

FACEBOOK DATA ANALYSIS

The background is a dark blue gradient with abstract data visualizations. On the left, a white line graph with three yellow circular markers is visible. In the center, a large, semi-transparent L-shaped graphic is positioned. To the right of the L-shape, there is a faint bar chart with blue bars. The number '289.33' is displayed in white text near the top of the L-shape and also at the bottom of the page.

289.33

289.33

1. Sanity Check:

Command:

```
export SPARK_MAJOR_VERSION=2
```

```
spark-submit SanityCheck.py
```

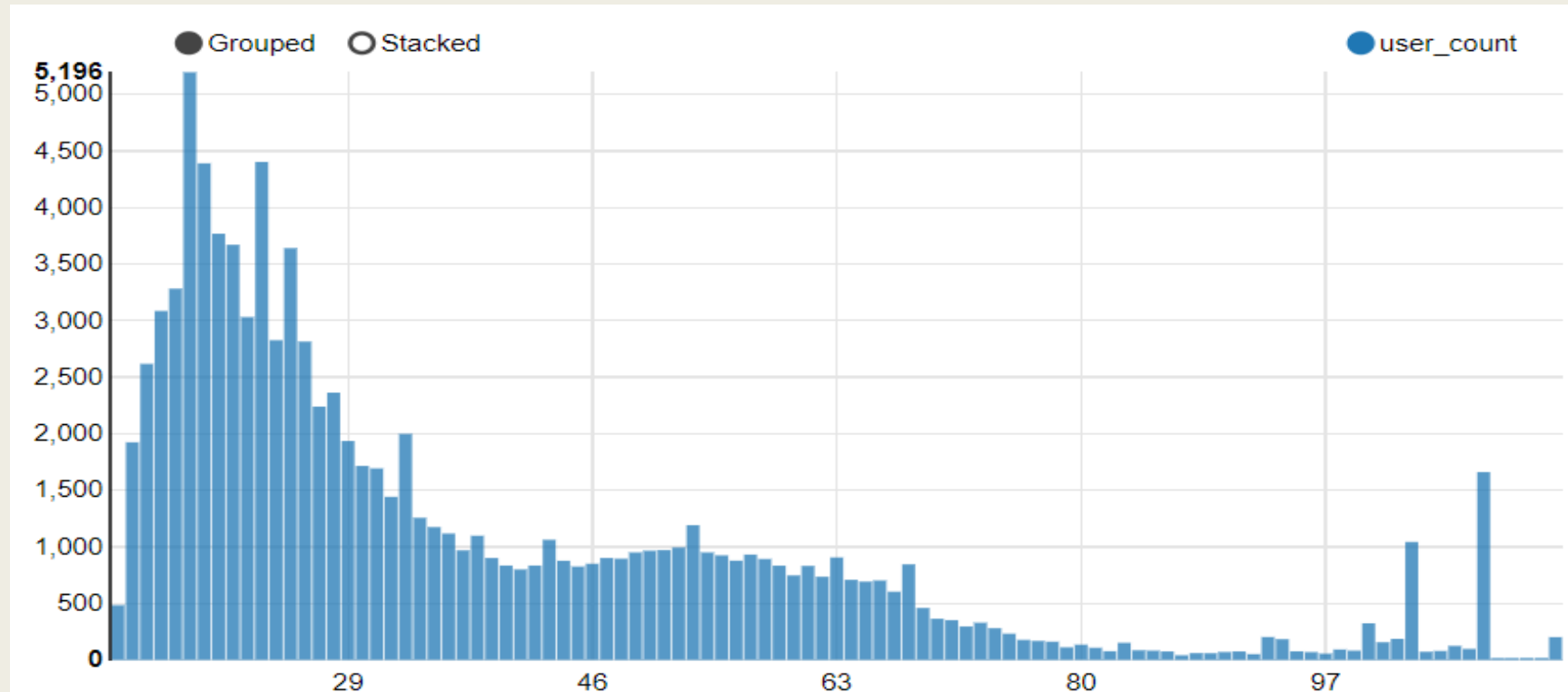
Output:

```
userid : 0
age : 0
dob_day : 0
dob_year : 0
dob_month : 0
gender : 175
tenure : 0
friend_count : 0
friendships_initiated : 0
likes : 0
likes_received : 0
mobile_likes : 0
mobile_likes_received : 0
www_likes : 0
www_likes_received : 0
```

