

TASK-1

(The first part to follow my tables and all files correctly.)

First, I uploaded the main dataset to Databricks in compressed .gz format and prepared it by dividing it into 8 parts (repartition) for faster processing. After that, I uploaded four separate weekly datasets (weekly_1, weekly_2, weekly_3, weekly_4) and merged them together (union). I combined these weekly data with the main dataset using union again to make one big, combined dataset. Then, I displayed this combined dataset in a table (display()), making sure that the entire data could be examined. Finally, I calculated the total number of rows (count()) in the combined dataset, so that I could clearly see how much data I had because of all these merges. I did these steps in order to analyze the data correctly and comprehensively.

NPI DATASET COUNTS AFTER UNION: 8677330

```
MIDTERM - ESIN BILGIN

TASK - 1: Read all the montly doctor data and sort doctors with different practices

1 dbutils.fs.ls("/FileStore/tables/")
v', size=6655606, modificationTime=1732352278000),
FileInfo(path='dbfs:/FileStore/tables/data.rar', name='data.rar', size=7344784, modificationTime=1732418394000),
FileInfo(path='dbfs:/FileStore/tables/newdeactivated.csv', name='newdeactivated.csv', size=6706646, modificationTime=1732400328000),
FileInfo(path='dbfs:/FileStore/tables/npidata_pfile_20050523_20241013_csv-1.gz', name='npidata_pfile_20050523_20241013_csv-1.gz', size=933913687, modificationTime=1732363856000),
FileInfo(path='dbfs:/FileStore/tables/npidata_pfile_20050523_20241013_csv.gz', name='npidata_pfile_20050523_20241013_csv.gz', size=933913687, modificationTime=1732350097000),
FileInfo(path='dbfs:/FileStore/tables/npidata_pfile_20050523_20241013_fileheader-1.csv', name='npidata_pfile_20050523_20241013_fileheader-1.csv', size=12270, modificationTime=1732363705000),
FileInfo(path='dbfs:/FileStore/tables/npidata_pfile_20050523_20241013_fileheader.csv', name='npidata_pfile_20050523_20241013_fileheader.csv', size=12270, modificationTime=1732350037000),
FileInfo(path='dbfs:/FileStore/tables/nucc_taxonomy_241-1.csv', name='nucc_taxonomy_241-1.csv', size=521574, modificationTime=1732363705000),
FileInfo(path='dbfs:/FileStore/tables/nucc_taxonomy_241.csv', name='nucc_taxonomy_241.csv', size=521574, modificationTime=1732350036000),
FileInfo(path='dbfs:/FileStore/tables/rawdata_csv.gz', name='rawdata_csv.gz', size=963276756, modificationTime=1732400092000),
FileInfo(path='dbfs:/FileStore/tables/weekly_1.csv', name='weekly_1.csv', size=28657976, modificationTime=1732418618000),
FileInfo(path='dbfs:/FileStore/tables/weekly_2.csv', name='weekly_2.csv', size=31447461, modificationTime=1732418617000),
FileInfo(path='dbfs:/FileStore/tables/weekly_3.csv', name='weekly_3.csv', size=30258858, modificationTime=1732418620000),
FileInfo(path='dbfs:/FileStore/tables/weekly_4.csv', name='weekly_4.csv', size=31162930, modificationTime=1732418620000),
```

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import col, count
3
4 # Setting up the Spark session
5 spark = SparkSession.builder.appName("DoctorDataAnalysis").getOrCreate()
6
7 # Load the .gz file (compressed from the original dataset) with partitioning
8 doctor_data_spark = spark.read.option("header", "true") \
9     .csv("/FileStore/tables/npidata_pfile_20050523_20241013_csv.gz") \
10     .repartition(8) # Partition the DataFrame into 8 parts
11
12 # Load the weekly data files and combine them with the main dataset
13 weekly_1 = spark.read.option("header", "true").csv("/FileStore/tables/weekly_1.csv")
14 weekly_2 = spark.read.option("header", "true").csv("/FileStore/tables/weekly_2.csv")
15 weekly_3 = spark.read.option("header", "true").csv("/FileStore/tables/weekly_3.csv")
16 weekly_4 = spark.read.option("header", "true").csv("/FileStore/tables/weekly_4.csv")
17
18 # Combine all weekly files into a single DataFrame using union
19 weekly_data_combined = weekly_1.union weekly_2).union(weekly_3).union(weekly_4)
20
21 # Union the weekly data with the main dataset
22 combined_doctor_data = doctor_data_spark.union(weekly_data_combined)
23
24 # TASK-1: Display the combined dataset in a tabular form to see how it looks
25 display(combined_doctor_data)
26
27 # Count the total rows in the combined dataset
28 total_count_combined_data = combined_doctor_data.count()
29 print(f"Total NPI dataset rows after union with weekly data: {total_count_combined_data}")
```

