
  Created wheel for pyspark: filename=pyspark-3.5.0-py3-none-any.whl size=317425344 sha256=ed99f2df847bcd86b63d2be3dbcacca951ba843bc8ed347c0563c66caf38493b
  Stored in directory: /root/.cache/pip/wheels/41/4e/10/c2cf2467f71c678cfc8a6b9ac9241e5e44a01940da8fbb17fc
Successfully built pyspark
Installing collected packages: pyspark
Successfully installed pyspark-3.5.0
```

1.2 - Llamado a librerías

```
1 from pyspark.sql import SparkSession
2 from pyspark.ml.regression import LinearRegression
3 from pyspark.ml import Pipeline
4 from pyspark.ml.classification import LogisticRegression
5 from pyspark.ml.feature import IndexToString, StringIndexer, VectorIndexer
6 from pyspark.ml.clustering import KMeans
7 from pyspark.ml.evaluation import ClusteringEvaluator
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 from sklearn.metrics import confusion_matrix
```

1.3 - Se levanta la sesión en Spark

```
1 spark = SparkSession \
2     .builder \
3     .appName("AlexGuzmn") \
4     .getOrCreate()
5
6 spark
```

2 - Una vez configurado tu ambiente de trabajo selecciona una base de datos que se caracterice por tener un gran volumen de datos.

Las bases de datos se obtuvieron de este [repositorio](#).

3 - Utilizando MLlib genera un modelo inteligente de clasificación, regresión o agrupamiento con la base de datos que seleccionaste

3.1 - Modelo de regresión lineal

```
1 training = spark.read.format("libsvm")\
2     .load("/content/sample_data/sample_linear_regression_data.txt")

1 # Se define el modelo de regresión lineal
2 lr = LinearRegression(maxIter=10, regParam=0.3, elasticNetParam=0.8)
3
4 # Ajuste del modelo
5 lrModel = lr.fit(training)
6
7 # Se imprimen los coeficientes y la intercepción del modelo
8 print("Coefficients: %s" % str(lrModel.coeficients))
9 print("Intercept: %s" % str(lrModel.intercept))
10
11 # Resumen del modelo
12 trainingSummary = lrModel.summary
13 print("numIterations: %d" % trainingSummary.totalIterations)
14 print("objectiveHistory: %s" % str(trainingSummary.objectiveHistory))
15 trainingSummary.residuals.show()
16 residuals = trainingSummary.residuals.toPandas()
17 rmse = trainingSummary.rootMeanSquaredError
18 r2 = trainingSummary.r2
19 print("RMSE: %f" % rmse)
20 print("r2: %f" % r2)

Coefficients: [0.0,0.3229251667740594,-0.3438548034562219,1.915601702345841,0.05288058680386255,0.765962720459771,0.0,-0.15105392669186676,-0.21587930360904645,0.2202536918881343]
Intercept: 0.15989368442397356
numIterations: 6
objectiveHistory: [0.49999999999999994, 0.4967620357443381, 0.49363616643404634, 0.4936351537897608, 0.4936351214177871, 0.49363512062528014, 0.4936351206216114]
+-----+
| residuals|
+-----+
|-9.889232683103197|
| 0.5533794340053553|
|-5.204019455758822|
|-20.56686715507908|
| -9.4497405180564|
|-6.909112502719487|
|-10.00431602969873|
| 2.0623978070504845|
| 3.1117508432954772|
|-15.89360822941938|
| -5.036284254673026|
| 6.4832158769943335|
| 12.429497299109002|
|-20.32003219007654|
| -2.0049838218725|
|-17.867901734183793|
| 7.646455887420495|
|-2.2653482182417406|
|-0.1030892043519545|
|-1.380034070385301|
+-----+
only showing top 20 rows