2. Facebook popularity based on ages:

Command: `python map_reduce1.py -r hadoop --hadoop-streaming-jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar hdfs:///tmp/facebook_data/pseudo_facebook.csv`

Age wise distribution of users:



Output: (Age,Count)

```
"109" "000009"
"110" "000015"
"112" "000018"
"111" "000018"
"87" "000042"
"92" "000052"
"97" "000056"
"89" "000060"
"88" "000061"
"96" "000070"
"90" "000071"
"104" "000073"
"86" "000076"
"91" "000076"
"95" "000077"
"82" "000078"
"105" "000080"
"99" "000083"
"85" "000083"
"84" "000086"
"98" "000093"
"107" "000098"
"81" "000108"
"79" "000108"
"106" "000125"
"80" "000136"
"83" "000152"
"101" "000157"
"78" "000162"
"77" "000169"
"76" "000178"
"94" "000184"
"102" "000187"
"42" "000835"
"68" "000846"
"42" "000835"
"68" "000846"
"46" "000851"
"44" "000877"
"56" "000878"
"58" "000893"
"48" "000896"
"47" "000902"
"39" "000902"
"63" "000907"
"55" "000925"
"57" "000932"
"54" "000951"
"49" "000951"
"50" "000966"
"37" "000969"
"51" "000971"
"52" "000995"
"103" "010444"
"43" "010663"
"38" "010999"
"36" "011118"
"35" "011175"
"53" "011192"
"34" "011257"
"32" "011443"
"108" "011661"
"31" "011694"
"30" "011716"
"14" "011925"
"29" "011936"
"33" "011999"
"27" "022240"
"28" "022364"
"15" "022618"
"26" "022815"
"24" "022827"
"22" "030332"
"16" "030886"
"17" "032283"
"25" "033641"
"21" "033671"
"20" "033769"
"19" "043391"
"23" "044404"
"18" "051196"
```

3. Likes Given:

CMD: apache-drill-1.12.0/bin/drillbit.sh start -Ddrill.exec.http.port=8765

Query 1: SELECT gender,avg(likes) AS AVG_Likes_Given

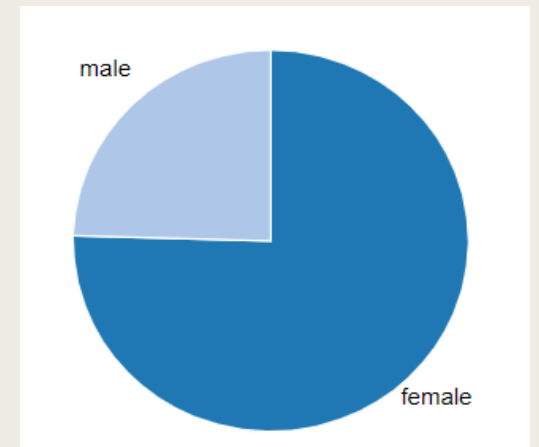
FROM hive.facebook_db.facebook

GROUP BY gender

ORDER BY AVG_Likes_Given DESC

Output: gender vs likes given :

gender	AVG_Likes_Given
female	260.0513240920157
NA	138.50857142857143
male	84.6778946290163



Query 2: SELECT userid, gender, likes AS Total_Likes_Given

FROM hive.facebook_db.facebook

ORDER BY Total_likes_Given DESC LIMIT 10

Output : Top 10 users with most likes received

userid	gender	Total_Likes_Given
1684195	male	25111
1656477	male	21652
1489463	female	16732
1429178	female	16583
1267229	female	14799
1783264	male	14355
1002588	female	14050
1412849	female	14039
1878566	female	13692
2104503	female	13622

