

Std_across user: 22.07944597647594
Mean_across user: 7.37114764172684

5. Provide the list of users with a mean tagging frequency within the two standard deviation from the mean frequency of all users.

```
In [167]: # defining udf which puts 1 in is_two_std_from_mean column if users mean frequency is under
2*StandardDeviation else 0
@udf("int")
def checkStdFromMean(freq,std,mean):
    if mean - 2 * std < freq < mean + 2 * std:
        return 1
    else:
        return 0

session_list_df=session_list_df.withColumn("is_two_std_from_mean",checkStdFromMean(session_list_df
session_list_df.select(['UserID','mean_Frequency_each_user','std_across_user','mean_across_user'],

-----+
|UserID|mean_Frequency_each_user|std_across_user|mean_across_user|is_two_std_from_mean|
-----+
|15|1.0|22.07944597647594|7.37114764172684|1|
|20|12.0|22.07944597647594|7.37114764172684|1|
|21|2.0|22.07944597647594|7.37114764172684|1|
|25|2.0|22.07944597647594|7.37114764172684|1|
|31|5.0|22.07944597647594|7.37114764172684|1|
|32|1.0|22.07944597647594|7.37114764172684|1|
|39|5.0|22.07944597647594|7.37114764172684|1|
|48|2.0|22.07944597647594|7.37114764172684|1|
|49|15.0|22.07944597647594|7.37114764172684|1|
|75|1.0|22.07944597647594|7.37114764172684|1|
|78|1.0|22.07944597647594|7.37114764172684|1|
|109|2.7777777|22.07944597647594|7.37114764172684|1|
|127|26.0|22.07944597647594|7.37114764172684|1|
|133|5.0|22.07944597647594|7.37114764172684|1|
|146|4.948949|22.07944597647594|7.37114764172684|1|
|147|2.0|22.07944597647594|7.37114764172684|1|
|170|1.0|22.07944597647594|7.37114764172684|1|
|175|1.0|22.07944597647594|7.37114764172684|1|
|181|4.0|22.07944597647594|7.37114764172684|1|
|190|6.5|22.07944597647594|7.37114764172684|1|
-----+
only showing top 20 rows
```

```
In [168]: # showing users within 2 std from mean frequency
session_list_df.filter(session_list_df.is_two_std_from_mean==1).select(['UserID','mean_Frequency_e

-----+
|UserID|mean_Frequency_each_user|std_across_user|mean_across_user|is_two_std_from_mean|
-----+
|15|1.0|22.07944597647594|7.37114764172684|1|
|20|12.0|22.07944597647594|7.37114764172684|1|
|21|2.0|22.07944597647594|7.37114764172684|1|
|25|2.0|22.07944597647594|7.37114764172684|1|
|31|5.0|22.07944597647594|7.37114764172684|1|
|32|1.0|22.07944597647594|7.37114764172684|1|
|39|5.0|22.07944597647594|7.37114764172684|1|
|48|2.0|22.07944597647594|7.37114764172684|1|
|49|15.0|22.07944597647594|7.37114764172684|1|
|75|1.0|22.07944597647594|7.37114764172684|1|
|78|1.0|22.07944597647594|7.37114764172684|1|
|109|2.7777777|22.07944597647594|7.37114764172684|1|
|127|26.0|22.07944597647594|7.37114764172684|1|
|133|5.0|22.07944597647594|7.37114764172684|1|
|146|4.948949|22.07944597647594|7.37114764172684|1|
|147|2.0|22.07944597647594|7.37114764172684|1|
|170|1.0|22.07944597647594|7.37114764172684|1|
|175|1.0|22.07944597647594|7.37114764172684|1|
|181|4.0|22.07944597647594|7.37114764172684|1|
|190|6.5|22.07944597647594|7.37114764172684|1|
-----+
only showing top 20 rows
```

```
In [169]: user_list_with_2_std = sorted((int(row.UserID) for row in
session_list_df.filter(session_list_df.is_two_std_from_mean==1).collect()))
user_list_with_2_std
```


