





## 1. Trip Count per Month

```
>>> from pyspark.sql.functions import month
>>>
>>> df.withColumn("month", month("tpep_pickup_datetime")) \
...     .groupBy("month") \
...     .count() \
...     .orderBy("month") \
...     .show()
+-----+-----+
|month|  count|
+-----+-----+
|    1|     42|
|    2|    362|
|    3| 7866111|
|    4|     78|
|    5|     11|
|    6|      1|
|    7|      5|
|    8|      1|
|   12|      9|
+-----+-----+
```

## 2. Trip Distance Distribution

```
>>> df.select("trip_distance").describe().show()
+-----+-----+
|summary|trip_distance|
+-----+-----+
|  count|      786620|
|   mean| 3.023191633764328|
| stddev|3.9280515814262924|
|    min|           0.0|
|    max|        237.17|
+-----+-----+
```

### 3. Fare vs. Distance Relationship

```
>>> df.select("fare_amount", "trip_distance") \
...   .filter("trip_distance > 0 and fare_amount > 0") \
...   .show(10)
```

fare_amount	trip_distance
13.0	3.7
41.0	14.1
27.0	9.6
5.5	0.8
6.0	1.2
5.5	0.6
17.0	5.65
6.0	1.16
5.0	0.71
10.5	2.63

only showing top 10 rows

### 4. Top Pickup & Drop-off Locations

```
>>> # Top Pickup Locations
>>> df.groupBy("PULocationID") \
...   .count() \
...   .orderBy(col("count").desc()) \
...   .show(10)
[Stage 17:=====] (3 + 1) / 4]25/07/24 10:46:23 WARN Utils: Truncated the string representation
nce it was too large. This behavior can be adjusted by setting 'spark.debug.maxToStringFields' in SparkEnv.conf.
[PULocationID] count
+-----+-----+
161|328433|
237|311896|
162|286072|
236|283119|
230|279261|
186|278438|
48|256545|
234|248504|
170|243858|
142|229982|
+-----+-----+
only showing top 10 rows
```

```

>>>
>>> # Top Drop-off Locations
>>> df.groupBy("DOLocationID") \
...   .count() \
...   .orderBy(col("count").desc()) \
...   .show(10)
+-----+-----+
|DOLocationID| count|
+-----+-----+
|          161|300488|
|          236|300236|
|          237|283692|
|          170|250462|
|          230|238802|
|          162|233992|
|           48|223673|
|          234|211565|
|          142|208331|
|          186|197788|
+-----+-----+
only showing top 10 rows

```

## 5. Average Trip Duration

```

>>> from pyspark.sql.functions import unix_timestamp
>>> df.withColumn("trip_duration_min",
...   (unix_timestamp("tpep_dropoff_datetime") - unix_timestamp("tpep_pickup_datetime")) / 60) \
...   .select("trip_duration_min") \
...   .describe() \
...   .show()
+-----+-----+
|summary| trip_duration_min|
+-----+-----+
| count|          7866620|
| mean| 17.729944578056063|
| stddev| 71.28184046473208|
| min| -4139.016666666666|
| max| 27454.783333333333|
+-----+-----+

```

## 6. Payment Type Distribution

```
>>> df.groupBy("payment_type") \
...     .count() \
...     .orderBy(col("count").desc()) \
...     .show()
+-----+-----+
|payment_type|  count|
+-----+-----+
|          1|5721775|
|          2|2057412|
|          3|   39281|
|          0|   33474|
|          4|   14648|
|          5|        30|
+-----+-----+
>>>
```

## 7. Tip Analysis

```

>>> df.select("tip_amount").describe().show()
+-----+-----+
|summary|    tip_amount|
+-----+-----+
|  count|         7866620|
|   mean| 2.2224419954258554|
| stddev| 50.52138315426142|
|    min|          -89.89|
|    max|        141492.02|
+-----+-----+

>>>
>>> # Tip-to-Fare Ratio (optional)
>>> df.withColumn("tip_ratio", col("tip_amount") / col("fare_amount")) \
...   .select("tip_ratio") \
...   .describe() \
...   .show()
+-----+-----+
|summary|    tip_ratio|
+-----+-----+
|  count|         7863649|
|   mean| 0.19918917663286123|
| stddev| 11.090063586122845|
|    min| -0.4523076923076923|
|    max|        22777.0|
+-----+-----+

```

## 8. Passenger Count Distribution

passenger_count	count
null	33474
0.0	140131
1.0	5496239
2.0	1176562
3.0	329128
4.0	152587
5.0	340192
6.0	198225
7.0	40
8.0	18
9.0	24

[illegible]