4. Likes Received:

CMD: apache-drill-1.12.0/bin/drillbit.sh start -Ddrill.exec.http.port=8765

Query 1: SELECT gender,avg(likes_received) AS AVG_Likes_Received

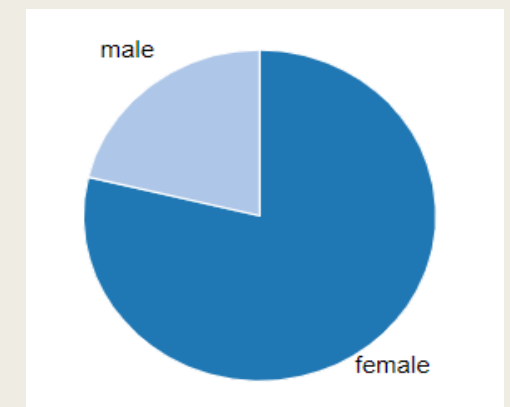
FROM hive.facebook_db.facebook

GROUP BY gender

ORDER BY AVG_Likes_Received DESC

Output: gender vs total likes received :

gender	AVG_Likes_Received
female	251.4354349878273
NA	157.38285714285715
male	67.91154778570697



Query 2: SELECT userid, gender, likes_received AS Total_Likes_Received
FROM hive.facebook_db.facebook
ORDER BY likes_received DESC
LIMIT 10

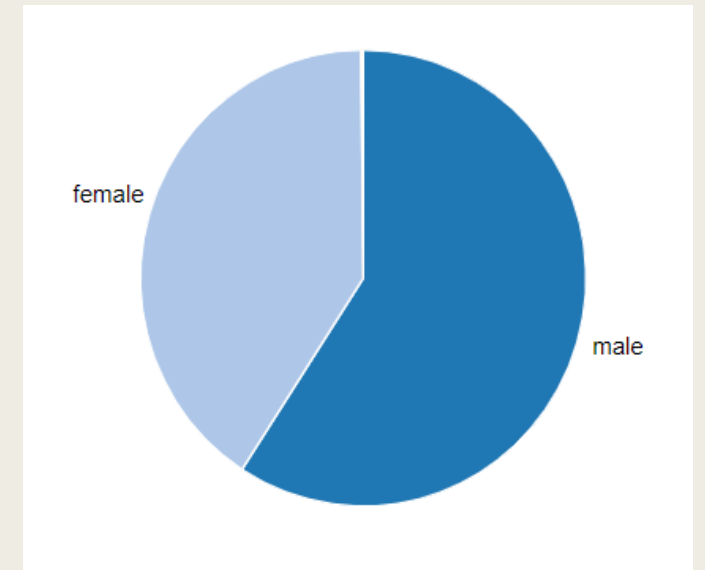
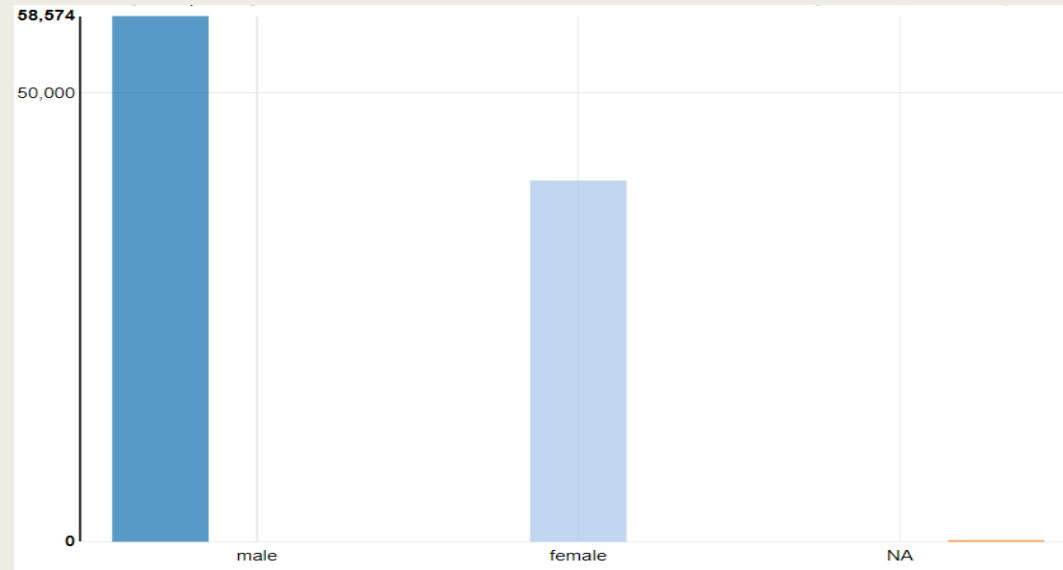
Output : Top 10 users with most likes received

userid	gender	Total_Likes_Received
1674584	female	261197
1441676	female	178166
1715925	female	152014
2063006	female	106025
1053087	male	82623
1432020	male	53534
2042824	male	52964
1559908	female	45633
1781243	female	42449
1015907	male	39536

