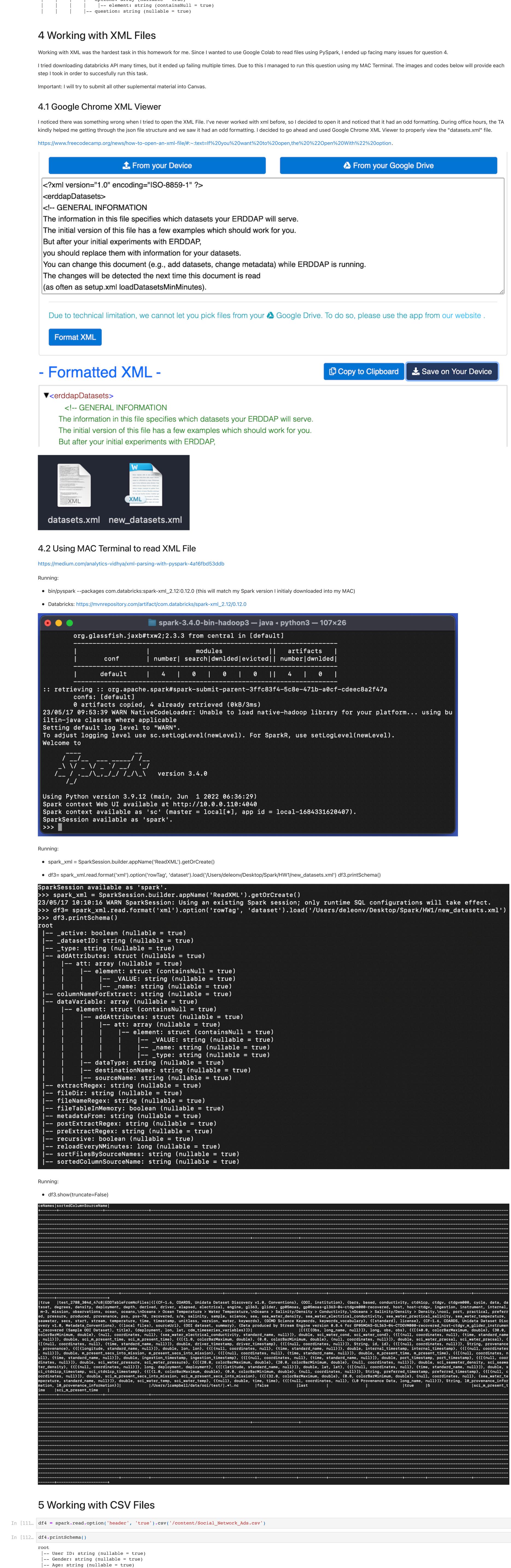
	1 Spark Installation 1.1 MAC Terminal
	Last login: Mon May 15 14:25:19 on ttys000 [(base) deleonv@Vicentes-MacBook-Air ~ % pwd /Users/deleonv [(base) deleonv@Vicentes-MacBook-Air Desktop [(base) deleonv@Vicentes-MacBook-Air Desktop [(base) deleonv@Vicentes-MacBook-Air Desktop % cd Spark [(base) deleonv@Vicentes-MacBook-Air Spark % cd spark-3.4.0-bin-hadoop3 [(base) deleonv@Vicentes-MacBook-Air spark-3.4.0-bin-hadoop3 % bin/pyspark Python 3.9.12 (main, Jun 1 2022, 06:36:29) [[clang 12.0.0] :: Anaconda, Inc. on darwin Type "help", "copyright", "credits" or "license" for more information. 23/05/16 17:57:95 WARN Utils: Your hostname, Vicentes-MacBook-Air.local resolves to a loopback address: 127.0.0.1; using 16.0.110 instead (on interface en0) 23/05/16 17:57:05 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address Setting default log level to "WARN". To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel). 23/05/16 17:57:06 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform using builtin-java class where applicable Welcome to
	/_////////
	 1.1.1 Code used for MAC Terminal pwd cd Desktop cd Spark cd spark-3.4.0-bin-hadoop3 bin/pyspark
	• quit() 1.2 PySpark + Google Colab Installation: I really like using Google Colab and after doing a couple of research I came across this link: https://www.youtube.com/watch?v=Ev_mwYGAbcg Due to this, I decided to try and to Homework 1 using Google Colab.
n []: n [2]: n [3]: n [4]:	<pre>!pip install pyspark py4j from pyspark.sql import SparkSession spark = SparkSession.builder.appName('DataFrame').getOrCreate() # https://stackoverflow.com/questions/69553072/reading-a-xml-file-in-pyspark spark_xml = SparkSession.builder.appName('ReadXML').getOrCreate()</pre>
	2 Working with Text Files Reading Text Files: This was the best source (youtube link) I came across for text files. It was extremely useful for me: https://www.youtube.com/watch?v=9Dpng3bDsPl
	<pre>https://www.geeksforgeeks.org/read-text-file-into-pyspark-dataframe/ https://spark.apache.org/docs/latest/sql-data-sources-text.html # https://spark.apache.org/docs/latest/api/python/reference/pyspark.sql/api/pyspark.sql.functions.col.html # https://spark.apache.org/docs/3.1.2/api/python/reference/api/pyspark.sql.functions.when.html from pyspark.sql.functions import col from pyspark.sql.functions import when df1 = spark.read.option('header', 'true') \</pre>
	<pre>.option('delimiter', '\t') \ .option('inferSchema', 'true') \ .csv('Raisin_Dataset.txt') # https://sparkbyexamples.com/spark/spark-show-full-column-content-dataframe/ dfl.show(n = 5, truncate = False) ++</pre>
[66]:	87524 442.2460114 253.291155 0.819738392 90546 0.758650579 1184.04 Kecimen
	root Area: integer (nullable = true) MajorAxisLength: double (nullable = true) MinorAxisLength: double (nullable = true) Eccentricity: double (nullable = true) ConvexArea: integer (nullable = true) Extent: double (nullable = true) Perimeter: double (nullable = true) Class: string (nullable = true) Class: string (nullable = true)
[67]:	new_df1 = df1.withColumnRenamed("Class", "Categorical Class") new_df1.show(n = 5, truncate=False) ++
[72]:	ttttttt
[74]:	Kecimen Kecime
[80]:	final_df.show(n = 5, truncate = False) +
[84]:	<pre># https://sparkbyexamples.com/pyspark/pyspark-select-distinct/ KvsB = final_df.select('Categorical Class', 'Numerical Class').distinct() KvsB.show() ++ Categorical Class Numerical Class ++ Kecimen</pre>
	<pre>#</pre>
	+
	quiz: struct (nullable = true)
	Working with XML was the hardest task in this homework for me. Since I wanted to use Google Colab to read files using PySpark, I ended up facing many issues for question 4. I tried downloading databricks API many times, but it ended up failing multiple times. Due to this I managed to run this question using my MAC Terminal. The images and codes below will provide ea step I took in order to successfully run this task. Important: I will try to submit all other suplemental material into Canvas. 4.1 Google Chrome XML Viewer
	I noticed there was something wrong when I tried to open the XML File. I've never worked with xml before, so I decided to open it and noticed that it had an odd formatting. During office hours, the kindly helped me getting through the json file structure and we saw it had an odd formatting. I decided to go ahead and used Google Chrome XML Viewer to properly view the "datasets.xml" file. https://www.freecodecamp.org/news/how-to-open-an-xml-file/#:~:text=If%20you%20want%20to%20open,the%20%22Open%20With%22%20option. From your Google Drive