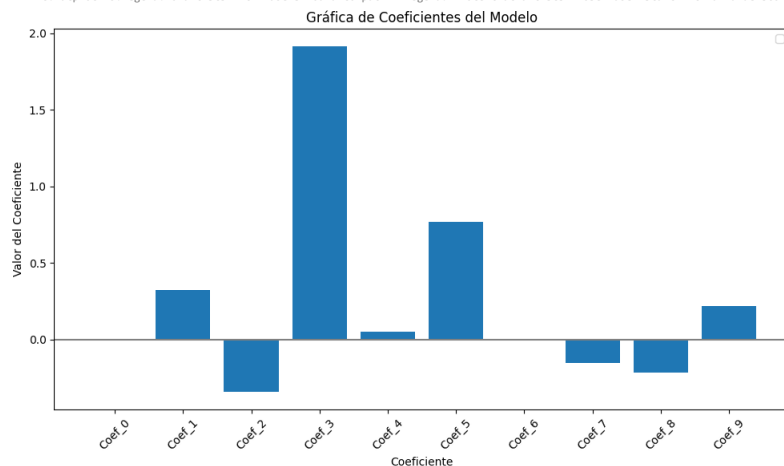
RMSE: 10.189077
r2: 0.022861
```

```

1 coef_labels = [f'Coef_{i}' for i in range(len(lrModel.coeficients))]
2
3 # Crear el gráfico de barras
4 plt.figure(figsize=(10, 6))
5 plt.bar(coef_labels, lrModel.coeficients)
6 plt.xlabel('Coeficiente')
7 plt.ylabel('Valor del Coeficiente')
8 plt.title('Gráfica de Coeficientes del Modelo')
9 plt.xticks(rotation=45)
10 plt.axhline(0, color='gray', linestyle='--', linewidth=1.5)
11
12 plt.tight_layout()
13 plt.legend()
14 plt.show()

```

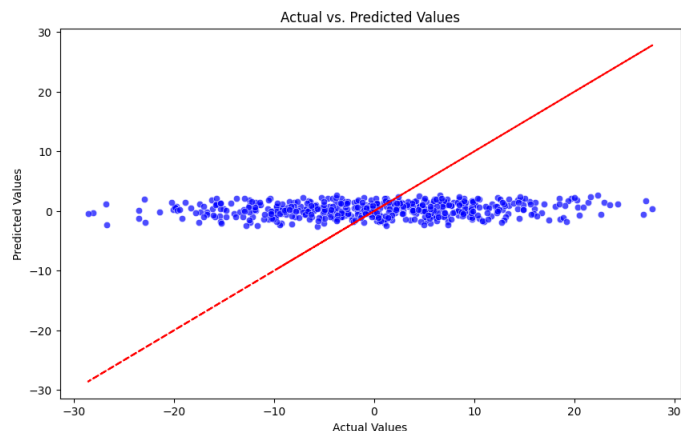
WARNING:matplotlib.legend:No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.



```

1 # Realizar predicciones utilizando el modelo y graficar la dispersión
2 predictions = lrModel.transform(training).toPandas()
3 plt.figure(figsize=(10, 6))
4 sns.scatterplot(x='label', y='prediction', data=predictions, color="b",alpha=0.7)
5 plt.plot(predictions['label'], predictions['label'], color='r', linestyle='--', label='Regression Line')
6 plt.title("Actual vs. Predicted Values")
7 plt.xlabel("Actual Values")
8 plt.ylabel("Predicted Values")
9 plt.show()

```



### 3.2 - Modelo de regresión logística

```

1 # Carga de datos
2 training = spark.read.format("libsvm").load("/content/sample_data/sample_libsvm_data.txt")
3
4 lr = LogisticRegression(maxIter=10, regParam=0.3, elasticNetParam=0.8)
5
6 # Ajuste del modelo
7 lrModel = lr.fit(training)
8
9 # Imprimir los coeficientes o intersecciones de la regresión logística
10 print("Coeficients: " + str(lrModel.coeficients))
11 print("Intercept: " + str(lrModel.intercept))
12
13 mlr = LogisticRegression(maxIter=10, regParam=0.3, elasticNetParam=0.8, family="multinomial")
14
15 mlrModel = mlr.fit(training)
16 print("Multinomial coeficients: " + str(mlrModel.coefficientMatrix))
17 print("Multinomial intercepts: " + str(mlrModel.interceptVector))

```