5. Gender Count:

Output:

```
+-----+-----+
|gender|count|
+-----+-----+
| male |58574|
|female|40254|
| NA   | 175 |
+-----+-----+
```



6. Likes Split Up:

Query 1:

```
SELECT gender,avg(mobile_likes) AS mobile_likes_given, avg(mobile_likes_received) AS  
mobile_likes_received, avg(www_likes) AS www_likes_given, avg(www_likes_received) AS  
www_likes_received
```

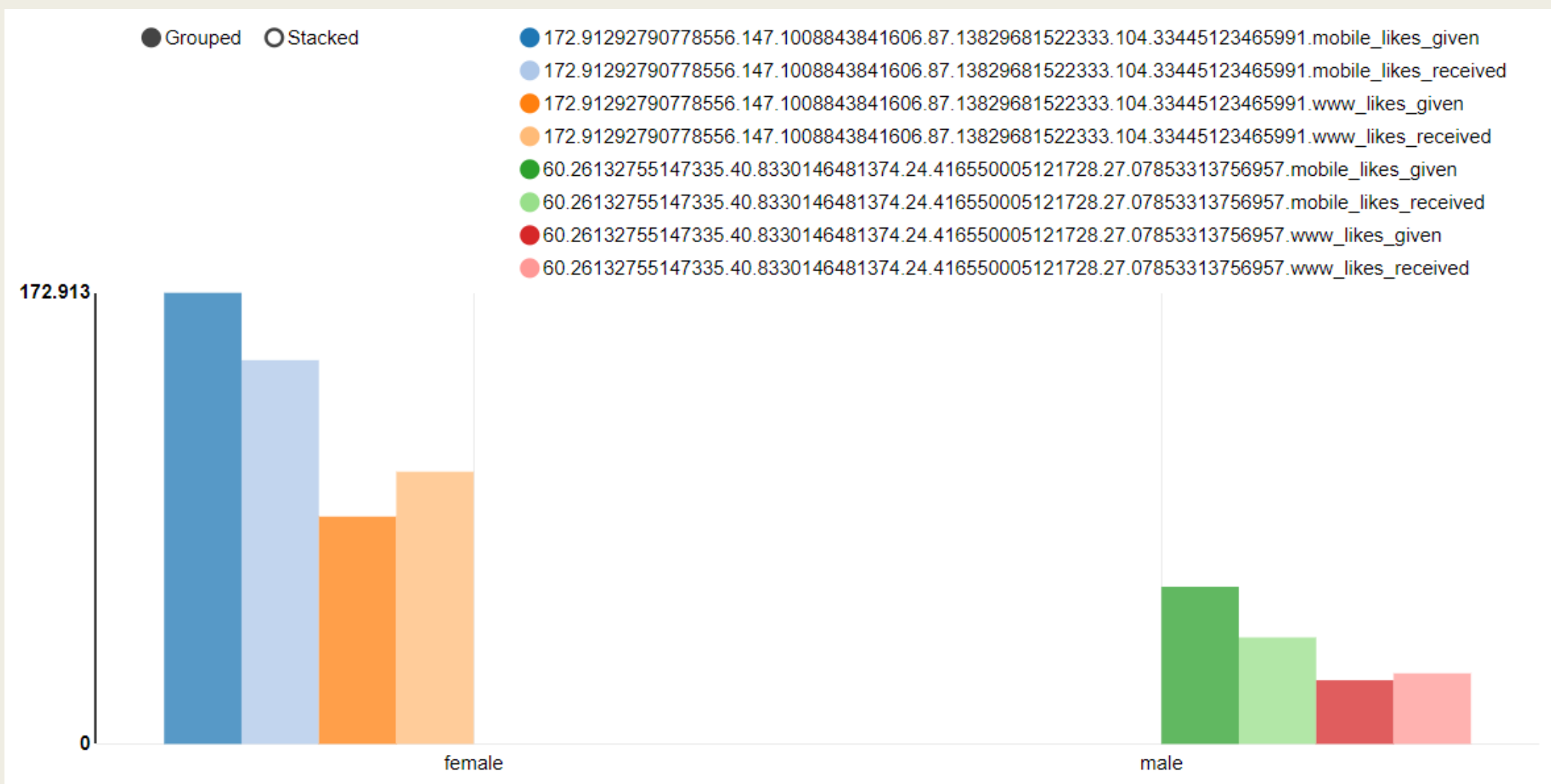
```
FROM fb
```

```
WHERE gender <> "NA"
```

```
GROUP BY gender
```