	xml version="1.0" encoding="ISO-8859-1" ? <erddapdatasets> <!-- GENERAL INFORMATION</p--> The information in this file specifies which datasets your ERDDAP will serve. The initial version of this file has a few examples which should work for you. But after your initial experiments with ERDDAP, you should replace them with information for your datasets. You can change this document (e.g., add datasets, change metadata) while ERDDAP is running. The changes will be detected the next time this document is read (as often as setup.xml loadDatasetsMinMinutes).</erddapdatasets>
	Due to technical limitation, we cannot let you pick files from your ♣ Google Drive. To do so, please use the app from our website . Format XML Copy to Clipboard ♣ Save on Your Device
	▼ <erddapdatasets> <!-- GENERAL INFORMATION</p--> The information in this file specifies which datasets your ERDDAP will serve. The initial version of this file has a few examples which should work for you. But after your initial experiments with ERDDAP,</erddapdatasets>
	The state of the s
	datasets.xml new_datasets.xml 4.2 Using MAC Terminal to read XML File https://medium.com/analytics-vidhya/xml-parsing-with-pyspark-4a16fbd53ddb Running:
	 bin/pysparkpackages com.databricks:spark-xml_2.12:0.12.0 (this will match my Spark version I initially downloaded into my MAC) Databricks: https://mvnrepository.com/artifact/com.databricks/spark-xml_2.12/0.12.0 spark-3.4.0-bin-hadoop3 — java < python3 — 107×26 org.glassfish.jaxb#txw2;2.3.3 from central in [default] modules artifacts
	conf number search dwnlded evicted number dwnlded
	Setting default log level to "WARN". To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel). Welcome to
	Using Python version 3.9.12 (main, Jun 1 2022 06:36:29) Spark context Web UI available at http://10.0.0.110:4040 Spark context available as 'sc' (master = local[*], app id = local-1684331620407). SparkSession available as 'spark'. >>> ■ Running:
	 spark_xml = SparkSession.builder.appName('ReadXML').getOrCreate() df3= spark_xml.read.format('xml').option('rowTag', 'dataset').load('/Users/deleonv/Desktop/Spark/HW1/new_datasets.xml') df3.printSchema() SparkSession available as 'spark'. >>> spark_xml = SparkSession.builder.appName('ReadXML').getOrCreate() 23/05/17 10:10:16 WARN SparkSession: Using an existing Spark session; only runtime SQL configurations will take effect. >>> df3= spark_xml.read.format('xml').option('rowTag', 'dataset').load('/Users/deleonv/Desktop/Spark/HW1/new_datasets.xml >>> df3.printSchema()
	root
	active: boolean (nullable = true) datasetID: string (nullable = true) type: string (nullable = true) addAttributes: struct (nullable = true)
	active: boolean (nullable = true) datasetID: string (nullable = true) type: string (nullable = true) addAttributes: struct (nullable = true)
	active: boolean (nullable = true) datasetID: string (nullable = true) type: string (nullable = true) addAttributes: struct (nullable = true)
	active: boolean (nullable = true) datasetID: string (nullable = true) type: string (nullable = true) addAttributes: struct (nullable = true) att: array (nullable = true) att: array (nullable = true) element: struct (containsNull = true) VALUE: string (nullable = true) OolumnNameForExtract: string (nullable = true) columnNameForExtract: string (nullable = true) dataVariable: array (nullable = true) element: struct (containsNull = true) addAttributes: struct (nullable = true) att: array (nullable = true) att: array (nullable = true) pelement: struct (containsNull = true) lement: struct (containsNull = true) lement: string (nullable = true) lement: string (nullable = true) dataType: string (nullable = true) destinationName: string (nullable = true) destinationName: string (nullable = true) fileDir: string (nullable = true) fileDir: string (nullable = true) fileNameRegex: string (nullable = true) fileNameRegex: string (nullable = true) metadataFrom: string (nullable = true) postExtractRegex: string (nullable = true) prefxtractRegex: string (nullable = true) prefxtractRegex: string (nullable = true)



-- EstimatedSalary: string (nullable = true)

+----+ |User ID |Gender|Age|EstimatedSalary|Purchased|

+----+

Select Columns 2 and 3 (Gender and Age respectively) and print them

0

0

0

0

0

0

0

| 1 0

0 |

Rename the column "EstimatedSalary" to "Annual Salary" and "Purchased" to "Score".

0 0

0 0

| 1

0

0

In [121... new_df4 = df4.withColumnRenamed('EstimatedSalary', 'Annual Salary') new_df4 = new_df4.withColumnRenamed('Purchased', 'Score')

-- Purchased: string (nullable = true)

Print 10 rows

In [115...] df4.show(n = 10, truncate = False)

|15624510|Male |19 |19000

|15810944|Male |35 |20000

|15668575|Female|26 |43000

|15603246|Female|27 |57000

|15804002|Male |19 |76000

|15728773|Male |27 |58000

|15598044|Female|27 |84000

|15600575|Male |25 |33000 |15727311|Female|35 |65000

only showing top 10 rows

In [117... df4.select('Gender', 'Age').show()

only showing top 20 rows

new_df4.show(n = 10, truncate=False)

|15624510|Male |19 |19000 |15810944|Male |35 |20000 |15668575|Female|26 |43000

|15603246|Female|27 |57000 |15804002|Male |19 |76000

|15728773|Male |27 |58000 |15598044|Female|27 |84000

|15694829|Female|32 |150000

|15600575|Male |25 |33000

|15727311|Female|35 |65000

only showing top 10 rows

+----+ |User ID |Gender|Age|Annual Salary|Score| +----+

+----+ |Gender|Age| +----+ Male| 19| Male | 35 |Female| 26| |Female| 27| Male| 19| Male| 27| |Female| 27| |Female| 32| Male| 25| |Female| 35| |Female| 26| |Female| 26| Male | 20 | Male| 32| Male| 18| Male 29 Male| 47| Male| 45| Male | 46| |Female| 48| +----+

|15694829|Female|32 |150000