



# US campaign Finance Analytics

By  
Harshitha Yentrapragada  
Shivani Badinehal  
Udveg Jukanti

**Problem:**

For analysis of expenditure for contesting in an election, most of the new candidates cannot analyze the raw data provided by the FEC(Federal Election Commission).

**Solution:**

By analyzing the data files of previous years, and generating curated data tables which will indeed create reports that describe overall understandings of expenditures

## Data Source

**Retrieved from :**

US Campaign Finance Data – Federal Election Commission(FEC)

**URL:**

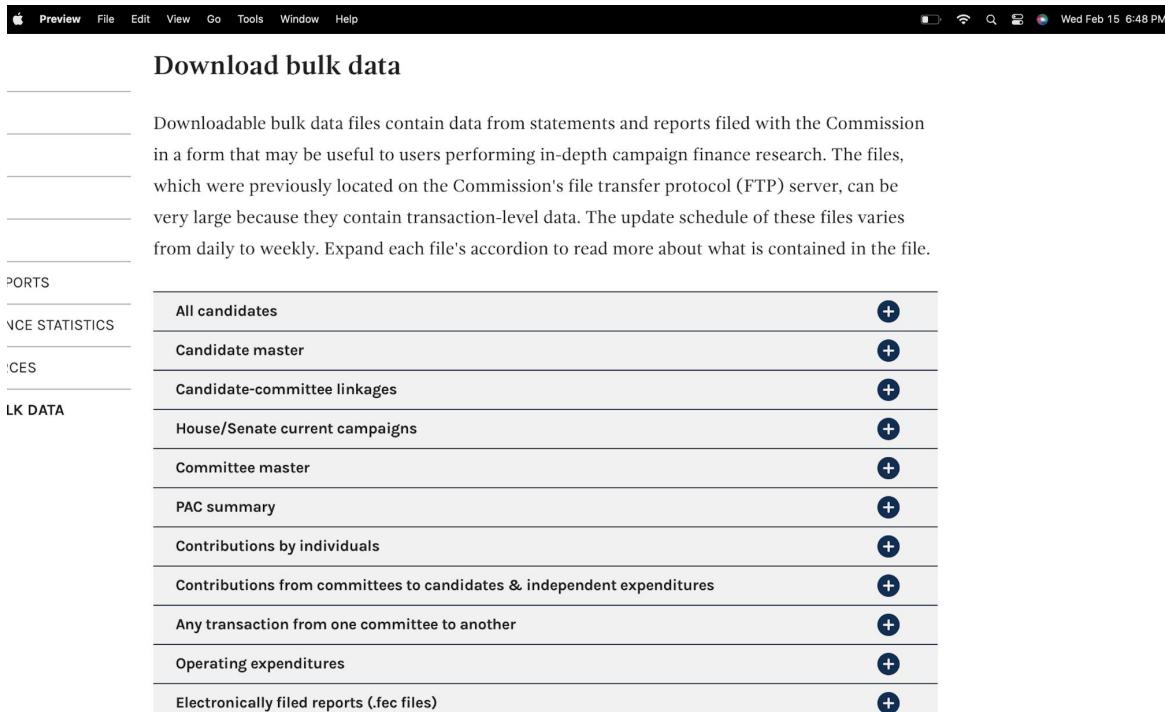
<https://www.fec.gov/data/browse-data/?tab=bulk-data>

**Description:**

- The size of the dataset is approximately 15 gb and has almost 28 million people information.
- This dataset includes info of candidates, committees, Political action committees, House and Senate campaigns, as well as transactions, individual and committee contributions, operations and independent expenditures for US elections.



# Files Used



The screenshot shows a Mac OS X desktop with a file browser window open. The window title is "Download bulk data". The menu bar at the top includes "Preview", "File", "Edit", "View", "Go", "Tools", "Window", and "Help". On the right side of the menu bar, there are icons for battery, signal strength, and a clock showing "Wed Feb 15 6:48 PM". The main content area of the window displays a list of downloadable bulk data files. The list is organized into sections with horizontal lines: "PORTS", "NCE STATISTICS", "ICES", and "LK DATA". Each section contains a list of files with a "+" icon to the right of each item. The "LK DATA" section is expanded, showing the following items:

- All candidates
- Candidate master
- Candidate-committee linkages
- House/Senate current campaigns
- Committee master
- PAC summary
- Contributions by individuals
- Contributions from committees to candidates & independent expenditures
- Any transaction from one committee to another
- Operating expenditures
- Electronically filed reports (.fec files)

## Fec Data Column Names and Abbreviations

<https://www.fec.gov/campaign-finance-data/committee-summary-file-description/>

```
'CAND_ID': 'Candidate identification',
'CAND_NAME': 'Candidate name',
'CAND_ICI': 'Incumbent challenger status',
'PTY_CD': 'Party code',
'CAND_PTY_AFFILIATION': 'Party affiliation',
'TTL_RECEIPTS': 'Total receipts',
'TRANS_FROM_AUTH': 'Transfers from authorized committees',
'TTL_DISB': 'Total disbursements',
'TRANS_TO_AUTH': 'Transfers to authorized committees',
'COH_BOP': 'Beginning cash',
'COH_COP': 'Ending cash',
'CAND CONTRIB': 'Contributions from candidate',
'CAND LOANS': 'Loans from candidate',
'OTHER LOANS': 'Other loans',
'CAND LOAN REPAY': 'Candidate loan repayments',
'OTHER LOAN REPAY': 'Other loan repayments',
'DEBTS_OWED_BY': 'Debts owed by',
'TTL_INDIV_CONTRIB': 'Total individual contributions',
'CAND OFFICE_ST': 'Candidate state',
'CAND OFFICE_DISTRICT': 'Candidate district',
'SPEC_ELECTION': 'Special election status',
'PRIM_ELECTION': 'Primary election status',
'RUN_ELECTION': 'Runoff election status',
'GEN_ELECTION': 'General election status',
'GEN_ELECTION_PRECENT': 'General election percentage',
'OTHER_POL_CMTE CONTRIB': 'Contributions from other political committees',
'POL_PTY CONTRIB': 'Contributions from party committees',
'CVG_END_DT': 'Coverage end date',
'INDIV_REFUNDS': 'Refunds to individuals',
'CMTE_REFUNDS': 'Refunds to committees'}
```