Output:

gender	mobile_likes_given	mobile_likes_received	www_likes_given	www_likes_received
female	172.91293	147.10088	87.1383	104.33445
male	60.26133	40.83301	24.41655	27.07853



Query2:

%sql

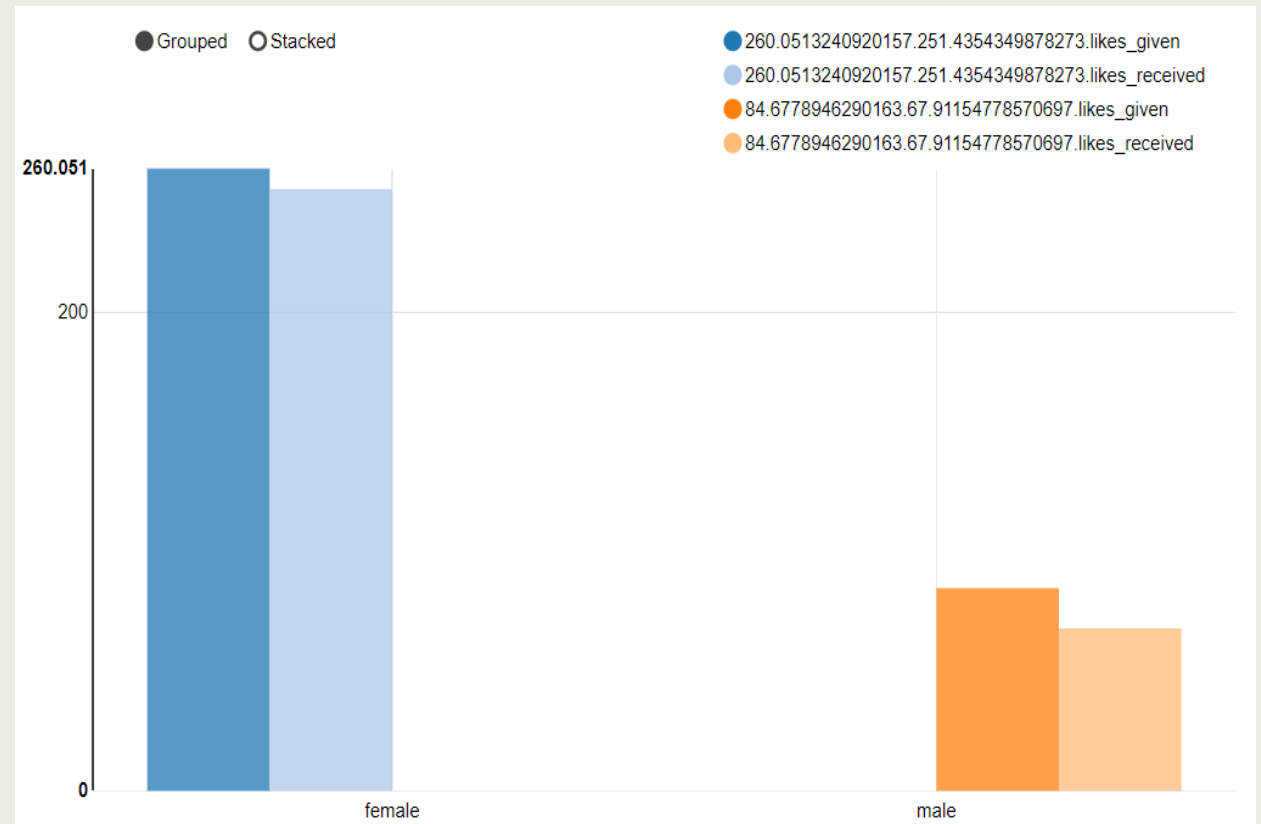
```
SELECT gender,avg(likes) AS likes_given ,avg(likes_received) AS likes_received
```

```
FROM fb
```

```
WHERE gender <> "NA"
```

```
GROUP BY gender
```