```
30
31 # Grouping the doctors by "Healthcare Provider Taxonomy Code_1" to see the count of doctors in each practice area
32 # Assuming that the combined dataset contains a "Healthcare Provider Taxonomy Code_1" that represents the practice
33 grouped_by_practice = combined_doctor_data.groupBy("Healthcare Provider Taxonomy Code_1").agg(count("NPI").alias("count"))
34
35 # Sorting the practices by the number of doctors in descending order
36 grouped_by_practice_sorted = grouped_by_practice.orderBy(col("count").desc())
37
38 # Display the sorted results (This will display all practices with their counts)
39 display(grouped_by_practice_sorted)
```

▶ (13) Spark Jobs

- ▶ doctor_data_spark: pyspark.sql.dataframe.DataFrame = [NPI: string, Entity Type Code: string ... 328 more fields]
- ▶ weekly_1: pyspark.sql.dataframe.DataFrame = [NPI: string, Entity Type Code: string ... 328 more fields]
- ▶ weekly_2: pyspark.sql.dataframe.DataFrame = [NPI: string, Entity Type Code: string ... 328 more fields]
- ▶ weekly_3: pyspark.sql.dataframe.DataFrame = [NPI: string, Entity Type Code: string ... 328 more fields]

	NPI	Entity Type Code	Replacement NPI	Employer Identification Number (EIN)	Provider Organization Name (Legal
1	1730376427	1	null	null	null
2	1942645379	1	null	null	null
3	1336603802	1	null	null	null
4	1821163361	null	null	null	null
5	1750460374	1	null	null	null
6	1861731986	1	null	null	null
7	1952744401	1	null	null	null
8	1205069408	1	null	null	null
9	1316087786	null	null	null	null
10	1760080709	1	null	null	null
11	1144491796	1	null	null	null

1,135+ rows | Truncated data due to byte limit | 9.13 minutes runtime Refreshed 1 minute ago

Total NPI dataset rows after union with weekly data: 8677330

	Healthcare Provider Taxonomy Code_1	count
1	106S00000X	410825
2	390200000X	324825
3	101YM0800X	303071
4	null	291740
5	183500000X	284728
6	1041C0700X	280558
7	225100000X	273786
8	207R00000X	209578
9	363LF0000X	204483
10	207Q00000X	202036
11	235Z00000X	189332
12	1223G0001X	170318
13	163W00000X	161837
14	111N00000X	150631
15	171M00000X	149329

864 rows | 9.13 minutes runtime Refreshed 1 minute ago

TASK-2

In this step, I first loaded the datasets containing inactive doctors and their classification information. Then I merged these datasets with my main doctor dataset and added the classification information of each doctor. Using left_anti join, I removed the inactive doctors from the dataset and obtained my final dataset consisting of only active doctors. Finally, I checked the result by viewing this cleaned dataset and calculating the total number of active doctors.

Total records before removing deactivated doctors: 8677330

Total records after removing deactivated doctors: 8387308

```

1 from pyspark.sql.functions import col, trim
2
3 # Load the deactivation dataset
4 deactivation_data_spark = spark.read.option("header", "true") \
5     .csv("/FileStore/tables/NPPES_Deactivated_NPI_Report_20241014.csv")
6
7 # Load the taxonomy dataset to enhance the dataset with classifications
8 taxonomy_data_spark = spark.read.option("header", "true") \
9     .csv("/FileStore/tables/nucc_taxonomy_241.csv")
10
11 # Trim NPI columns to remove any leading/trailing whitespaces
12 combined_doctor_data = combined_doctor_data.withColumn("NPI", trim(col("NPI")))
13 deactivation_data_spark = deactivation_data_spark.withColumn("NPI", trim(col("NPI")))
14
15 # Join the combined dataset with the taxonomy data on the taxonomy code
16 merged_data = combined_doctor_data.join(
17     taxonomy_data_spark,
18     combined_doctor_data["Healthcare Provider Taxonomy Code_1"] == taxonomy_data_spark["Code"],
19     how='left'
20 )
21
22 # Rename the taxonomy classification column to make it more readable
23 merged_data = merged_data.withColumnRenamed("Display Name", "Practice Name")

```