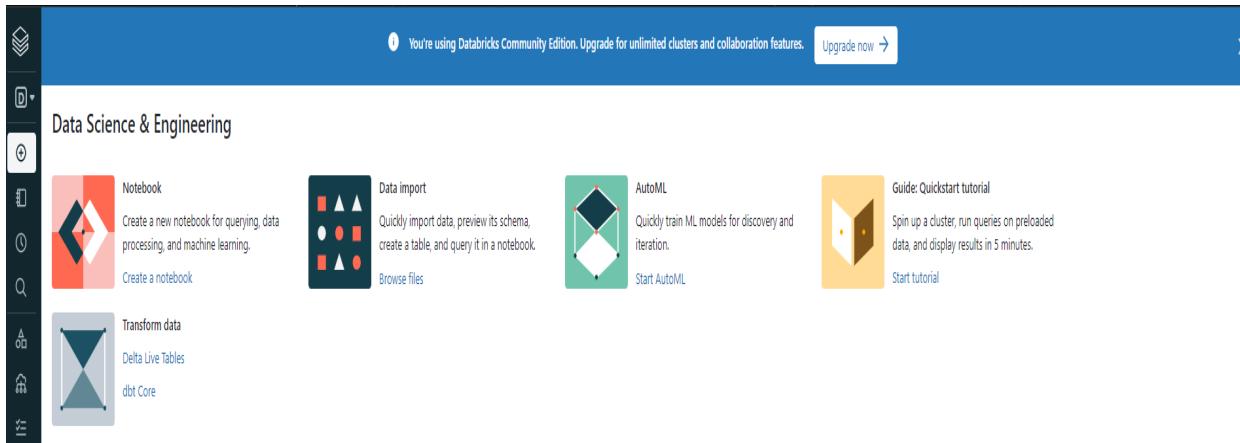
## Software Tools and Technologies:

- Databricks
- Spark
- Tableau



## Data Bricks

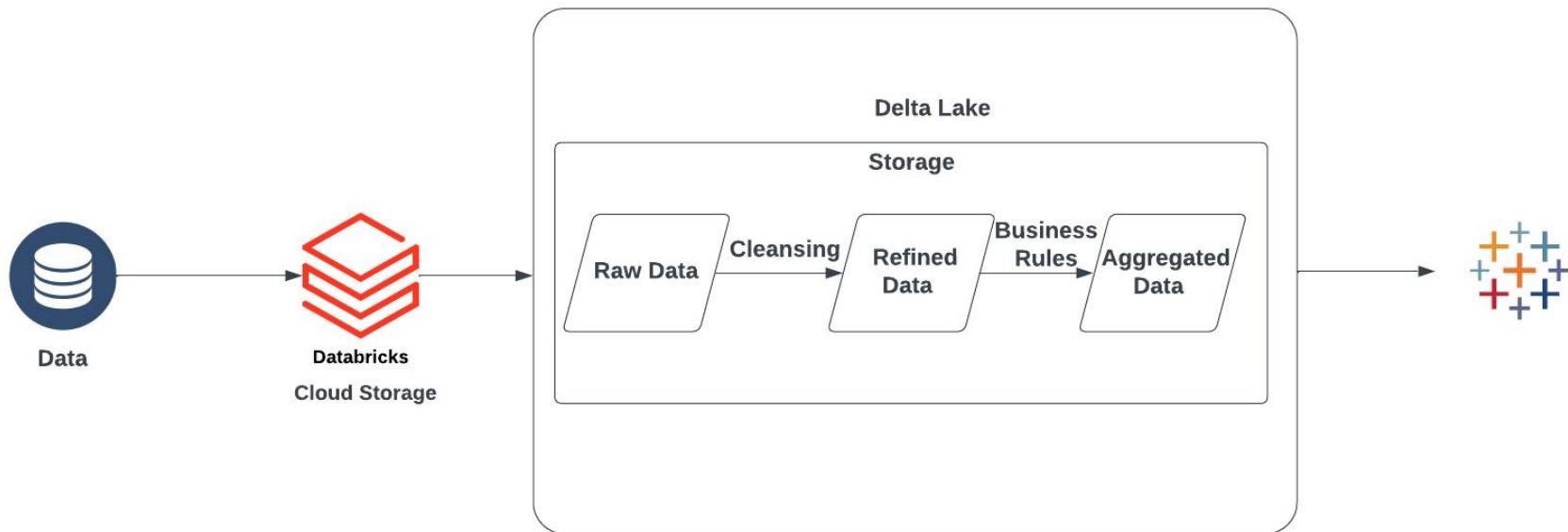
- The AWS, Azure, and Google clouds can all support the use of Databricks to handle enormous volumes of raw, unprocessed data.
- The computing resources and settings that make up a Databricks cluster are used to run data engineering, data science, and data analytics workloads such production ETL pipelines, streaming analytics, ad-hoc analytics, and machine learning.



# Parquet

- The very best encoding and compression techniques are supported by Apache Parquet. When used with modern cloud technologies like Amazon Athena, Redshift Spectrum, BigQuery, and Azure Data Lakes, Apache Parquet reduces the cost of data file storage and increases the efficiency of data queries.
- With the ability to read and write Parquet files, Spark SQL maintains the original data's schema automatically. All columns are automatically changed to be nullable when creating Parquet files for compatibility reasons.
- Reading and writing Parquet files is supported by Spark SQL, which automatically preserves the schema of the original data. All columns are automatically altered to be nullable when producing Parquet files due to compatibility concerns.

# Architecture



# Importing Data

Importing Data Files(.txt) files and .csv(Header files) files into Data Bricks.

The screenshot shows the Databricks Data browser interface. The left sidebar has a dark theme with a navigation menu:

- Create
- Workspace
- Recents
- Search
- Data** (selected)
- Compute
- Workflows

The main area displays a file tree under the path `/FileStore/tables`. The tree structure includes `jars` and `tables` nodes, with a search bar above them. A large list of files is shown on the right, many of which have dropdown arrows indicating they are expandable:

- cm\_header\_file.csv
- cm.txt
- cn.txt
- cont\_and\_ind\_exp.txt
- cont\_indiv\_header\_file.csv
- import\_illustration\_4151f89c.svg
- Innovating\_with\_Big\_Data\_Hass...
- itho.txt
- oppexp\_header\_file.csv
- oppexp.txt
- oth\_header\_file-1.csv
- oth\_header\_file.csv
- PAC\_summary.txt
- pas2\_header\_file-1.csv
- pas2\_header\_file.csv
- Senate\_current\_campaigns.txt
- US\_Campaign\_Finance\_Analytic...
- weball20\_2\_.zip
- weball20.txt

## Read all the files

- Creating a custom schema for the files which don't have headers

```
1 from pyspark.sql.types import *
2 custom_schema_all = StructType([
3     StructField("CAND_ID", StringType(), True),
4     StructField("CAND_NAME", StringType(), True),
5     StructField("CAND_ICI", StringType(), True),
6     StructField("PTY_CD", IntegerType(), True),
7     StructField("CAND_PTY_AFFILIATION", StringType(), True),
8     StructField("TTL_RECEIPTS", DoubleType(), True),
9     StructField("TRANS_FROM_AUTH", DoubleType(), True),
10    StructField("TTL_DISB", DoubleType(), True),
11    StructField("TRANS_TO_AUTH", DoubleType(), True),
12    StructField("COH_BOP", DoubleType(), True),
13    StructField("COH_COP", DoubleType(), True),
14    StructField("CAND_CONTRIB", DoubleType(), True),
15    StructField("CAND_LOANS", DoubleType(), True),
16    StructField("OTHER_LOANS", DoubleType(), True),
17    StructField("CAND_LOAN_REPAY", DoubleType(), True),
18    StructField("OTHER_LOAN_REPAY", DoubleType(), True),
19    StructField("DEBTS_OWE_BY", DoubleType(), True),
20    StructField("TTL_INDIV_CONTRIB", DoubleType(), True),
21    StructField("CAND_OFFICE_ST", StringType(), True),
22    StructField("CAND_OFFICE_DISTRICT", StringType(), True),
23    StructField("SPEC_ELECTION", StringType(), True),
24    StructField("PRIM_ELECTION", StringType(), True),
25    StructField("RUN_ELECTION", StringType(), True),
26    StructField("GEN_ELECTION", StringType(), True),
27    StructField("GEN_ELECTION_PRECENT", StringType(), True),
28    StructField("OTHER_POL_CMTE CONTRIB", DoubleType(), True),
29    StructField("POL_PTY CONTRIB", DoubleType(), True),
30    StructField("CVG_END_DT", StringType(), True),
31    StructField("INDIV_REFUNDS", DoubleType(), True),
32    StructField("CMTE_REFUNDS", DoubleType(), True)
33 ])
```



# Reading CSV files

```
from pyspark.sql.types import StructType

read_format = 'csv'
load_path = 'dbfs:/FileStore/tables/weball20.txt'

# Load the data from its source.
all = (spark.read.format("csv").schema(custom_schema_all).option("header", True).option("delimiter", "|").load("dbfs:/FileStore/tables/weball20.txt"))
```

```
all.show()
```

# Cleaning The Data

# Converting file type and loading

```
write_format = 'parquet'  
save_path = '/FinalProject/all'  
  
# Writing the data in binary format into the source folder.  
all.write \  
    .mode("overwrite") \  
    .format(write_format) \  
    .save(save_path)
```

```
load_path = '/FileStore/tables/candidate_master_header_file.csv'  
cand_master_header = spark \  
    .read \  
    .format(read_format) \  
    .load(load_path , header=True, inferSchema=True)
```

```
cand_master_header.printSchema()  
  
root  
|-- CAND_ID: string (nullable = true)  
|-- CAND_NAME: string (nullable = true)  
|-- CAND_PTY_AFFILIATION: string (nullable = true)  
|-- CAND_ELECTION_YR: string (nullable = true)  
|-- CAND_OFFICE_ST: string (nullable = true)  
|-- CAND_OFFICE: string (nullable = true)  
|-- CAND_OFFICE_DISTRICT: string (nullable = true)  
|-- CAND_ICI: string (nullable = true)  
|-- CAND_STATUS: string (nullable = true)  
|-- CAND_PCC: string (nullable = true)  
|-- CAND_ST1: string (nullable = true)  
|-- CAND_ST2: string (nullable = true)  
|-- CAND_CITY: string (nullable = true)  
|-- CAND_ST: string (nullable = true)  
|-- CAND_ZIP: string (nullable = true)
```



# UMBC

# Checking for null values

databricks

Cleaning (Python)

[Import Notebook](#)

```
import pyspark
def Shape(dataFrame):
    return (dataFrame.count(), len(dataFrame.columns))
pyspark.sql.dataframe.DataFrame.shape = Shape
print(All_candidates.shape())

(3981, 30)
```

```
# Distinct checks whether the dataframe has duplicate values
print("distinct count: " + str(All_candidates.distinct().count()))

distinct count: 3981
```

```
from pyspark.sql.functions import col,isnan, when, count
```

```
All_candidates.select([count(when(isnan(i) | col(i).isNull(), i)).alias(i) for i in All_candidates.columns[:10]]).show()
All_candidates.select([count(when(isnan(i) | col(i).isNull(), i)).alias(i) for i in All_candidates.columns[10:20]]).show()
All_candidates.select([count(when(isnan(i) | col(i).isNull(), i)).alias(i) for i in All_candidates.columns[20:30]]).show()
```

```
+-----+-----+-----+-----+-----+-----+-----+-----+
|CAND_ID|CAND_NAME|CAND_ICI|PTY_CD|CAND_PTY_AFFILIATION|TTL_RECEIPTS|TRANS_FROM_AUTH|TTL_DISB|TRANS_TO_AUTH|COH_BOP|
+-----+-----+-----+-----+-----+-----+-----+-----+
|      0|       0|     78|      0|           1|        0|        0|        0|        0|       0|
+-----+-----+-----+-----+-----+-----+-----+-----+
```

```
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|COH_COP|CAND_CONTRIB|CAND_LOANS|OTHER_LOANS|CAND_LOAN_REPAY|OTHER_LOAN_REPAY|DEBTS_OWED_BY|TTL_INDIV_CONTRIB|CAND_OFFICE_ST|CAND_OFFICE_DISTRICT|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      0|       0|       0|       0|       0|       0|       0|       0|       0|       0|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
```

```
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|SPEC_ELECTION|PRIM_ELECTION|RUN_ELECTION|GEN_ELECTION|GEN_ELECTION_PRECENT|OTHER_POL_CMTE CONTRIB|POL_PTY CONTRIB|CVG_END_DT|INDIV_REFUNDS|CMTE_REFUNDS|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      3981|      3981|      3981|      3981|      3981|        0|        0|        0|        0|       0|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
```

# Dropping Null Values

```
# Dropping the columns which have null values
null_columns=["SPEC_ELECTION","PRIM_ELECTION","RUN_ELECTION","GEN_ELECTION","GEN_ELECTION_PRECENT"]
All_candidates=All_candidates.drop(*null_columns)
len(All_candidates.columns)
```

```
Out[5]: 25
```

```
display(All_candidates.describe())
```

Table

	summary	CAND_ID	CAND_NAME	CAND_ICI	PTY_CD	CAND_PTY_AFFILIATION	TTL_RECEIPTS	TRANS_FROM_AUTH	TTL_DISB
1	count	3981	3981	3903	3981	3980	3981	3981	3981
2	mean	null	null	null	1.6541070082893745	null	3490903.2603391143	199196.3327103742	3182462.0778698847
3	stddev	null	null	null	0.6416658107355816	null	8.233354043721126E7	5849210.709572181	6.819973962700449E7
4	min	H0AL01055	ANGELA GLASS	C	1	AMP	-674132.5	0.0	-674157.5
5	max	S8WV00143	ZUNKER, TRICIA	O	3	WFP	4.824617973E9	2.7517652812E8	3.848951085E9

Showing all 5 rows.

# Replacing the Values

```
TTL_RECEIPTS_correct = when(((All_candidates['TRANS_FROM_AUTH'] != 0.0) & (All_candidates['TRANS_TO_AUTH'] != 0.0)), ((All_candidates['TTL_RECEIPTS'] -  
All_candidates['TRANS_FROM_AUTH'])))  
All_candidates = All_candidates.withColumn('TTL_RECEIPTS_CORRECTED', TTL_RECEIPTS_correct)  
All_candidates.show()  
  
Out[8]: (3981, 26)
```

```
from pyspark.sql.functions import when  
from pyspark.sql.functions import regexp_replace  
  
PAC_file=PAC_file.withColumn('CMTE_DSGN',  
    when((PAC_file.CMTE_DSGN == 'A'),regexp_replace(PAC_file.CMTE_DSGN,'A','Authorized by a candidate')) \  
.when((PAC_file.CMTE_DSGN == 'B'),regexp_replace(PAC_file.CMTE_DSGN,'B','Lobbyist/Registrant PAC')) \  
.when((PAC_file.CMTE_DSGN == 'D'),regexp_replace(PAC_file.CMTE_DSGN,'D','Leadership PAC')) \  
.when((PAC_file.CMTE_DSGN == 'J'),regexp_replace(PAC_file.CMTE_DSGN,'J','Joint fundraiser')) \  
.when((PAC_file.CMTE_DSGN == 'P'),regexp_replace(PAC_file.CMTE_DSGN,'P','Principal campaign committee of a candidate')) \  
.when((PAC_file.CMTE_DSGN == 'U'),regexp_replace(PAC_file.CMTE_DSGN,'U','Unauthorized')))  
  
PAC_file.select("CMTE_DSGN").distinct().show()  
  
+-----+  
| CMTE_DSGN |  
+-----+  
| Unauthorized |  
| Lobbyist/Registrant PAC |  
| null |  
| Leadership PAC |  
| Principal campaign committee of a candidate |  
| Joint fundraiser |  
+-----+
```

# List of files after loading

```
%fs ls "/FinalProject"
```

Table

	path	name	size	modificationTime
1	dbfs:/FinalProject/All_Candidates_Cleaned/	All_Candidates_Cleaned/	0	0
2	dbfs:/FinalProject/Cont_by_indv/	Cont_by_indv/	0	0
3	dbfs:/FinalProject/PAC/	PAC/	0	0
4	dbfs:/FinalProject/all/	all/	0	0
5	dbfs:/FinalProject/cand_master/	cand_master/	0	0
6	dbfs:/FinalProject/ccl/	ccl/	0	0
7	dbfs:/FinalProject/cm/	cm/	0	0

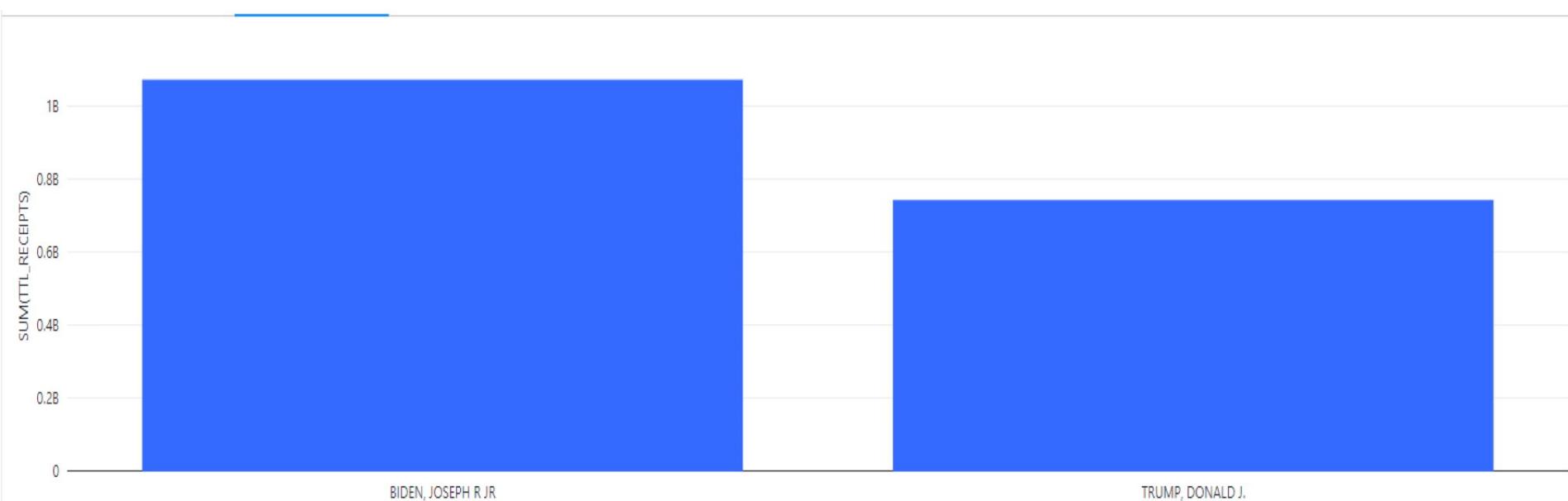
Showing all 13 rows.

# Data Visualizations

Merging the necessary dataframes and creating tables for visualizing the data.

## Visualization For Total Receipts for Each Candidate

- Total Receipts are anything of value (money, goods, services or property) received by a political committee.



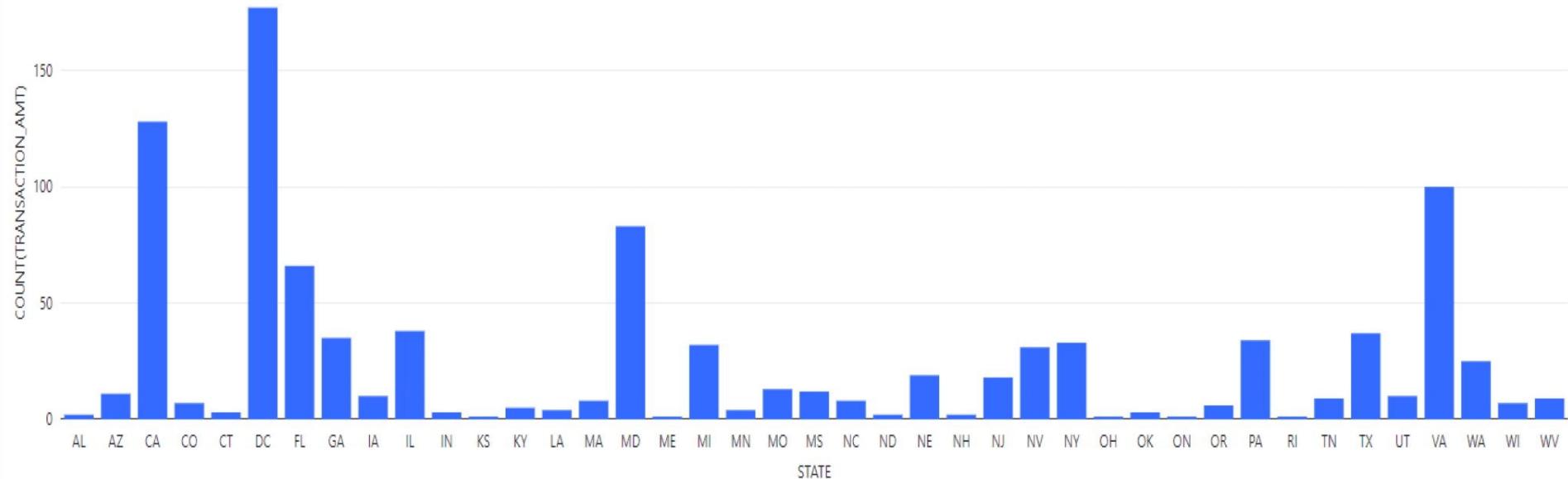
# Transaction amount in each state

Table

Visualization 1

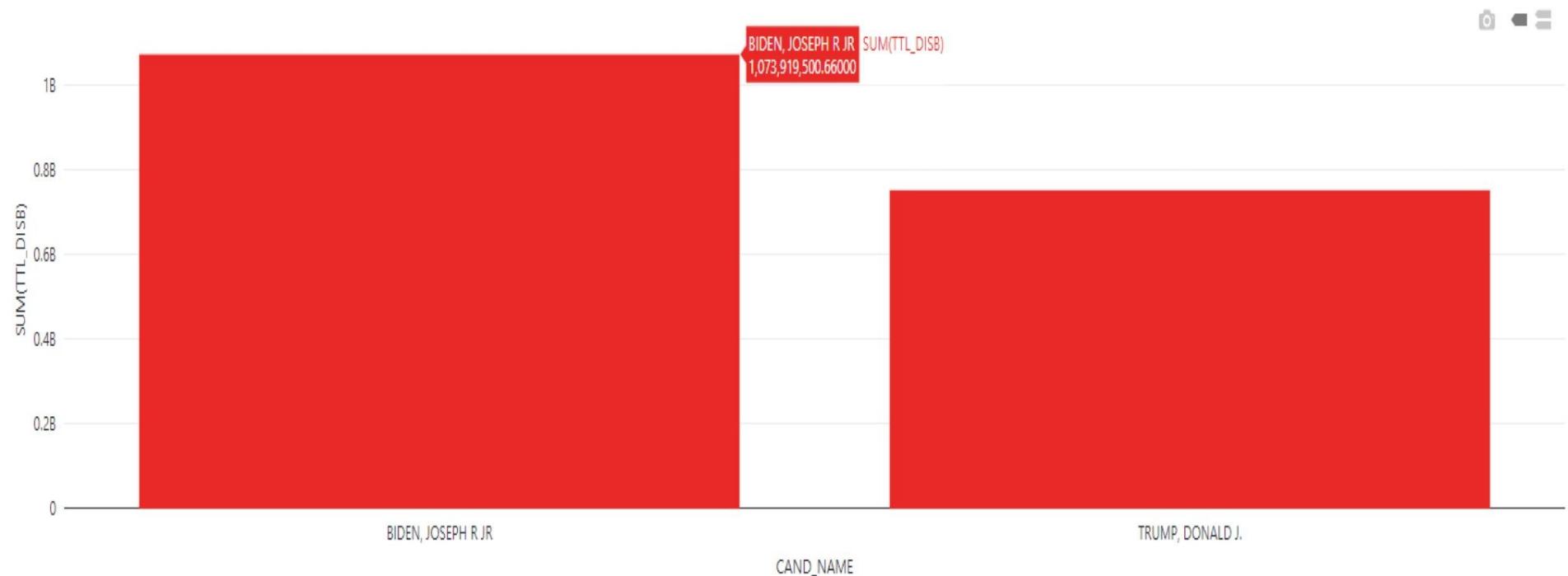
Visualization 2

+

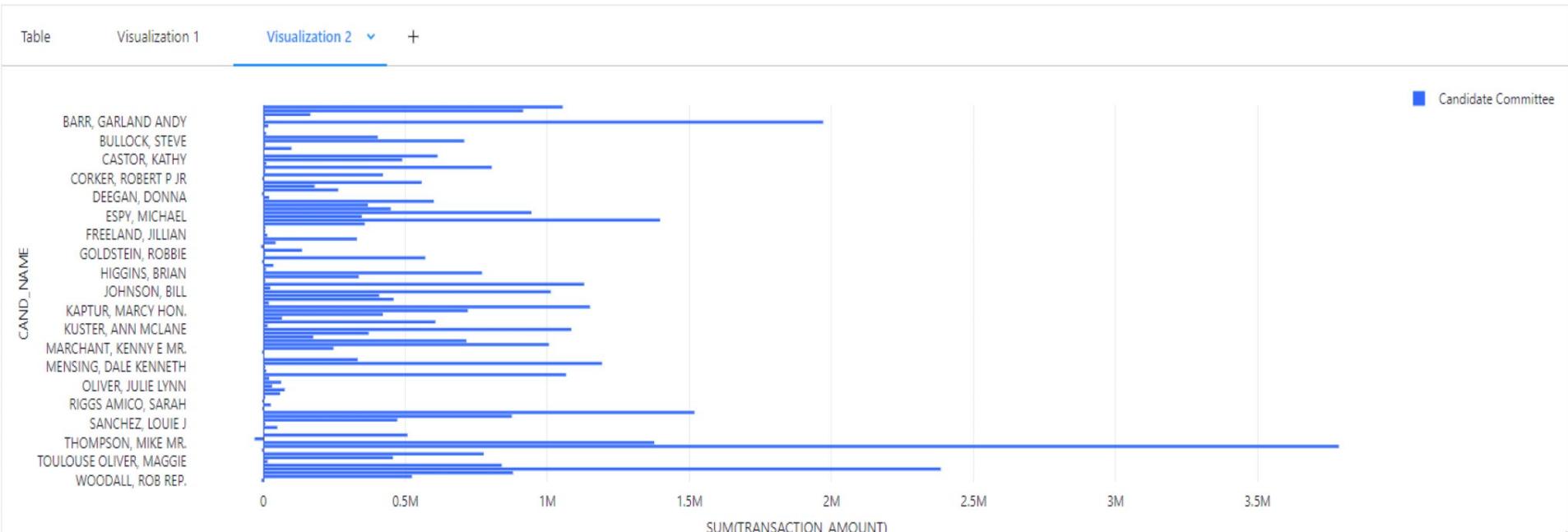


## Total disbursements by Candidates

- Disbursement is a border term that covers both expenditure and other kinds of payments(those not made to influence federal election)



# Transaction amount for each candidate contesting

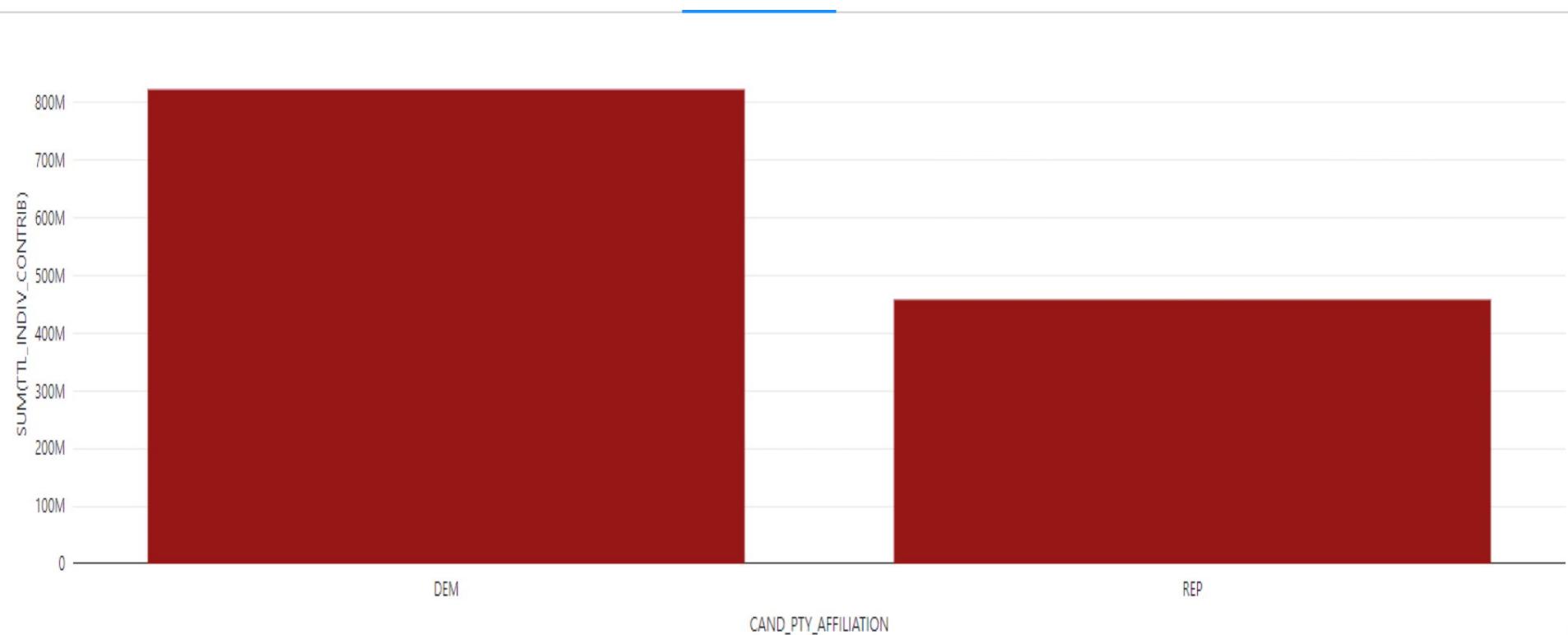


Edit Visualization

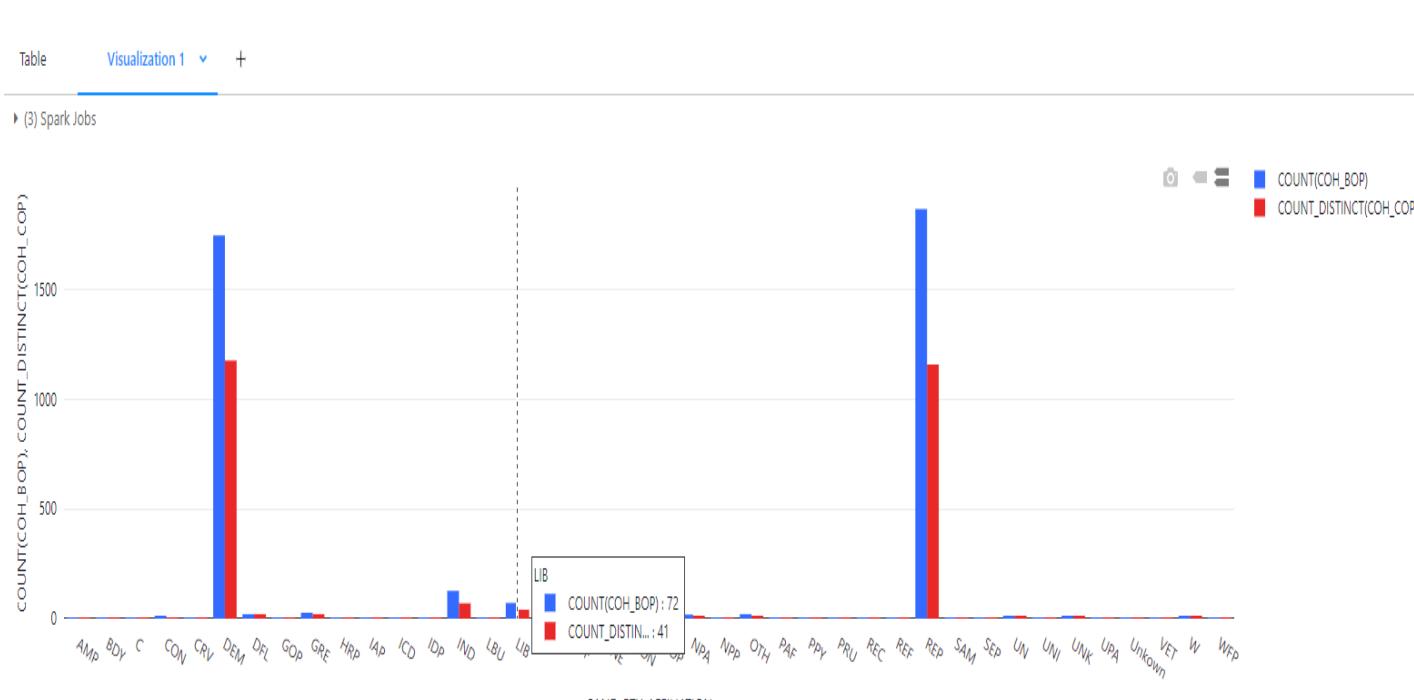
Showing all 100 rows.

Refreshed 14 hours ago

## Total Individual Contributions made by candidates

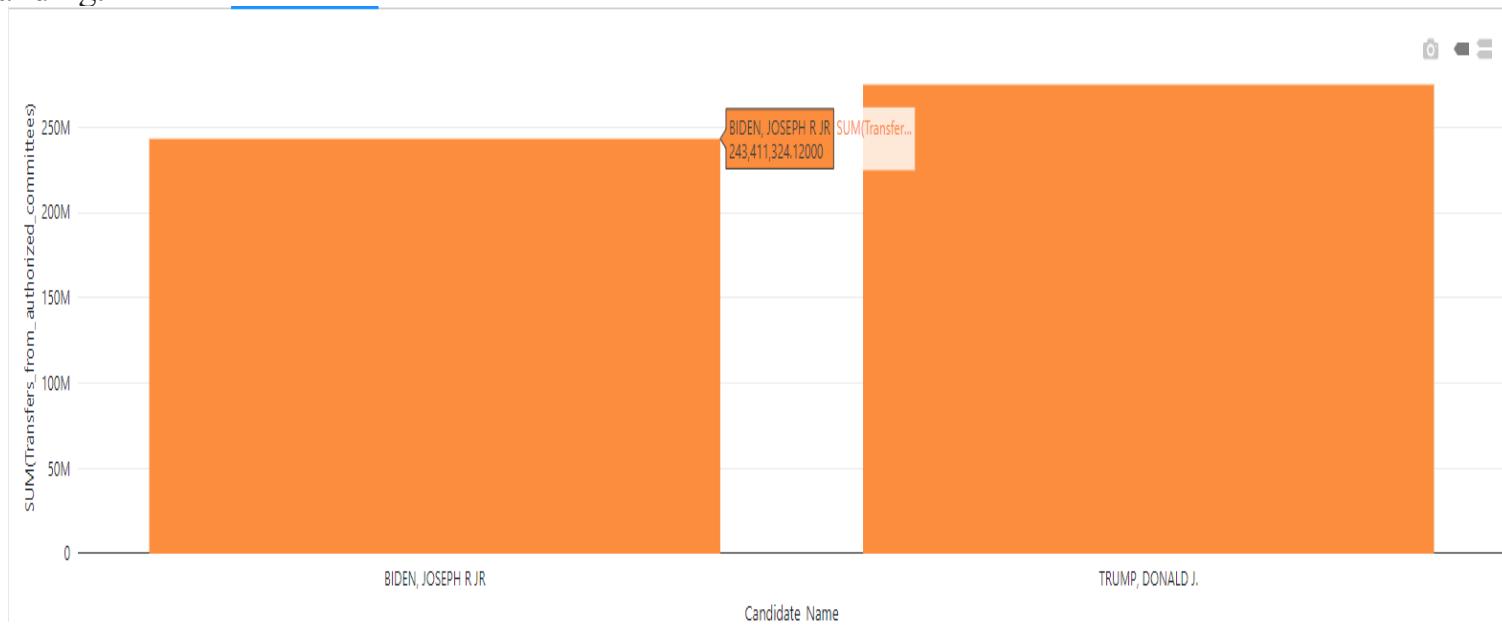


# Visualization of COH\_BOP and COH\_COP for each party



## Visualization for transfers from authorized committees to the candidates:

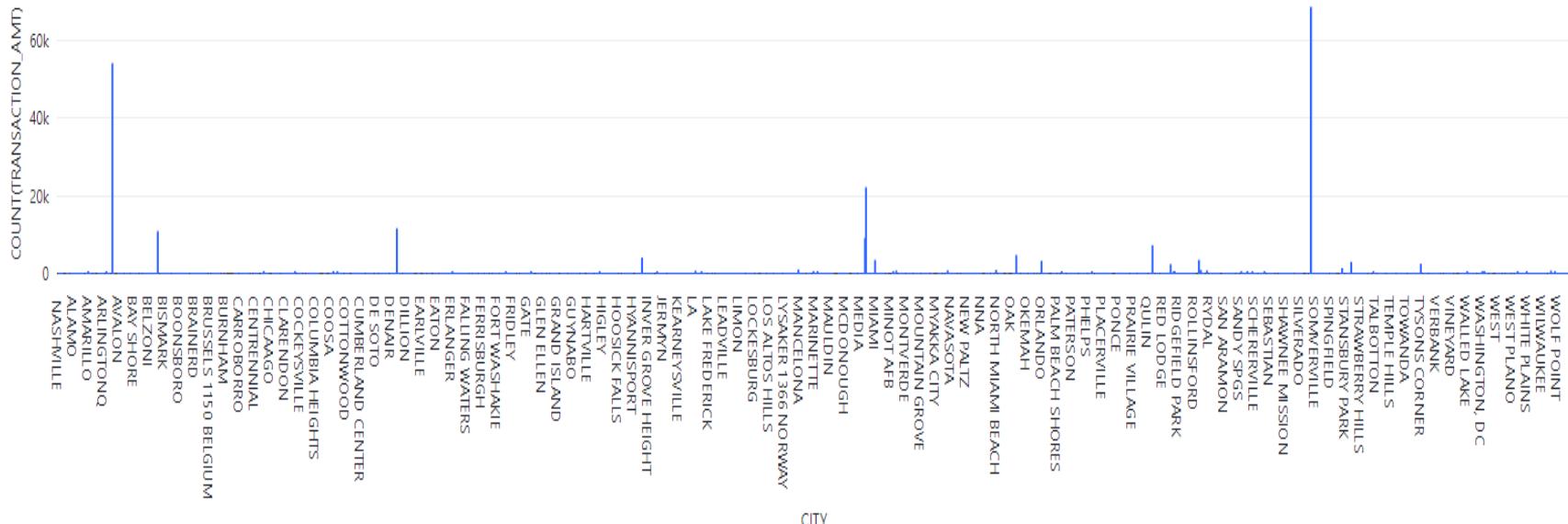
- In general, funds may be transferred between authorized committees of the same candidate (for example, from a previous