Output(Likes vs Likes Recived by gender):



7. Friends Counts & Friendships initiated:

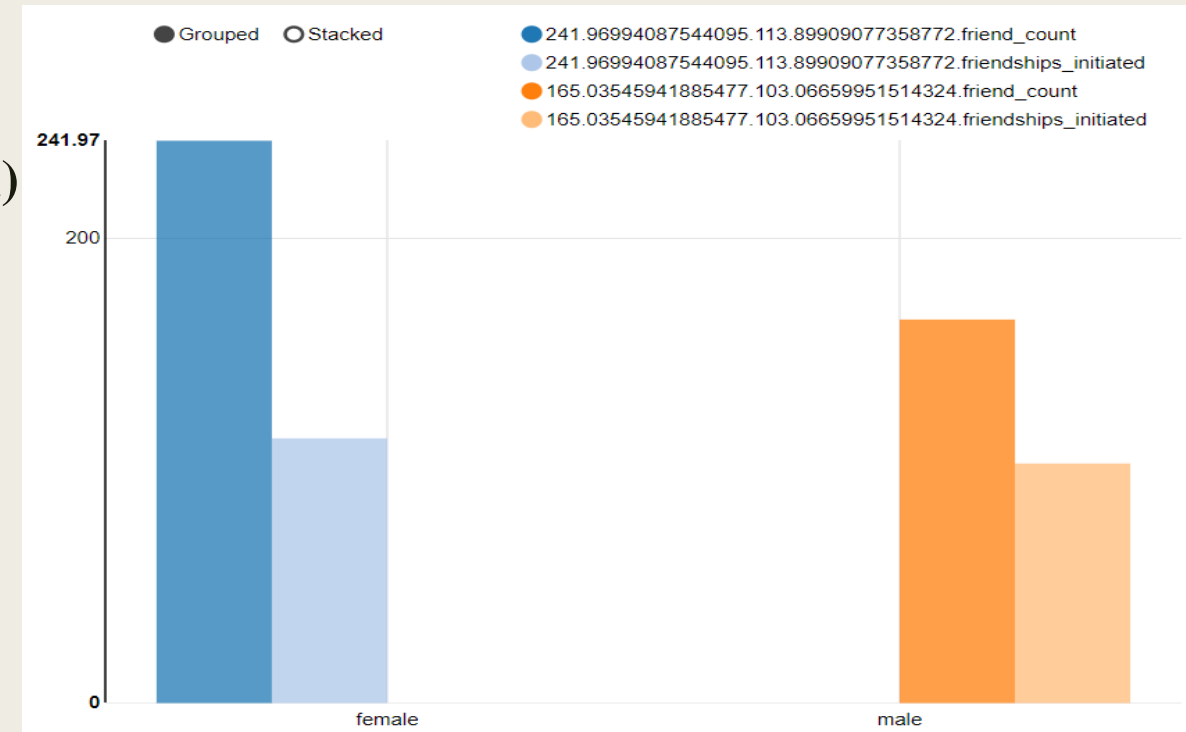
Query :

```
SELECT gender,avg(friend_count) AS friend_count ,avg(friendships_initiated) AS friendships_initiated  
FROM fb
```

```
WHERE gender <> "NA"
```

```
GROUP BY gender
```

Output : (Friends Count vs Friendships Initiated)



8. Users w.r.t birth year:

Query: SELECT dob_year,count(userid) AS users_count
FROM fb
GROUP BY dob_year

Output:

