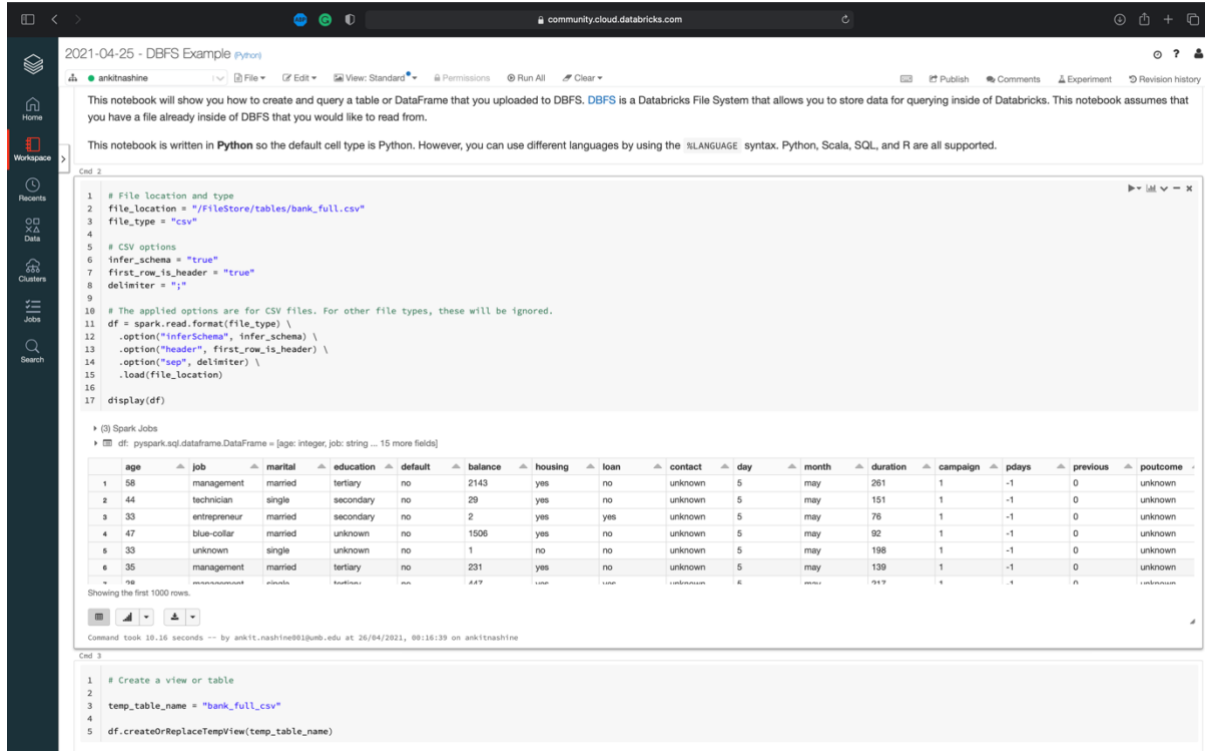# Programming Assignment – 1
# ANKIT NASHINE
# Data – Bank-full.csv

**Importing Bank-full.csv and converting it into table.**

Loading bank-full.csv into databricks and writing codes to convert it to table and viewing it.

# HYPOTHESIS

## 1.1 Total Count of bank Clients

```sql
1   %sql
2
3   select count(*) from `bank_full_csv`
4
```

▸ (2) Spark Jobs

| | count(1) ▲ | |
|---|---|---|
| 1 | 45211 | |

Showing all 1 rows.

Command took 2.30 seconds -- by ankit.nashine001@umb.edu at 26/04/2021, 00:39:04 on ankitnashine

There are total 45211 clients. These clients will be used for further analysis.