campaign to a current campaign committee) without limit as long as the committee making the transfer has no net debts outstanding.



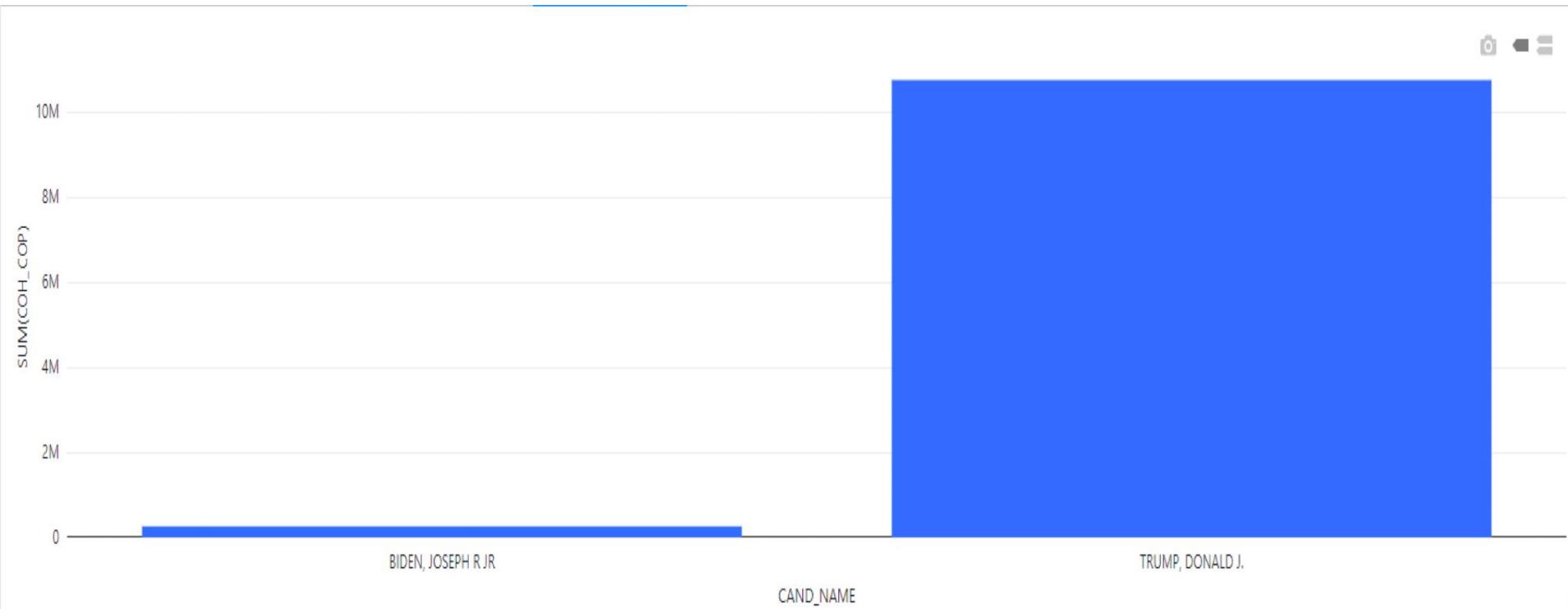


# Visualization of Transaction amount in each City

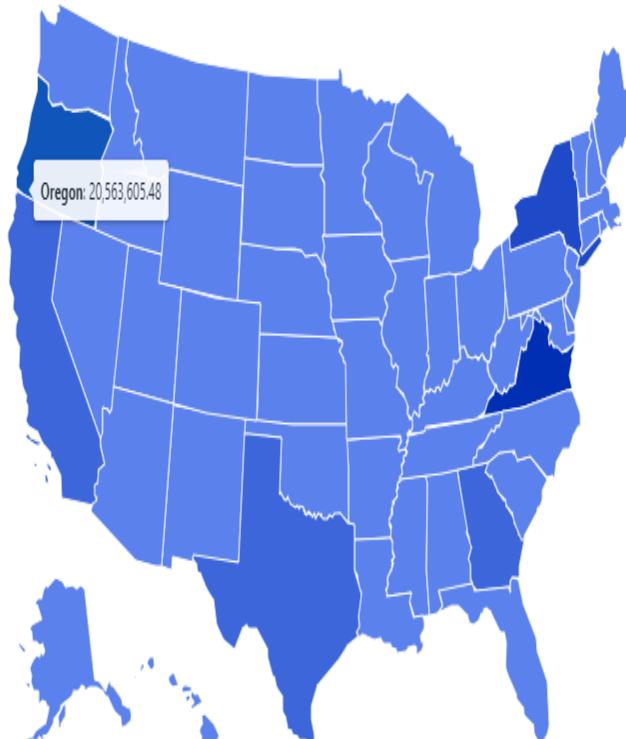
▶ (2) Spark Jobs



## Visualization for COH for the candidates



## Visualization of Transaction amount by state for campaigning





Thank  
you