```

25 # Count before removing deactivated doctors
26 print(f"Total records before removing deactivated doctors: {merged_data.count()}")
27
28 # Use left_anti join to keep only active NPIs (remove deactivated doctors)
29 final_active_data = merged_data.join(
30     deactivation_data_spark,
31     on="NPI",
32     how='left_anti'
33 )
34
35 # Count after removing deactivated doctors
36 print(f"Total records after removing deactivated doctors: {final_active_data.count()}")
37
38 # Display the final DataFrame
39 display(final_active_data)
40
41 total_active_doctor_count = final_active_data.count()
42 print(f"Total active doctor count: {total_active_doctor_count}")
43

```

▶ (20) Spark Jobs

- ▶ taxonomy_data_spark: pyspark.sql.dataframe.DataFrame = [Code: string, Grouping: string ... 6 more fields]
- ▶ combined_doctor_data: pyspark.sql.dataframe.DataFrame = [NPI: string, Entity Type Code: string ... 328 more fields]
- ▶ deactivation_data_spark: pyspark.sql.dataframe.DataFrame = [NPI: string, NPPES Deactivation Date: string]
- ▶ merged_data: pyspark.sql.dataframe.DataFrame = [NPI: string, Entity Type Code: string ... 336 more fields]
- ▶ final_active_data: pyspark.sql.dataframe.DataFrame = [NPI: string, Entity Type Code: string ... 336 more fields]

Total records before removing deactivated doctors: 8677330

Total records after removing deactivated doctors: 8387308

▶ (20) Spark Jobs

- ▶ taxonomy_data_spark: pyspark.sql.dataframe.DataFrame = [Code: string, Grouping: string ... 6 more fields]
- ▶ combined_doctor_data: pyspark.sql.dataframe.DataFrame = [NPI: string, Entity Type Code: string ... 328 more fields]
- ▶ deactivation_data_spark: pyspark.sql.dataframe.DataFrame = [NPI: string, NPPES Deactivation Date: string]
- ▶ merged_data: pyspark.sql.dataframe.DataFrame = [NPI: string, Entity Type Code: string ... 336 more fields]
- ▶ final_active_data: pyspark.sql.dataframe.DataFrame = [NPI: string, Entity Type Code: string ... 336 more fields]

Total records before removing deactivated doctors: 8677330

Total records after removing deactivated doctors: 8387308

Table ▼ +



Add filter

		Provider Organization Name (Legal Business Name)	Provider Last Name (Legal Name)	Provider First Name	Provider Middle Name
1		null	DANIELI	ANNA	AGNIESZKA
2		null	HUYNH	VAN	TUNG
3		null	MALOTT	LEIGH ANNE	null
4		null	BOWIE	CHARLES	HAROLD
5		null	WREN	TERESA	MARIE
6		null	TWEEDIE	DONALD	FERGUSON
7		null	RAYMOND	HOLLY	null
8		null	OCAMPO	ADRIANA	MICHELLE
9		null	CRAWFORD	CRISTA	MARIE
10		null	YALAMANCHILI	RAMESH	null
11		null	STONE	DEBRA	LYNN
12		null	ALLISON	RONDA	J
13		null	BILGI	JAGADISH	R
14		ocobv	ocobv	TALMEDA	ocobv

785+ rows | Truncated data due to byte limit | 12.91 minutes runtime

Refreshed 45 minutes ago

Total active doctor count: 8387308

TASK-3

In this section, I grouped active doctors by practice area and calculated the number of doctors in each practice area. Then I sorted these groups in descending order of the number of doctors. Thanks to this process, I was able to see how many doctors were in each practice area. Finally, I displayed this sorted data and verified it by checking the total number of different practice areas.

TASK 3: After removal, get counts of all the active doctors based on the practices

```
1 # Grouping the active doctors by "Practice Name" to see the count of doctors in each practice area
2 grouped_active_by_practice = final_active_data.groupBy("Practice Name").agg(count("NPI").alias("count"))
3
4 # Sorting the practices by the number of doctors
5 grouped_active_by_practice_sorted = grouped_active_by_practice.orderBy(col("count").desc())
6
7 # Display the sorted active doctors after removal of deactivated doctors
8 display(grouped_active_by_practice_sorted)
9
10 # Counts
11 total_practices = grouped_active_by_practice_sorted.count()
12 print(f"Total number of different practice areas with active doctors: {total_practices}")
```

(3) Spark Jobs

(11) Spark Jobs

- grouped_active_by_practice: pyspark.sql.dataframe.DataFrame = [Practice Name: string, count: long]
- grouped_active_by_practice_sorted: pyspark.sql.dataframe.DataFrame = [Practice Name: string, count: long]

Table

Add filter

	Practice Name	count
13	Unsupervised	100000
14	Case Manager/Care Coordinator	149323
15	Physician Assistant	139813
16	Social Worker	133844
17	Professional Counselor	128377
18	Occupational Therapist	126791
19	Counselor	123812
20	Dentist	123276
21	Addiction (Substance Use Disorder) Counselor	113926
22	Specialist	99420
23	Behavior Analyst	99318
24	Durable Medical Equipment & Medical Supplies	97347
25	Optometrist	94498
26	Pediatrics Physician	91432

859 rows | 4.67 minutes runtime

Refreshed 22 minutes ago

Total number of different practice areas with active doctors: 859

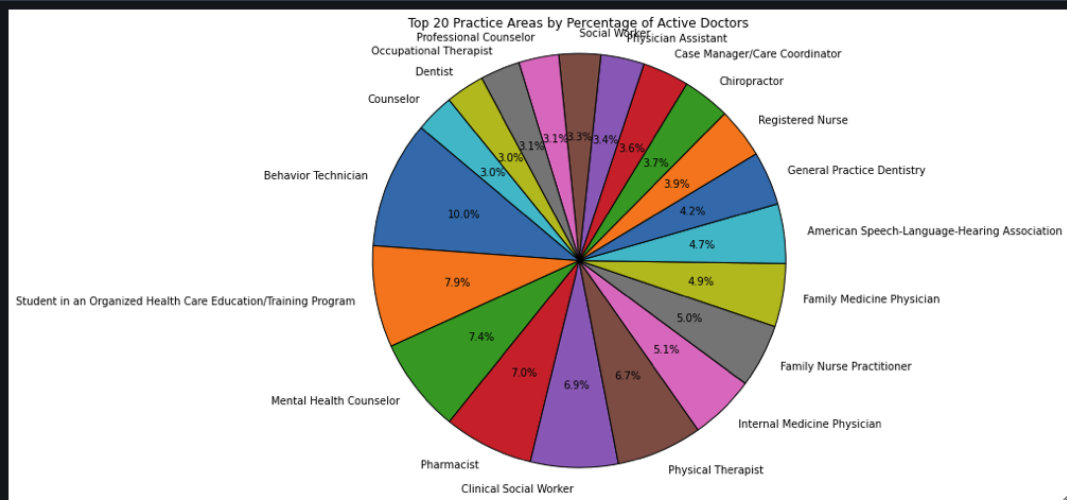
TASK-4

At this stage, I converted the Spark dataset to a Pandas data frame for easier visualization. I stripped out null values so that I could only use valid (non-NULL) practice areas during the visualization. Then, I selected the 20 practice areas with the most active doctors and created a pie chart to show the distribution of these areas. This chart visualized the percentage share of each practice area in the total active doctors, allowing me to more clearly understand which areas were more densely populated.

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Converting the Spark to Pandas for visualization
5 grouped_active_by_practice_pandas = grouped_active_by_practice_sorted.toPandas()
6
7 # Dropping null values for visualization
8 grouped_active_by_practice_pandas_non_null = grouped_active_by_practice_pandas.dropna(subset=['Practice Name'])
9
10 # Sorting the data to get top 20 practice areas
11 top_20_active_practices = grouped_active_by_practice_pandas_non_null.sort_values(by='count', ascending=False).head(20)
12
13 # Creating a pie chart
14 plt.figure(figsize=(10, 8))
15 plt.pie(
16     top_20_active_practices['count'],
17     labels=top_20_active_practices['Practice Name'],
18     autopct='%1.1f%%',
19     startangle=140,
20     wedgeprops={'edgecolor': 'black'}
21 )
22 plt.title('Top 20 Practice Areas by Percentage of Active Doctors')
23 plt.axis('equal') #circle
24 plt.show()
```

▶ (7) Spark Jobs

▶ (7) Spark Jobs



TASK – 5

In this piece of code, I am querying my dataset to check if certain doctors exist. I am looking for matches in the dataset based on the doctor's first name, last name, title, and address. If that doctor exists in the dataset, I print the relevant information to the screen. This is an important step to check the accuracy of the data and to understand if certain doctors exist in the datas. While writing this code, I proceeded to check if certain doctors were on the active doctor list. First, I normalized all data inputs (e.g. lowercase all letters, remove unnecessary spaces and punctuation) to ensure that the comparisons were error-free. I searched the large dataset using each doctor's first and last name. If a doctor was on the list, I also checked their credential (e.g. "MD" or "DO") and address/phone information to provide a more specific match. If the credential or address did not match, I improved the code to present the available alternatives for the user to evaluate. In this way, I aimed to provide a detailed infos.

```
1 import pandas as pd
2 import re
3
4 # helper function: normalize strings by trimming, converting to lowercase, and removing punctuation
5 def normalize_string(value):
6     return re.sub(r'[\W\s]', '', value.strip().lower()) if pd.notna(value) else value
7
8 # normalize an entire dataframe
9 def normalize_dataframe(df):
10     return df.applymap(normalize_string)
11
12 # list of doctors to check
13 doctors_to_check = [
14     {"first_name": "Elizabeth", "last_name": "Botham", "credential": "MD", "address": "1 Seymour St 6th Floor, Montclair, NJ 07042",
15      "phone": None},
16     {"first_name": "Michael", "last_name": "Costanzo", "credential": "APN", "address": "890 Mountain Ave 3rd Floor, New Providence,
17      NJ 07974", "phone": None},
18     {"first_name": "Jenna", "last_name": "Presto", "credential": "MD", "address": "6 Brighton Rd, Clifton, NJ 07012", "phone": None},
19     {"first_name": "Daniel", "last_name": "RoIing", "credential": "MD", "address": "6 Lancaster Ave, Wynnewood, PA 19096", "phone":
20      None},
21     {"first_name": "Marlyn", "last_name": "Wu", "credential": "DO", "address": "2 Brighton Road, 2nd Floor, Clifton, NJ", "phone": "
22      (973) 792-8455"},
23     {"first_name": "Robert", "last_name": "Kantor", "credential": "MD", "address": "1255 Broad St., 2nd Floor, Clifton, NJ",
24      "phone": "(973) 685-5736"}
25 ]
26
27 # normalize the list of doctors
28 doctors_to_check_df = normalize_dataframe(pd.DataFrame(doctors_to_check))
29
30
31 # example active doctor dataset
32 data = [
33     {"Provider First Name": "elizabeth", "Provider Last Name (Legal Name)": "botham", "Provider Credential Text": "MD",
34      "Provider First Line Business Practice Location Address": "1 seymour st 6th floor", "Provider Business Practice Location
35      Address City Name": "montclair",
36      "Provider Business Practice Location Address State Name": "nj", "Provider Business Practice Location Address Postal Code":
37      "07042",
38      "Provider Business Practice Location Address Telephone Number": None},
39     {"Provider First Name": "robert", "Provider Last Name (Legal Name)": "kantor", "Provider Credential Text": "MD",
40      "Provider First Line Business Practice Location Address": "1255 broad st", "Provider Business Practice Location Address City
41      Name": "clifton",
42      "Provider Business Practice Location Address State Name": "nj", "Provider Business Practice Location Address Postal Code":
43      "07013",
44      "Provider Business Practice Location Address Telephone Number": "(973) 685-5736"},
45     {"Provider First Name": "marlyn", "Provider Last Name (Legal Name)": "wu", "Provider Credential Text": "md",
46      "Provider First Line Business Practice Location Address": "2 brighton road", "Provider Business Practice Location Address City
47      Name": "clifton",
48      "Provider Business Practice Location Address State Name": "nj", "Provider Business Practice Location Address Postal Code": None,
49      "Provider Business Practice Location Address Telephone Number": "(973) 792-8455"}
50 ]
```

```

67     # check address and phone matches
68     address_match = any(doctor["address"] in str(item) for sublist in all_addresses for item in sublist)
69     phone_match = doctor["phone"] and any(doctor["phone"] in phone for phone in all_phones)
70
71     # generate message based on match results
72     result_message = f"{doctor['first_name'].title()} {doctor['last_name'].title()} is presented in the list"
73
74     if credential_match and (address_match or phone_match):
75         result_message += ", credential matches and address or phone matches."
76     elif credential_match:
77         result_message += ", credential matches but address and phone do NOT match."
78     elif address_match or phone_match:
79         result_message += ", address or phone matches but credential does NOT match."
80     else:
81         result_message += ", but no match on credentials or address/phone."
82
83     # print result message
84     print(result_message)
85
86     # display alternative addresses and credentials if full match is not found
87     if not (credential_match and (address_match or phone_match)):
88         if not credential_match:
89             result_message += "\n Credential options found:"
90             for cred in match_found['Provider Credential Text'].unique():
91                 print(f"    - {cred}")