## 1.2 Subscribers and non-subscribers (clients) in a term deposit.

```sql
1   %sql
2
3   select Y,count(*) from `bank_full_csv` group by Y
4
```

▸ (2) Spark Jobs



Command took 2.15 seconds -- by ankit.nashine001@umb.edu at 26/04/2021, 01:11:21 on ankitnashine

Out of total clients, only 12% have subscribed for a term deposit. So, to understand subscribers, further analysis will be done on 12% population (i.e. subscribers).

## 2. Subscriber's age group analysis.

```
1  %sql
2
3  select age,count(*) as subscribed_count from `bank_full_csv` where Y='yes' group by age order by subscribed_count desc
4
```

▶ (2) Spark Jobs



⊞  〰 ▾  Plot Options...  ⬇

Command took 2.22 seconds -- by ankit.nashine001@umb.edu at 26/04/2021, 01:56:10 on ankitnashine

It can be inferred from the chart above that most of the subscribers fall in the age range of 25 – 45.

## 3. Contact month of subscribers.

```
1  %sql
2
3  select month,count(*) from `bank_full_csv` where Y='yes' group by month order by count(*) desc
4
```

▶ (2) Spark Jobs



⊞  ▦ ▾  Plot Options...  ⬇

Command took 1.79 seconds -- by ankit.nashine001@umb.edu at 26/04/2021, 02:11:06 on ankitnashine

It can be inferred from the above chart that the best month to contact clients for making them subscribe is May. Should avoid contacting in Jan, Dec to be more feasible.

## 4. Best contact mode clients.

```sql
%sql

select contact,count(contact) from `bank_full_csv` where Y='yes' group by contact order by count(contact) desc
```

▸ (2) Spark Jobs

contact

- cellular
- unknown
- telephone

7%
10%
83%

⊞ ◐ ▾ Plot Options... ⬇

Command took 1.39 seconds -- by ankit.nashine001@umb.edu at 26/04/2021, 02:17:48 on ankitnashine

It can be inferred that clients with cellular mode of contact are more likely to subscribe.

## 5. Marital Status

```sql
%sql


Select marital, count(marital) from `bank_full_csv` where Y='yes' group by marital

```

▸ (2) Spark Jobs



⊞ ▊ ▾ Plot Options... ⬇

Command took 1.60 seconds -- by ankit.nashine001@umb.edu at 26/04/2021, 02:25:23 on ankitnashine

Married and single people are more likely to subscribe.

## 6. Credit History Analysis

```
1  %sql
2
3  select default,count(*) from `bank_full_csv` where Y='yes' group by default order by count(*) desc
4
```

▸ (2) Spark Jobs

| | default ▲ | count(1) ▲ |
|---|---|---|
| 1 | no | 5237 |
| 2 | yes | 52 |

Showing all 2 rows.

⊞  📊 ▾  📥

Command took 1.49 seconds -- by ankit.nashine001@umb.edu at 26/04/2021, 03:04:38 on ankitnashine

Clients with credit history not default are potential subscribers.

## 7. Subscriber's Job Analysis

```
1  %sql
2
3  select job,count(*) from bank_full_csv where Y='yes' group by job order by count(*) desc
4
```

▸ (2) Spark Jobs



⊞  📊 ▾  Plot Options...  📥

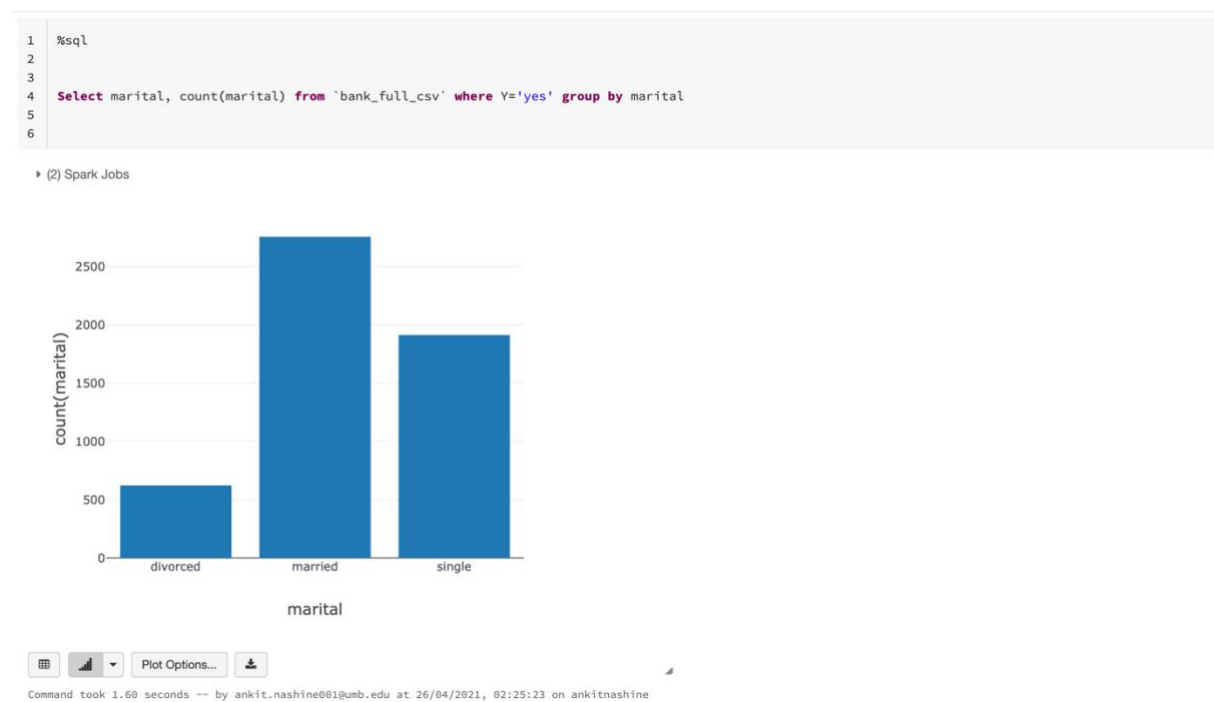Command took 0.79 seconds -- by ankit.nashine001@umb.edu at 26/04/2021, 03:05:44 on ankitnashine

Management job people are most likely to subscribe.
Best is to contact people falls under first half of jobs and ignore second half such as housemaid, unemployed, etc.

## 8. Analyzing subscribers based on Loan Status

```sql
%sql

select loan,count(*) from `bank_full_csv` where Y='yes'  group by loan order by count(*) desc
```

▸ (2) Spark Jobs

| | loan | count(1) |
|---|---|---|
| 1 | no | 4805 |
| 2 | yes | 484 |

Showing all 2 rows.

Command took 0.73 seconds -- by ankit.nashine001@umb.edu at 26/04/2021, 03:10:10 on ankitnashine

Based on the above results, clients with no loans are potential subscribers.

## 9. Final model to predict Potential Subscribers

```sql
%sql

select score, count(*) from
(
select Individual, default_flag+marital_flag+job_flag+loan_flag+contact_flag+housing_flag as score from
(
select
concat(CAST(age AS STRING), CAST(balance AS STRING),marital) as Individual,
case when default='no' then 1 else 0 end default_flag,
case when marital in ('married', 'single') then 1 else 0 end marital_flag,
case when job in ('management', 'technician', 'blue-collar', 'admin', 'retired', 'services') then 1 else 0 end job_flag,
case when loan='no' then 1 else 0 end loan_flag,
case when housing='no' then 1 else 0 end housing_flag,
case when contact='cellular' then 1 else 0 end contact_flag

from `bank_full_csv` where y='no' and age between 25 and 45
) bankindflag )
bankrollup group by score order by score desc

```

▸ (2) Spark Jobs



Command took 0.82 seconds -- by ankit.nashine001@umb.edu at 26/04/2021, 06:09:17 on ankitnashine

Based on previous analysis, predicting clients (non-subscribers) which are most likely to subscribe and contacting then to offer term deposit.
Should contact clients with score 6 and 5.