

# REPORT

---

Mohammed Khursheed Ali Khan / ID# 0662357

April 12, 2020

## Motivation

The reason behind using the ACS PUMS US census data was that it is a data for housing and population of us citizens for the years 2013-2017, and includes information of persons in all the states, and the resulting data has a good set of attributes with a mix of both numerical, categorical and missing values. There by helps in good analysis of the data in a distributed in manner in spark. Also being a vast data of 14GB in size, with 100M records and 1000 columns, it truly defines the big Data we are talking about and really helps in learning the spark framework. Also being data hungry the domain of Machine Learning, It served as a good purpose of ML in distributed environment with Spark.

## 1 Problem Addressed

The problem I am addressing is broad in nature as I am trying to gain insights on the data trying to discover some patterns on the data based on demographics of the population. Once I have gained those insights, I am working on further trying to apply ML algorithms on the data for further predictive analysis.

As mentioned in this project I am trying to carry a predictive analysis on the Industry the people are working in based on their age, schooling attained, their field of degree and the sectors they are working in. Also not to forget the gender that is an important attribute for determining the industry.

## 2 Approach Overview

Lets get started with the environment I have deployed a standalone Spark cluster on a *Azure VM* with *14GB RAM* and *4vcpu cores*. And as for the executors I have used local cluster of 4 cores such that the driver/executors are run as multiple threads across 4 cores locally.

As for the data, the API for data loading was not responding, so I resorted to transferring files from my local system to remote work environment via `SCP transfer`.

Once the initial environment was setup, and data was loaded, then came the data analysis part,

since I have used `pyspark` API for working, I accordingly have used the corresponding constructs, for EDA, which involved carrying out some analysis for understanding the trends in data and feature engineering. Note that I haven't used any plots as I started working with `python 2` and once I reached the end of my work only to realize that a lot of `pyspark` visualizations supported work on `python 3`. However, I have included some histograms that define the data trends.

Once that was done and feature engineering was done, a machine learning pipeline was created to fit a model on our data and achieve the above mentioned problem statement.

Further in the next section on how each component was done and results was achieved such as how EDA, what core concepts were used for feature engineering and what kind a pipeline was developed with algorithms being used.

### 3 Approach Explained

The first step I achieved was data loading as data frames, because those built on top of RDD provide much easier munging as the data is given and structure in dataframes which is justified by the fact that we have a highly organized dataset structured in nature. Also the recent changes and deprecation of the `MLlib` package which focused on RDD to dataframe based ml package further motivated me to use dataframes for analysis and working.

Along with the data, I also had to load the data dictionary as some of the attributes were coded .

Once that was done the next step involved some data analysis, EDA, this involved getting some insights by data filtering on dataframes which are distributed and immutable in nature. So some analysis I did was regarding the Class of workers above and below age 50 and the industry of work among the genders. The analysis gave some interesting and realistic results, wherein the x industry was male dominant and y industry was female dominant. Also for the workers above age 50, they were a workers and below 50 where mostly b workers.

For EDA, as defined in previous section, one way I could develop visualizations such as correlation and scatter and histograms etc. would be that to use `toPandas` and convert the dataframe to pandas dataframe and carry on the pandas visualizations, but pandas doesn't work the same ways as `pyspark` dataframes work and so I didn't use that Idea as it defies the whole big Data definition.

That being said, the next stage of feature engineering was the crucial phase. First, I decided on the features to work on for developing the model. The feature selection included selecting features that could help in predicting the industry the persons are working in. The best attributes that could define this were, schooling attainment, field of degree they have their gender as we had seen earlier in EDA how gender effects the class of worker and industry the persons are working in, and mostly importantly the class of workers an individual is in an industry. These all were selected as features,

along with age. Attributes selected in a nutshell are as follows:

- SCHL - Schooling attained
- AGEP - Age
- FOD1P - Field of Degree first
- FOD2P - Field of Degree second
- COW - Class of Workers
- INDP - Industry working in

An important step I carried before I proceeded ahead was the data cleaning, this involved filtering the records we need and we feel are relevant to our analysis. One such filter was the individuals with age less than 17, as we need individuals who have worked in industry, and obviously anyone less than age of 17 hasn't worked without the loss of generality.

The next step was the feature engineering phase. Involved transforming the categorical features as one-hot vectors for machine learning model. The categorical features were first converted to category indices via `stringIndexer` and then using a `oneHotEstimator` they were transformed to one-hot encoding. Similarly labels were also which are `INDP` transformed into categorical indices using `stringIndexer`.

Once this is done, the next step involved combining these features into a feature vector for building a ML model. This is helped via the `VectorAssembler` that creates a vector for selected features.

This is then followed by the core ML `pipeline`. All the above transformations and estimations such as the feature engineering and one-hot encoding along with the estimation such as model fitting and predictions are designed as a ML pipeline with multiple stages in a pipeline that define it. These stages are either transformers or estimators. The final stage in the pipeline is the model fitting phase, wherein I select a ML algorithm. I will be using and getting a model by training the data based on a 60-30m random split of training and test data.

Following is a break down of how the pipeline is defined for the above model we are working on.

1. Stage – string indexing and encoding Feature engineering – transformer
2. Stage – string indexing labels – transformer
3. Stage – model fitting - estimator
4. Stage - predictions - estimator

Once the above pipeline was completed we had the evaluation phase, for which I used multiple metrics such as  $F1$ , and *area under the curve*. The model can be further improved based on further hyperparameter tuning, but getting a better accuracy wasn't the main purpose of this work in spark. It was leveraging the power of spark and understanding the architecture via hands-on.

The screenshots below of the DAG which defines how the tasks were executed by the executors gives a good understanding of the distributed and resilient spark design. Since we have 4 cores as can be seen the task are divided as threads locally among the 4 cores. These are then mapped as partitions further after a union is taken from the 4 executor tasks of the results and further stages in DAG involved similar to RDD operations such as `flatmap` and `groupByKey` etc.

While I did manage to get decent results, though not that appealing for a production use case, however it was a good starting with  $ROC$  and  $F1$  scores reaching 0.4 on a scale of 1. I used decision trees and random forest algorithms for model building and as expected got better accuracy with `RandomForest`.

Screenshots in the next section with brief descriptions of each of the above steps.

## 4 Spark webUI

This section briefs on the components (with screenshots) on the *Spark webUI* involved with the understanding of this project and some visualizations of the data.

As you can see in the screenshot below the `master` is set to `local[*]` instructing to use all the local cores of the system.