```

Coeficients: (692,[272,300,323,350,351,378,379,405,406,407,428,433,434,435,455,456,461,462,483,484,489,490,496,511,512,517,539,540,568],[ -7.520689871384125e-05, -8.115773146847006e-05, 3.814692771846427e-05, 0.0003776490540424338, 0.00
Intercept: -0.5991460286401438
Multinomial coeficients: 2 X 692 CSRMatrix
(0,272) 0.0001
(0,300) 0.0001
(0,350) -0.0002
(0,351) -0.0001
(0,378) -0.0003
(0,379) -0.0002
(0,405) -0.0002
(0,406) -0.0004
(0,407) -0.0002
(0,433) -0.0003
(0,434) -0.0005
(0,435) -0.0001
(0,456) 0.0
(0,461) -0.0002
(0,462) -0.0004
(0,483) 0.0001

```

```

...
...
Multinomial intercepts: [0.27505875857180895,-0.27505875857180895]

1 '''
2 Se extrae el resumen del modelo de Regresión Logística devuelto en la
3 instancia anterior, entrenado en el ejemplo anterior.
4 '''
5 trainingSummary = lrModel.summary
6
7 # Se obtiene el objetivo por iteración.
8 objectiveHistory = trainingSummary.objectiveHistory
9 print("objectiveHistory:")
10 for objective in objectiveHistory:
11     print(objective)
12
13 trainingSummary.roc.show()
14 print("areaUnderROC: " + str(trainingSummary.areaUnderROC))
15
16 # Se establece el umbral del modelo para maximizar la puntuación F.
17 fMeasure = trainingSummary.fMeasureByThreshold
18 maxfMeasure = fMeasure.groupBy().max('F-Measure').select('max(F-Measure)').head()
19 bestThreshold = fMeasure.where(fMeasure['F-Measure'] == maxfMeasure['max(F-Measure)']) \
20     .select('threshold').head()['threshold']
21 lr.setThreshold(bestThreshold)

objectiveHistory:
0.6833149135741672
0.6661906127558117
0.6207433672479603
0.6131541253123869
0.6059149689952393
0.5923656241678249
0.589823308283802
0.5868012627420282
0.5844432058719142
0.5830790068041746
0.5807015754032354
+---+-----+
|FPR|          TPR|
+---+-----+
|0.0|          0.0|
|0.0| 0.017543859649122806|
|0.0| 0.03508771929824561|
|0.0| 0.05263157894736842|
|0.0| 0.07017543859649122|
|0.0| 0.08771929824561403|
|0.0| 0.10526315789473684|
|0.0| 0.12280701754385964|
|0.0| 0.14035087719298245|
|0.0| 0.15789473684210525|
|0.0| 0.17543859649122806|
|0.0| 0.19298245614035087|
|0.0| 0.21052631578947367|
|0.0| 0.22807017543859648|
|0.0| 0.24561403508771928|
|0.0| 0.26315789473684211|
|0.0| 0.2807017543859649|
|0.0| 0.2982456140350877|
|0.0| 0.3157894736842105|
|0.0| 0.3333333333333333|
+---+-----+
only showing top 20 rows

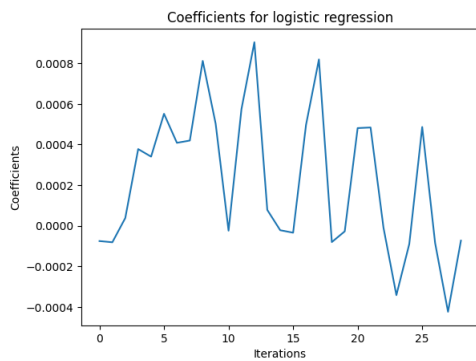
areaUnderROC: 1.0
LogisticRegression_175bffb3cef

1 len(lrModel.coefficients)

692

1 plt.plot(lrModel.coefficients)
2 plt.xlabel("Iterations")
3 plt.ylabel("Coefficients")
4 plt.title("Coefficients for logistic regression")
5 plt.show()

```



### 3.3 - Modelo de clusters

```

1 # Carga de datos
2 dataset = spark.read.format("libsvm").load("/content/sample_data/sample_kmeans_data.txt")
3
4 # Se entrena un modelo de k-means.
5 kmeans = KMeans().setK(2).setSeed(1)
6 model = kmeans.fit(dataset)
7
8 # Se realizan las predicciones
9 predictions = model.transform(dataset)
10
11 # Se evalúa el agrupamiento calculando el puntaje Silhouette.
12 evaluator = ClusteringEvaluator()
13
14 silhouette = evaluator.evaluate(predictions)
15 print("Silhouette with squared euclidean distance = " + str(silhouette))
16
17 # Se muestra el resultado
18 centers = model.clusterCenters()
19 print("Cluster Centers: ")
20 for center in centers:
21     print(center)
22

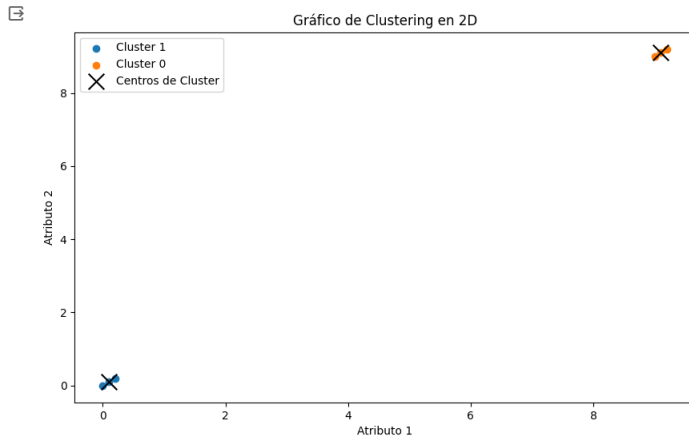
```

```

Silhouette with squared euclidean distance = 0.9997530305375207
Cluster Centers:
[9.1 9.1 9.1]
[0.1 0.1 0.1]

1 import matplotlib.pyplot as plt
2
3 # Extrae los datos de entrada (atributos) y las etiquetas de clúster de las predicciones.
4 features = predictions.select("features").rdd.map(lambda row: row.features)
5 cluster_labels = predictions.select("prediction").rdd.map(lambda row: row.prediction)
6
7 # Convierte los datos a una lista de Python.
8 features_list = features.collect()
9 cluster_labels_list = cluster_labels.collect()
10
11 # Convierte los centros de clúster a una lista de Python.
12 cluster_centers = model.clusterCenters()
13
14 # Separa los datos en función de las etiquetas del clúster.
15 cluster_data = {}
16 for i, label in enumerate(cluster_labels_list):
17     if label in cluster_data:
18         cluster_data[label].append(features_list[i])
19     else:
20         cluster_data[label] = [features_list[i]]
21
22 # Grafica los datos en 2D.
23 plt.figure(figsize=(10, 6))
24
25 for label, data in cluster_data.items():
26     data = list(zip(*data)) # Desempaqueta los datos
27     plt.scatter(data[0], data[1], label=f'Cluster {label}')
28
29 # Grafica los centros de clúster.
30 center_data = list(zip(*cluster_centers)) # Desempaqueta los centros de clúster
31 plt.scatter(center_data[0], center_data[1], c='black', marker='x', s=200, label='Centros de Cluster')
32
33 plt.xlabel('Atributo 1')
34 plt.ylabel('Atributo 2')
35 plt.legend()
36 plt.title('Gráfico de Clustering en 2D')
37 plt.show()
38

```



```

1 # Grafica los datos en 3D.
2 fig = plt.figure()
3 ax = fig.add_subplot(111, projection='3d')
4
5 for label, data in cluster_data.items():
6     data = list(zip(*data)) # Desempaqueta los datos
7     ax.scatter(data[0], data[1], data[2], label=f'Cluster {label}')
8
9 # Grafica los centros de clúster.
10 center_data = list(zip(*cluster_centers)) # Desempaqueta los centros de clúster
11 ax.scatter(center_data[0], center_data[1], center_data[2], c='black', marker='x', s=200, label='Centros de Cluster')
12
13 ax.set_xlabel('Atributo 1')
14 ax.set_ylabel('Atributo 2')
15 ax.set_zlabel('Atributo 3')
16 ax.legend()
17 plt.show()
18

```