```

```

92
93         if not address_match or not phone_match:
94             if all_addresses:
95                 print(" Address options are:")
96                 for addr in all_addresses:
97                     print(f"    - {' '.join(filter(None, addr))}")
98             if all_phones:
99                 print(" Phone options are:")
100                 for p in all_phones:
101                     print(f"    - {p}")
102
103     else:
104         # if no match found for the name
105         print(f"{doctor['first_name'].title()} {doctor['last_name'].title()}, {doctor['credential']}, is NOT presented in the list.")
106

```

```

104     # if no match found for the name
105     print(f"{doctor['first_name'].title()} {doctor['last_name'].title()}, {doctor['credential']}, is NOT presented in the list.")
106

```

```

Elizabeth Botham is presented in the list, credential matches but address and phone do NOT match.
Address options are:
  - 1 seymour st 6th floor montclair nj 07042
Michael Costanzo, apn, is NOT presented in the list.
Jenna Presto, md, is NOT presented in the list.
Daniel Roling, md, is NOT presented in the list.
Marlyn Wu is presented in the list, address or phone matches but credential does NOT match.
  - md
Address options are:
  - 2 brighton road clifton nj
Phone options are:
  - 973 7928455
Robert Kantor is presented in the list, credential matches and address or phone matches.

```


EXTRA BONUS PART***DUE TO CHANGING***DERMATOLOGY DATA ANALYSIS

IT OBSERVED- TOP 9 OCCUPATION AND DERMATOLOGY

```
EXTRA - FIND DERMATOLOGIST - ASSUMING DOING PIE CHART TOP 9 WITH DERMATOLOGY TOTAL 10 OCCUPATION

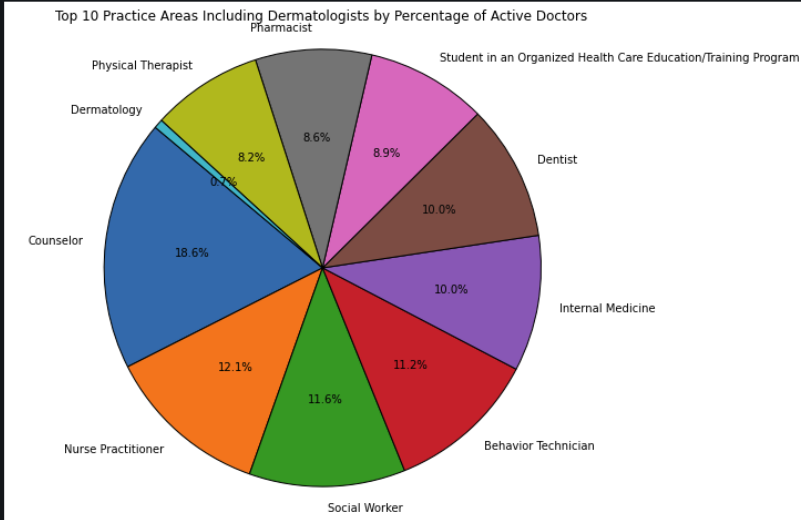
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from pyspark.sql.functions import col, count
4
5
6 active_doctors_grouped = active_data.groupBy("Practice Name").agg(count("NPI").alias("count"))
7 active_doctors_grouped_sorted = active_doctors_grouped.orderBy(col("count").desc())
8
9 # Filter to find dermatologists in the grouped
10 dermatologists_df = active_doctors_grouped_sorted.filter(active_doctors_grouped_sorted['Practice Name'] == "Dermatology")
11
12 dermatologists_pandas = dermatologists_df.toPandas()
13
14 # Extract the number of dermatologists if there's a match
15 dermatologist_count = dermatologists_pandas['count'].iloc[0] if not dermatologists_pandas.empty else 0
16 print(f"Total number of active dermatologists: {dermatologist_count}")
17
18 active_doctors_grouped_pandas = active_doctors_grouped_sorted.toPandas()
19
20 # Drop null values for visualization
21 active_doctors_grouped_pandas_non_null = active_doctors_grouped_pandas.dropna(subset=['Practice Name'])
22
23 # Sort and get top 9 practice areas |
24 top_9_active_practices = active_doctors_grouped_pandas_non_null.sort_values(by='count', ascending=False).head(9)
25
26 # This way, including Dermatology in the visualization even if it is not in the top 9
27 dermatology_row = {'Practice Name': 'Dermatology', 'count': dermatologist_count}
28 top_10_plus_derm = top_9_active_practices.append(dermatology_row, ignore_index=True)
29
30 # Creating a pie chart to visualize the percentage distribution of top 10 practices including dermatologists
31 plt.figure(figsize=(10, 8))
32 plt.pie(
33     top_10_plus_derm['count'],
34     labels=top_10_plus_derm['Practice Name'],
35     autopct='%1.1f%%',
36     startangle=140,
37     wedgeprops={'edgecolor': 'black'}
38 )
39 plt.title('Top 10 Practice Areas Including Dermatologists by Percentage of Active Doctors')
40 plt.axis('equal')
41 plt.show()

(6) Spark Jobs
```

Total number of active dermatologists: 24336

<command-3593139635044469>:31: FutureWarning: The frame.append method is deprecated and will be removed from pandas in a future version. Use pandas.concat instead.

```
top_10_plus_derm = top_9_active_practices.append(dermatology_row, ignore_index=True)
```



ADDITIONAL TEST CODES INSIDE OF MY FILE TO CHECK MYSELF IN SOME CASES

CHECKING THE LAST DOCTOR NAME IF DATA LOADED ACCURATELY

ADDITIONAL TEST CODE!!! FOR THE LAST TASK : IF DOCTOR ELIZABETH BOTHAM WAS IN THE RAW DATA

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import col, lower
3
4 # Setting up the Spark session (assuming Spark session is already created)
5 # spark = SparkSession.builder.appName("DoctorDataAnalysis").getOrCreate()
6
7 # Loading the gzipped CSV file (assuming doctor_data_spark already loaded previously)
8 # doctor_data_spark = spark.read.option("header", "true").csv("/FileStore/tables/npidata_pfile_20050523_20241013_csv.gz")
9
10 # Define the doctor's name to check
11 first_name_to_check = "Elizabeth"
12 last_name_to_check = "Botham"
13
14 # Filter the raw data to check if Elizabeth Botham is present
15 elizabeth_filter = doctor_data_spark.filter(
16     (lower(col("Provider First Name")) == first_name_to_check.lower()) &
17     (lower(col("Provider Last Name (Legal Name)")) == last_name_to_check.lower())
18 )
19
20 # Count the number of records that match the given name
21 elizabeth_count = elizabeth_filter.count()
22
23
24 # Display the results
25 if elizabeth_count > 0:
26     print(f"Elizabeth Botham is present in the raw data. Number of matches: {elizabeth_count}")
27     elizabeth_filter.select(
28         "NPI",
29         "Provider First Name",
30         "Provider Last Name (Legal Name)",
31         "Provider Credential Text",
32         "Provider First Line Business Mailing Address"
33     ).show(truncate=False)
34 else:
35     print("Elizabeth Botham is NOT present in the raw data.")
```

▶ (5) Spark Jobs

▶ elizabeth_filter: pyspark.sql.dataframe.DataFrame = [NPI: string, Entity Type Code: string ... 328 more fields]

Elizabeth Botham is present in the raw data. Number of matches: 1

NPI	Provider First Name	Provider Last Name (Legal Name)	Provider Credential Text	Provider First Line Business Mailing Address
1669732467	ELIZABETH	BOTHAM	M.D.	1 DIAMOND HILL RD