The screenshot shows a Safari browser window with the following details:

- Address Bar:** Not Secure — 168.62.181.173
- Tab Bar:** hdp-2.6-vm-dev-eus - Microsoft Azure, Microsoft Azure Sponsorships | Usage, Visualization — pandas 1.0.3 documentation, ACS PUMS 2013-2017 Data Analysis - Environment
- Header:** Apache Spark 2.4.5, Jobs, Stages, Storage, Environment, Executors, SQL, ACS PUMS 2013-2017 Data Analysis application UI
- Section: Environment**
- Section: Runtime Information**

Name	Value
Java Version	1.8.0_242 (Oracle Corporation)
Java Home	/usr/lib/jvm/java-1.8.0-openjdk-1.8.0.242.b08-0.el7_7.x86_64/jre
Scala Version	version 2.11.12

- Section: Spark Properties**

Name	Value
spark.app.id	local-1586551535960
spark.app.name	ACS PUMS 2013-2017 Data Analysis
spark.driver.host	sandbox-host.hortonworks.com
spark.driver.port	40278
spark.executor.id	driver
spark.master	local[*]
spark.rdd.compress	True
spark.scheduler.mode	FIFO
spark.serializer.objectStreamReset	100
spark.submit.deployMode	client
spark.ui.showConsoleProgress	true

- Section: System Properties**
- Mac OS Dock:** Icons for various applications including Finder, Mail, Safari, and Visual Studio Code.

Figure 1: Runtime Environment configurations

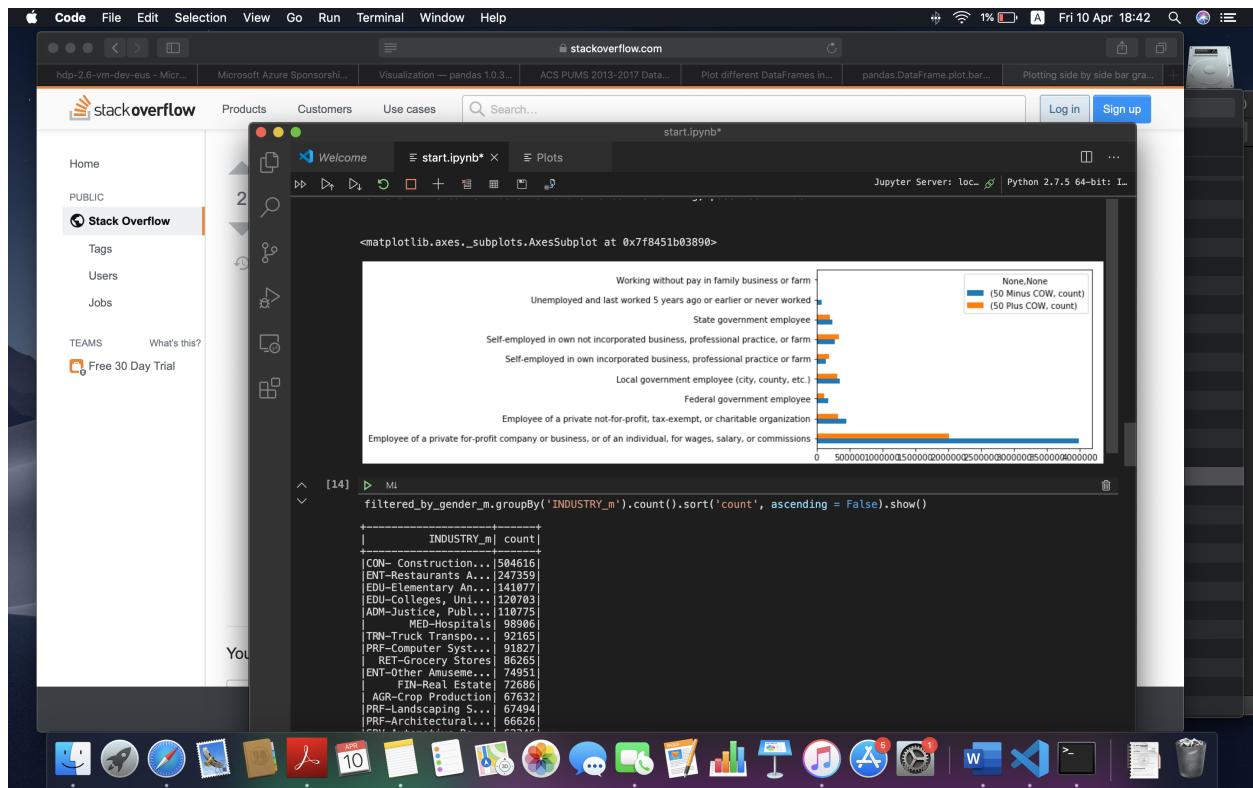


Figure 2: Class of workers for age  $> & < 50$

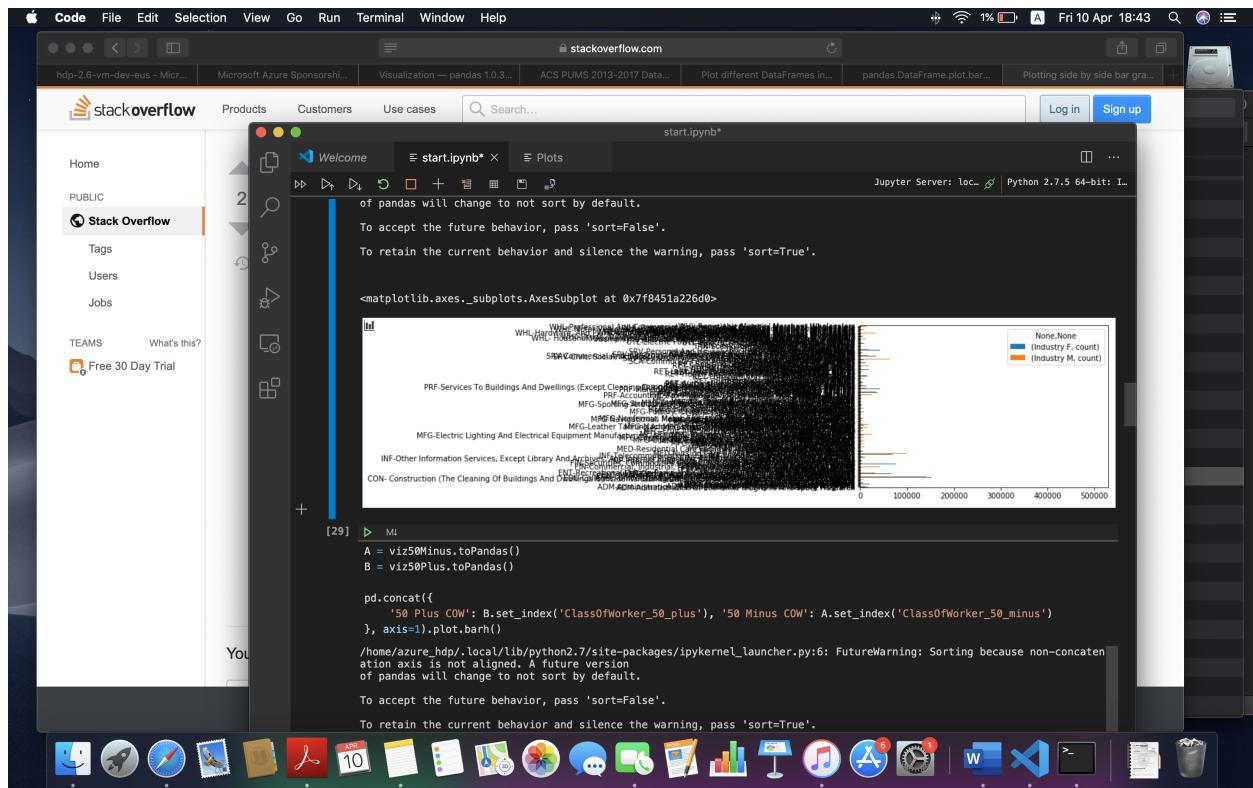


Figure 3: Industry workers working in

The screenshot shows a Safari browser window with the following tabs:

- hdp-2.6-vm-dev-eus - Microsoft Azure
- Microsoft Azure Sponsorships | Usage
- Visualization — pandas 1.0.3 documentation
- ACS PUMS 2013-2017 Data Analysis - Executors

The main content area is titled "Executors". It includes a "Summary" table and a detailed "Executors" table.

**Summary Table:**

RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Blacklisted
Active(1) 0	630.3 KB / 384.1 MB	0.0 B	4	5	0	3612	3617	2.4 h (48 s)	87.9 GB	11.1 MB	5.6 MB	0
Dead(0) 0	0.0 B / 0.0 B	0.0 B	0	0	0	0	0	0 ms (0 ms)	0.0 B	0.0 B	0.0 B	0
Total(1) 0	630.3 KB / 384.1 MB	0.0 B	4	5	0	3612	3617	2.4 h (48 s)	87.9 GB	11.1 MB	5.6 MB	0

**Executors Table:**

Executor ID	Address	Status	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Thread Dump
driver	sandbox-host.hortonworks.com:42210	Active	0	630.3 KB / 384.1 MB	0.0 B	4	5	0	3612	3617	2.4 h (48 s)	87.9 GB	11.1 MB	5.6 MB	<a href="#">Thread Dump</a>

Showing 1 to 1 of 1 entries

Previous 1 Next

Figure 4: Executors tasks running on

The below figures 5, 6, 7, 8 describe the different timeline about how tasks are executed one after the other. Note that in a lot of places you might see that there exists more than one time same tasks repeated, this is because we are using dataframes, which are essentially build over RDD. So any say `df.show()` is placed it is first converted/placed to RDD `first(21)` and further converted/placed as `take(21)`.

**Spark Jobs (?)**

User: azure\_hdp  
 Total Uptime: 4.1 h  
 Scheduling Mode: FIFO  
 Completed Jobs: 73

Event Timeline

Completed Jobs (73)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
72	countByValue at MulticlassMetrics.scala:42 countByValue at MulticlassMetrics.scala:42	2020/04/11 00:46:49	6.3 min	2/2	230/230
71	collectAsMap at MulticlassMetrics.scala:48 collectAsMap at MulticlassMetrics.scala:48	2020/04/11 00:40:33	6.3 min	2/2	230/230
70	showString at NativeMethodAccessoRImpl.java:0 showString at NativeMethodAccessoRImpl.java:0	2020/04/11 00:40:31	2 s	1/1	1/1
69	first at RandomForestClassifier.scala:145 first at RandomForestClassifier.scala:145	2020/04/11 00:40:28	2 s	1/1	1/1
68	collectAsMap at RandomForest.scala:567 collectAsMap at RandomForest.scala:567	2020/04/11 00:38:39	1.8 min	2/2	230/230
67	collectAsMap at RandomForest.scala:567 collectAsMap at RandomForest.scala:567	2020/04/11 00:36:55	1.7 min	2/2	230/230
66	collectAsMap at RandomForest.scala:567 collectAsMap at RandomForest.scala:567	2020/04/11 00:35:18	1.6 min	2/2	230/230
65	collectAsMap at RandomForest.scala:567 collectAsMap at RandomForest.scala:567	2020/04/11 00:33:41	1.6 min	2/2	230/230
64	collectAsMap at RandomForest.scala:567 collectAsMap at RandomForest.scala:567	2020/04/11 00:27:10	6.5 min	2/2	230/230

Figure 5: Jobs 1

Safari File Edit View History Bookmarks Develop Window Help

Not Secure — 168.62.181.173 Fri 10 Apr 19:42

hdp-2.6-vm-dev-eus - Microsoft Azure Microsoft Azure Sponsorships | Usage ACS PUMS 2013-2017 Data Analysis - Stages for All Jobs

APACHE Spark 2.4.5 Jobs Stages Storage Environment Executors SQL ACS PUMS 2013-2017 Data Analysis application UI

## Stages for All Jobs

Active Stages: 1  
Pending Stages: 1  
Completed Stages: 69  
Skipped Stages: 36

**Active Stages (1)**

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
105	flatMap at RandomForest.scala:919	+details (kill)	2020/04/10 23:41:01	1.5 min	26/115 (5 running)	3.3 GB		10.2 KB

**Pending Stages (1)**

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
106	collectAsMap at RandomForest.scala:927	+details	Unknown	Unknown	0/1			

**Completed Stages (69)**

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
104	count at DecisionTreeMetadata.scala:118	+details	2020/04/10 23:34:49	6.2 min	115/115	14.0 GB		
103	take at DecisionTreeMetadata.scala:112	+details	2020/04/10 23:34:45	4 s	1/1	128.1 MB		
102	count at NativeMethodAccessormImpl.java:0	+details	2020/04/10 23:34:44	36 ms	1/1		6.6 KB	
101	count at NativeMethodAccessormImpl.java:0	+details	2020/04/10 23:28:28	6.3 min	115/115	14.0 GB		6.6 KB
100	count at NativeMethodAccessormImpl.java:0	+details	2020/04/10 23:28:27	0.3 s	1/1		6.6 KB	
99	count at NativeMethodAccessormImpl.java:0	+details	2020/04/10 23:22:09	6.3 min	115/115	14.0 GB		6.6 KB
98	showString at NativeMethodAccessormImpl.java:0	+details	2020/04/10 23:22:09	31 ms	1/1	64.0 KB		
97	countByValue at StringIndexer.scala:140	+details	2020/04/10 23:22:07	0.5 s	115/115		788.1 KB	
96	countByValue at StringIndexer.scala:140	+details	2020/04/10 23:15:50	6.3 min	115/115	14.0 GB		788.1 KB

Figure 6: Jobs - ML

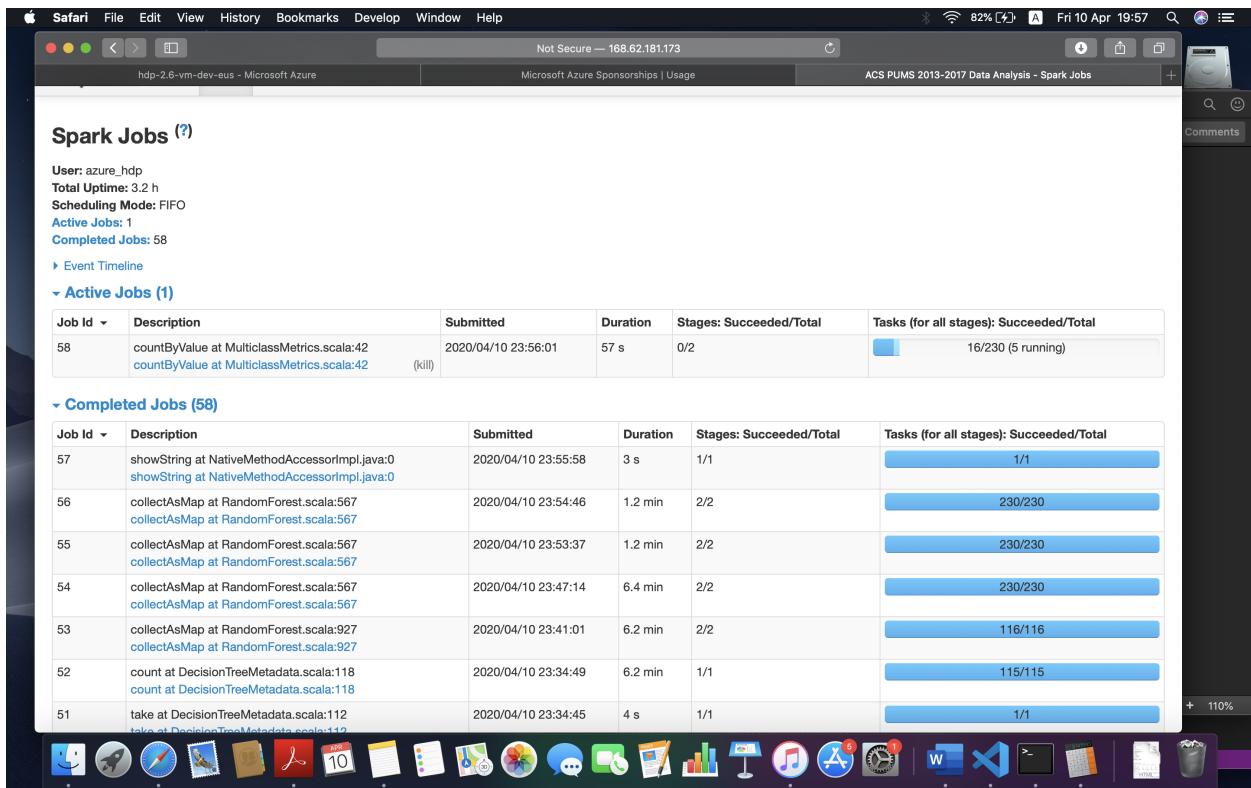


Figure 7: Jobs - metric

Safari File Edit View History Bookmarks Develop Window Help Not Secure — 168.62.181.173 Fri 10 Apr 18:47

hdp-2.6-vm-dev-eus - Microsoft Azure Microsoft Azure Sponsorships | Usage Visualization — pandas 1.0.3 documentation ACS PUMS 2013-2017 Data Analysis - Stages for All J... ACS PUMS 2013-2017 Data Analysis application UI

APACHE Spark 2.4.5 Jobs Stages Storage Environment Executors SQL

### Stages for All Jobs

Active Stages: 1  
Pending Stages: 1  
Completed Stages: 50  
Skipped Stages: 36

**Active Stages (1)**

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
86	countByValue at StringIndexer.scala:140	+details (kill)	2020/04/10 22:44:19	3.5 min 60/115 (5 running)	7.3 GB			68.8 KB

**Pending Stages (1)**

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
87	countByValue at StringIndexer.scala:140	+details	Unknown	Unknown	0/115			

**Completed Stages (50)**

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input	Output	Shuffle Read	Shuffle Write
85	showString at NativeMethodAccessorImpl.java:0	+details	2020/04/10 22:44:11	15 ms 1/1	64.0 KB			
84	toPandas at <ipython-input-30-92958661ac79>:4	+details	2020/04/10 22:39:36	0.4 s 200/200			28.1 KB	
83	toPandas at <ipython-input-30-92958661ac79>:4	+details	2020/04/10 22:39:35	1.0 s 200/200			2.5 MB	28.1 KB
81	toPandas at <ipython-input-30-92958661ac79>:4	+details	2020/04/10 22:39:34	0.8 s 200/200			2.5 MB	
80	toPandas at <ipython-input-30-92958661ac79>:4	+details	2020/04/10 22:33:16	6.3 min 115/115	14.0 GB			2.5 MB
79	toPandas at <ipython-input-30-92958661ac79>:4	+details	2020/04/10 22:33:16	22 ms 1/1	275.8 KB			
78	toPandas at <ipython-input-30-92958661ac79>:3	+details	2020/04/10 22:33:15	0.5 s 200/200			28.1 KB	
75	toPandas at <ipython-input-29-2e42891e34d2>:2	+details	2020/04/10 22:29:15	20 ms 9/9			1167.0 B	
72	toPandas at <ipython-input-29-2e42891e34d2>:1	+details	2020/04/10 22:29:15	22 ms 9/9	1168.0 B			

Figure 8: Jobs - stringIndexer

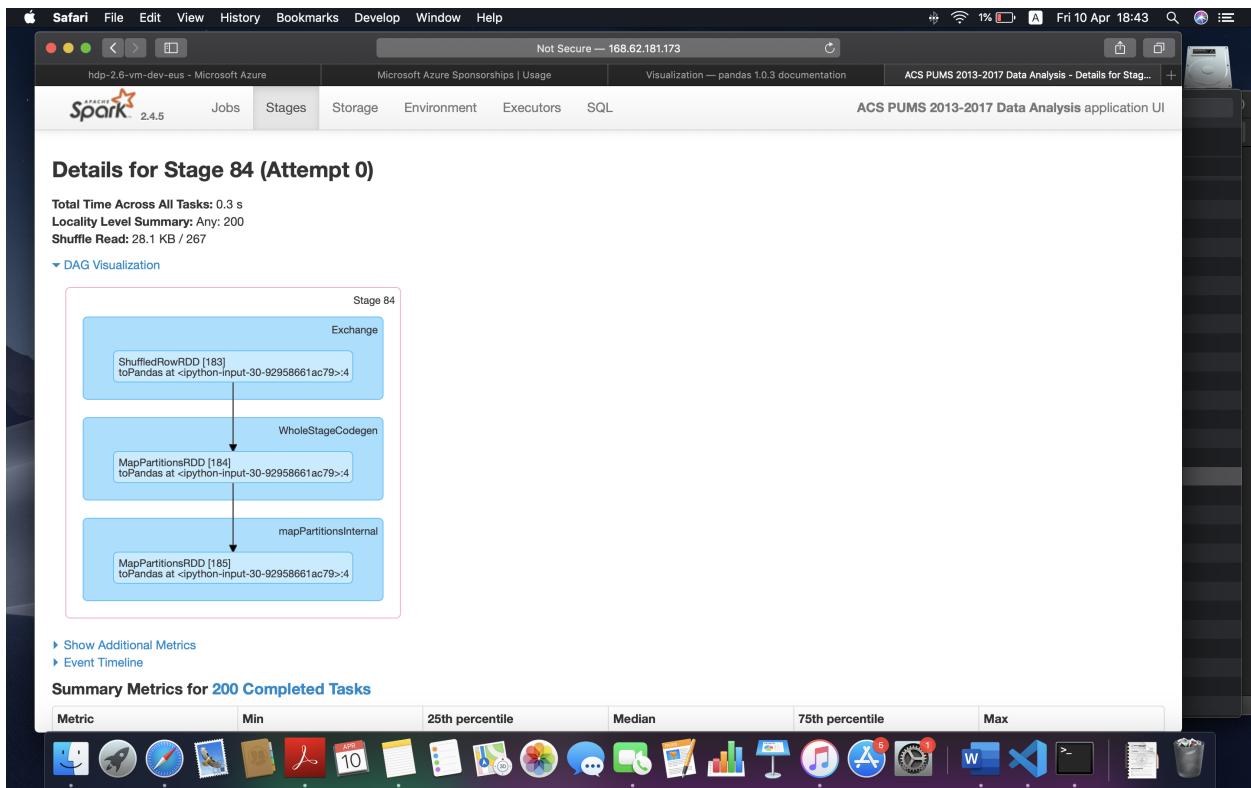


Figure 9: DAG - `toPandas()`

Also observe in the above figure 9, as to how the `toPandas()` DAG validates the above mentioned point, wherein the actions and transformations are performed on RDD after calling them internally.

An important point to mention about the spark design is it caches the tasks executed especially those that involved shuffling such as `reduceBy`, and as a result of this it fetches it from cache and skips that task as can be seen in below figure 10 where tasks are skipped.

**Spark Jobs (?)**

User: azure\_hdp  
Total Uptime: 2.0 h  
Scheduling Mode: FIFO  
Completed Jobs: 41

Event Timeline

Completed Jobs (41)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
40	toPandas at <ipython-input-30-92958661ac79>:4 toPandas at <ipython-input-30-92958661ac79>:4	2020/04/10 22:39:35	1 s	2/2 (1 skipped)	400/400 (115 skipped)
39	toPandas at <ipython-input-30-92958661ac79>:4 toPandas at <ipython-input-30-92958661ac79>:4	2020/04/10 22:33:16	6.3 min	2/2	315/315
38	toPandas at <ipython-input-30-92958661ac79>:4 toPandas at <ipython-input-30-92958661ac79>:4	2020/04/10 22:33:16	26 ms	1/1	1/1
37	toPandas at <ipython-input-30-92958661ac79>:3 toPandas at <ipython-input-30-92958661ac79>:3	2020/04/10 22:33:15	0.5 s	1/1 (2 skipped)	200/200 (315 skipped)
36	toPandas at <ipython-input-29-2e42891e34d2>:2 toPandas at <ipython-input-29-2e42891e34d2>:2	2020/04/10 22:29:15	24 ms	1/1 (2 skipped)	9/9 (315 skipped)
35	toPandas at <ipython-input-29-2e42891e34d2>:1 toPandas at <ipython-input-29-2e42891e34d2>:1	2020/04/10 22:29:15	26 ms	1/1 (2 skipped)	9/9 (315 skipped)
34	toPandas at <ipython-input-28-3387f19da76f>:2 toPandas at <ipython-input-28-3387f19da76f>:2	2020/04/10 22:28:50	24 ms	1/1 (2 skipped)	9/9 (315 skipped)
33	toPandas at <ipython-input-28-3387f19da76f>:1 toPandas at <ipython-input-28-3387f19da76f>:1	2020/04/10 22:28:50	29 ms	1/1 (2 skipped)	9/9 (315 skipped)
32	toPandas at <ipython-input-27-c39a76872c26>:2 toPandas at <ipython-input-27-c39a76872c26>:2	2020/04/10 22:25:47	23 ms	1/1 (2 skipped)	9/9 (315 skipped)

Figure 10: Jobs - Caching

In the below figure 11 the DAG for `groupByKey()` is done with union of all tasks performed on 4 cores and then flatmapping and partitioning these tasks results and finally `countByValue` which defines the `groupByKey()` functionality.

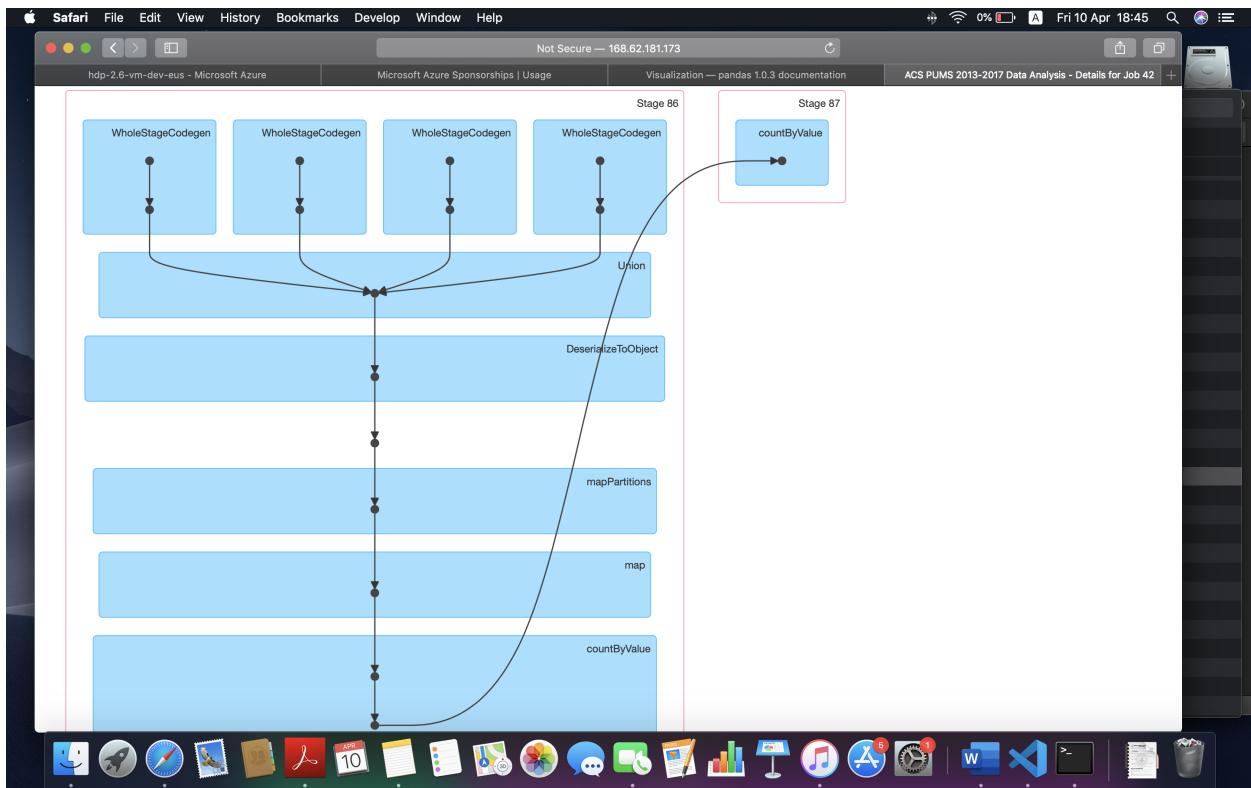


Figure 11: DAG - `groupBy()`

Similar to above, the below figure 12, once string indexing of categorical variables is done, we essentially have indices for each category which are then reduced by key `reduceByKey`, and essentially act as key-value pair.

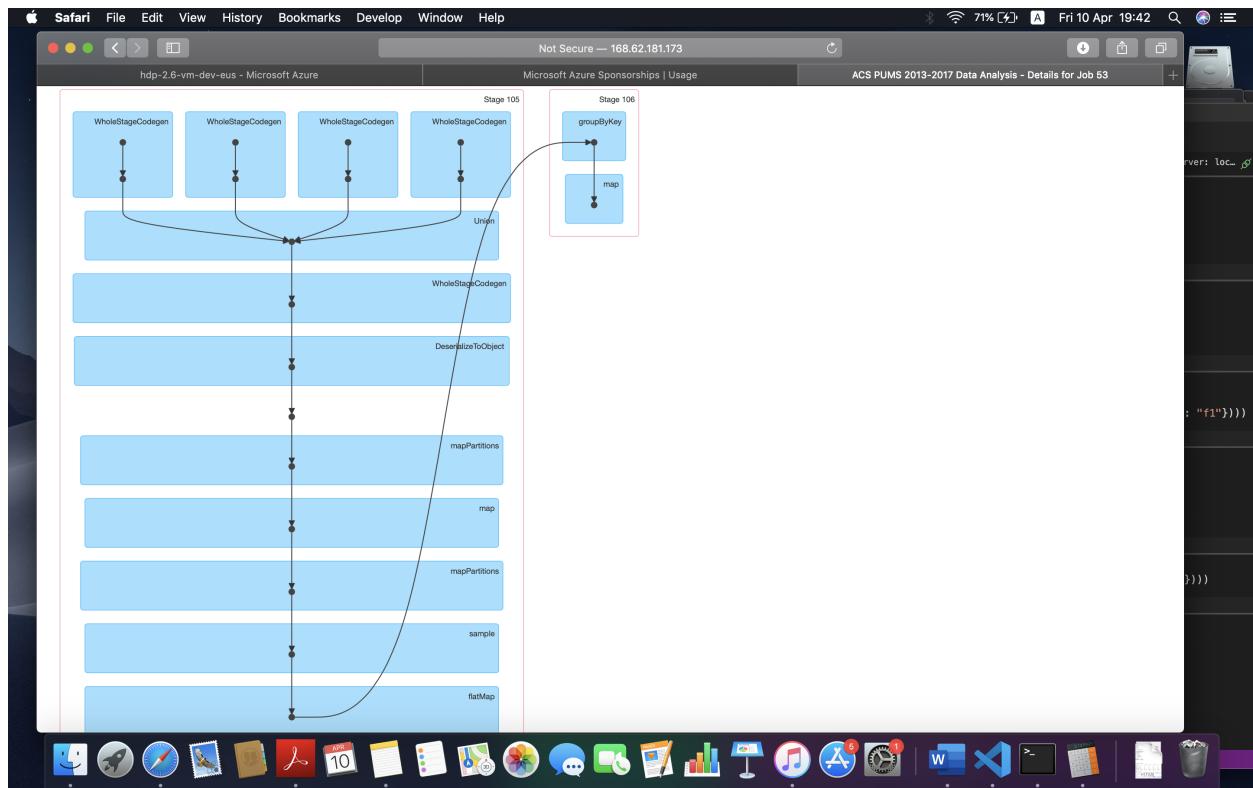


Figure 12: DAG - stringIndexer

The final figure 13 below describes the DAG for a ML algorithm task we used, involving multiple `map` and `flatMap` operations with `mapPartitions` carried out for the model building.

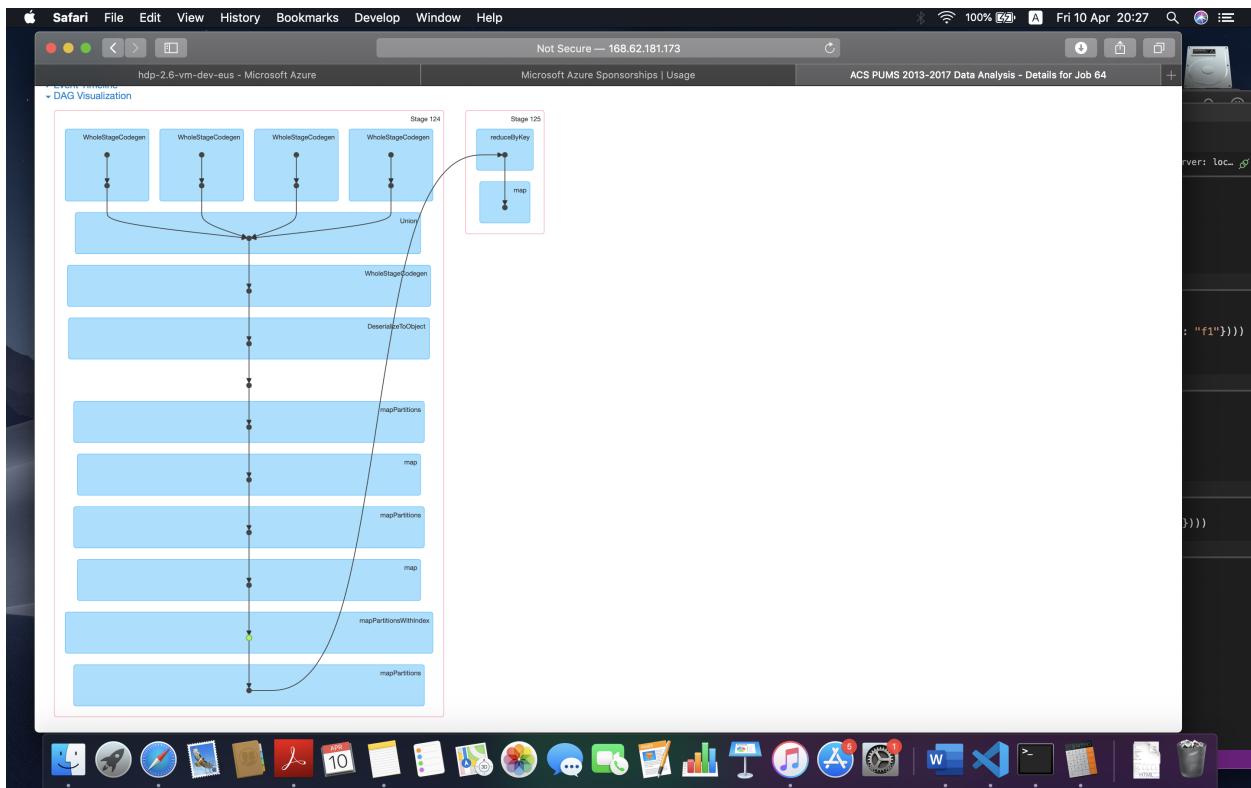


Figure 13: DAG - Final

## Appendix

This section includes the utility commands used such as those for file/data transfers and other miscellaneous stuff.

[GitHub](#) repository for this project can be found here.

Data can be found [here](#)

- Installing Spark standalone

```
wget https://www.apache.org/dyn/closer.lua/spark/spark-2.4.5/spark-2.4.5-bin-hadoop2.7.tgz
```

```
tar xzvf spark-2.4.5-bin-hadoop2.7.tgz
```

- Data transfer from local system to remote VM

```
scp -r <local_dir> <user>@<remote_host>:<remote_dir>
```

Below is the python implementation of the concepts addressed in above sections.

```

1 from IPython import get_ipython
2
3 # %%
4 from pyspark.sql import SparkSession
5 from pyspark import SparkContext
6 from pyspark.sql import functions as f
7 from pyspark.ml.feature import OneHotEncoderEstimator, StringIndexer,
8     VectorAssembler
9 from pyspark.ml import Pipeline
10 from pyspark.ml.classification import DecisionTreeClassifier,
11     RandomForestClassifier
12 from pyspark.ml.evaluation import MulticlassClassificationEvaluator
13 from matplotlib import pyplot as plt
14 from pyspark_dist_explore import hist, Histogram, distplot
15 import seaborn as sns
16 import pandas as pd
17
18 # %%
19 spark = SparkSession.builder.master(
20     "local[*]").appName("ACS PUMS 2013-2017 Data Analysis").getOrCreate()
21
22 # %%
23 # loading the dataset as DataFrame API, one at a time cleaning and removing it
24 # until next load
24 dataA = spark.read.csv('data/csv_pus/csv_pus/psam_pusa.csv',
25                         header=True, inferSchema=True)
26 dataB = spark.read.csv('data/csv_pus/csv_pus/psam_pusb.csv',
27                         header=True, inferSchema=True)
28 dataC = spark.read.csv('data/csv_pus/csv_pus/psam_pusc.csv',
29                         header=True, inferSchema=True)
30 dataD = spark.read.csv('data/csv_pus/csv_pus/psam_pusd.csv',
31                         header=True, inferSchema=True)
32
33 # %%
34 # selecting the columns we need - transformation
35 cols = ['REGION', 'ST', 'AGEP', 'CIT', 'CITWP', 'COW', 'DDRS', 'DEAR', 'DEYE', ,
36         'DREM', 'ENG', 'FER', 'GCL', 'HINS1', 'HINS2', 'HINS3', 'HINS4', 'HINS5', 'HINS6',
37         ',',
38         'HINS7', 'JWMNP', 'JWRIP', 'JWTR', 'LANX', 'MAR', 'RELP', 'SCH', 'SCHG', 'SCHL',
39         'SEX', 'FOD1P', 'FOD2P', 'INDP', 'JWAP', 'JWDP', 'LANP', 'POVPIP', 'POWSP']
40
41 dataA = dataA.select(*cols)
42 dataB = dataB.select(*cols)
43 dataC = dataC.select(*cols)
44 dataD = dataD.select(*cols)

```

```

42
43
44 # %%
45 # load data dictionary
46 # join with data Dictionary only during plotting as legend
47 dataDict = spark.read.csv(
48     'data/csv_pus/csv_pus/PUMS_Data_Dictionary_2013-2017.csv', header=True,
49     inferSchema=True)
50
51 # %%
52 dataDict = dataDict.dropna(how='any')
53 dataDict = dataDict.select('RT', 'Record Type', '_c6')
54
55 # %%
56 dataDict = dataDict.filter(f.col('RT').isin(cols))
57
58 # %%
59 dataDict = dataDict.withColumnRenamed('RT', 'colName')
60 dataDict = dataDict.withColumnRenamed('Record Type', 'code')
61 dataDict = dataDict.withColumnRenamed('_c6', 'abbrev')
62 dataDict.show()
63
64 # %%
65 # get the session details
66 # spark.sparkContext.getConf().getAll()
67 dataDict.filter(dataDict.colName == 'SEX').show()
68
69 # %%
70 # get the spark webUI URL
71 spark.sparkContext.uiWebUrl
72
73 # %%
74 # merge all the dataframes into one dataframe
75 df = dataA.union(dataB)
76 df = df.union(dataC)
77 df = df.union(dataD)
78
79 # %%
80 # clearing the dataframes to free up the memory
81 dataA.unpersist(True)
82 dataB.unpersist(True)
83 dataC.unpersist(True)
84 dataD.unpersist(True)
85
86 # %%
87 # finding the insights on data,
88 # the vision etc difficulties demographics and mostly on high age groups
89 # Create some selections on this data

```

```

89 filtered_by_gender_m = df.filter(f.col('SEX') == '1').select('INDP').join(dataDict
90     ,
91     dataDict.colName == 'INDP', dataDict.code.cast('int') == df.INDP], 'inner').
92     select(f.col('abbrev').alias('INDUSTRY_m'))
93 filtered_by_gender_f = df.filter(f.col('SEX') == '2').select('INDP').join(dataDict
94     ,
95     dataDict.colName == 'INDP', dataDict.code.cast('int') == df.INDP], 'inner').
96     select(f.col('abbrev').alias('INDUSTRY_f'))
97
98 filtered_by_age_50_plus = df.filter(f.col('AGEP') > 50).select('COW').join(
99     dataDict,
100    [
101      dataDict.colName == 'COW', dataDict.code.cast('int') == df.COW], 'inner').
102      select(f.col('abbrev').alias('ClassOfWorker_50_plus'))
103 filtered_by_age_50_minus = df.filter(f.col('AGEP') <= 50).select('COW').join(
104     dataDict,
105    [
106      dataDict.colName == 'COW', dataDict.code.cast('int') == df.COW], 'inner').
107      select(f.col('abbrev').alias('ClassOfWorker_50_minus'))
108
109 # %%
110 # aggregating the data for easier and memory efficient pandas plotting
111 vizM = filtered_by_gender_m.groupBy(
112     'INDUSTRY_m').count().sort('count', ascending=False)
113 vizF = filtered_by_gender_f.groupBy(
114     'INDUSTRY_f').count().sort('count', ascending=False)
115
116 viz50Minus = filtered_by_age_50_minus.groupBy(
117     'ClassOfWorker_50_minus').count().sort('count', ascending=False)
118 viz50Plus = filtered_by_age_50_plus.groupBy(
119     'ClassOfWorker_50_plus').count().sort('count', ascending=False)
120
121 # %%
122 # panadas plotting visualizations
123 C = vizM.toPandas()
124 D = vizF.toPandas()
125
126 pd.concat({
127     'Industry M': C.set_index('INDUSTRY_m'), 'Industry F': D.set_index('INDUSTRY_f'
128     ),
129 }, axis=1).plot.barch()
130
131 # %%
132 A = viz50Minus.toPandas()
133 B = viz50Plus.toPandas()
134
135 pd.concat({
136     '50 Plus COW': B.set_index('ClassOfWorker_50_plus'), '50 Minus COW': A.
137     set_index('ClassOfWorker_50_minus')
138 }, axis=1).plot.barch()

```

```

127
128 # %%
129 filtered_by_gender_m.groupBy('INDUSTRY_m').count().sort(
130     'count', ascending=False).show()
131
132 # %%
133 filtered_by_gender_f.groupBy('INDUSTRY_f').count().sort(
134     'count', ascending=False).show()
135
136 # %%
137 filtered_by_age_50_minus.groupBy('ClassOfWorker_50_minus').count().sort(
138     'count', ascending=False).show()
139
140 # %%
141 filtered_by_age_50_plus.groupBy('ClassOfWorker_50_plus').count().sort(
142     'count', ascending=False).show()
143
144 # %%
145 # determining the industry of work using Machine Learning
146 #features = [AGEP, SCH, SCHG, SCHL, SEX, FOD1P, FOD2P, COW]
147 #target = INDP
148
149 # AGEP > 17 takes care of COW and Schooling
150 # not doing join with the dataDict as during one hot encoding the column names
151 # might have spaces and doesn't make sense, can use join during plotting and
152 # graphing
153 dfX = df.filter((f.col('AGEP') > 17) & (f.col('INDP').isNotNull())).select(
154     'AGEP', 'SCHL', 'SEX', 'FOD1P', 'FOD2P', 'COW', 'INDP')
155 dfX = dfX.fillna({'FOD1P': 0, 'FOD2P': 0})
156
157 # %%
158 dfX.show()
159
160 # onehotencoding of categorical variables for classification
161 categCol = ['SCHL', 'SEX', 'FOD1P', 'FOD2P', 'COW']
162 numericCol = ['AGEP']
163 stages = []
164
165 # StringIndexing and OneHotEncoding of categorical features
166 for cols in categCol:
167     stringIndexer = StringIndexer(inputCol=cols, outputCol=cols+'Index')
168     encoder = OneHotEncoderEstimator(
169         inputCols=[stringIndexer.getOutputCol()], outputCols=[cols + 'classVec'])
170     stages += [stringIndexer, encoder]
171
172 # String indexing of label
173 labelStringIndexer = StringIndexer(inputCol='INDP', outputCol='label')

```

```

173 stages += [labelStringIndexer]
174
175 assemblerInput = [c + 'classVec' for c in categCol] + numericCol
176 assembler = VectorAssembler(inputCols=assemblerInput, outputCol="features")
177 stages += [assembler]
178
179 # %%
180 # creating a ML pipeline
181 pipeline = Pipeline(stages=stages)
182 pipelineModel = pipeline.fit(dfX)
183 dfX = pipelineModel.transform(dfX)
184 dfX = dfX.select('label', 'features')
185
186 # %%
187 dfX.show()
188
189 # %%
190 # train test split
191 train, test = dfX.randomSplit([0.8, 0.2], seed=2020)
192 print("train count: " + str(train.count()))
193 print("test count: " + str(test.count()))
194
195 # %%
196 dctClf = DecisionTreeClassifier(
197     featuresCol='features', labelCol='label', maxDepth=3)
198 dctModel = dctClf.fit(train)
199 predictions = dctModel.transform(test)
200 predictions.select('prediction', 'probability').show(10)
201
202 # %%
203 # peerformance metric, evaluation
204 evaluator = MulticlassClassificationEvaluator()
205 print("Area under Curve ROC: " +
206       str(evaluator.evaluate(predictions, {evaluator.metricName: "f1"})))
207
208 # %%
209 # RandomForesrt classifier
210 rfClf = RandomForestClassifier(featuresCol='features', labelCol='label')
211 rfModel = rfClf.fit(train)
212 predictions = rfModel.transform(test)
213 predictions.select('prediction', 'probability').show(10)
214
215 # %%
216 evaluator = MulticlassClassificationEvaluator()
217 print("Accuracy: " + str(evaluator.evaluate(predictions,
218                               {evaluator.metricName: "accuracy"})))
219
220 # %%

```