## 10. Trip Duration Outliers

```

>>>> # Duration > 6 hours
>>> df.filter("trip_duration_min > 360").show()
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|VendorID|tpep_pickup_datetime|tpep_dropoff_datetime|passenger_count|trip_distance|RatecodeID|store_and_fwd_flag|PULocationID|DOLocationID|payment_type|fare_amount|extra|mta_tax|tip_amount|tolls_amount|improvement_surcharge|total_amount|congestion_surcharge|airport_fee|trip_duration_min|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|7.5|2|2019-02-28 21:34:06|2019-03-01 20:41:16|1.0|1.79|1.0|N|65|61|2|0.5|0.5|0.0|0.0|0.3|8.8|1.0|0.0|null|1387.1666666666667|
|5.5|2|2019-03-01 00:49:35|2019-03-02 00:27:36|2.0|0.35|1.0|N|79|79|2|0.5|0.5|0.0|0.0|0.3|9.3|1.0|2.5|null|1418.8166666666667|
|11.0|2|2019-03-01 00:32:11|2019-03-01 23:45:36|2.0|2.44|1.0|N|79|246|1|0.5|0.5|2.96|0.0|0.3|17.76|1.0|2.5|null|1393.4166666666667|
|30.5|2|2019-02-28 23:58:18|2019-03-01 23:56:12|2.0|6.23|1.0|N|246|255|2|0.5|0.5|0.0|0.0|0.3|34.3|1.0|2.5|null|1437.9|
|19.5|2|2019-02-28 09:38:07|2019-03-01 09:04:26|1.0|4.33|1.0|N|90|256|1|0.5|0.5|4.66|0.0|0.3|27.96|1.0|2.5|null|1486.3166666666666|
|4.5|2|2019-02-28 22:32:43|2019-03-01 21:36:17|2.0|0.78|1.0|N|162|233|1|0.5|0.5|2.49|0.0|0.3|10.79|1.0|2.5|null|1383.5666666666666|
|8.0|2|2019-02-28 16:30:50|2019-03-01 15:38:48|1.0|1.5|1.0|N|244|244|2|0.5|0.5|0.0|0.0|0.3|9.3|1.0|0.0|null|1387.9666666666667|
|30.5|2|2019-03-01 00:01:32|2019-03-01 23:59:18|5.0|5.93|1.0|N|229|65|1|0.5|0.5|0.0|0.0|0.3|34.3|1.0|2.5|null|1437.7666666666667|
|10.0|2|2019-02-28 15:04:48|2019-03-01 14:18:50|1.0|1.71|1.0|N|229|48|1|0.5|0.5|3.45|0.0|0.3|17.25|1.0|2.5|null|1394.0333333333333|
|25.0|2|2019-03-01 00:48:09|2019-03-02 00:04:31|5.0|8.67|1.0|N|138|17|1|0.5|0.5|5.26|0.0|0.3|31.56|1.0|0.0|null|1396.3666666666666|
|5.5|2|2019-03-01 00:28:59|2019-03-02 00:25:10|3.0|0.88|1.0|N|246|68|1|0.5|0.5|0.93|0.0|0.3|10.23|1.0|2.5|null|1436.1833333333334|
|19.0|2|2019-02-28 16:37:05|2019-03-01 16:29:19|2.0|5.32|1.0|N|164|13|1|0.5|0.5|0.0|0.0|0.3|22.8|1.0|2.5|null|1432.2333333333333|
|6.0|2|2019-03-01 00:42:38|2019-03-02 00:40:12|2.0|1.03|1.0|N|166|238|2|0.5|0.5|0.0|0.0|0.3|7.3|1.0|0.0|null|1437.5666666666666|
|5.0|2|2019-03-01 00:56:51|2019-03-02 00:55:01|5.0|0.73|1.0|N|230|186|2|0.5|0.5|0.0|0.0|0.3|8.8|1.0|2.5|null|1438.1666666666667|
|18.0|2|2019-02-28 02:16:08|2019-03-01 01:39:14|1.0|4.63|1.0|N|79|142|1|0.5|0.5|4.36|0.0|0.3|26.16|1.0|2.5|null|1403.1|

```

[illegible][illegible]

```

172.5| 0.0| 0.5| 22.0| 0.0| 0.3| 195.3| 0.0| 0.0| null|27.216666666666665|
| 1| 2019-03-01 07:19:43| 2019-03-01 07:30:38| 1.0| 0.0| 1.0|
10.5| 2.5| 0.5| 2.75| 0.0| 0.3| 16.55| 2.5| null|10.916666666666666|
| 2| 2019-03-01 06:28:52| 2019-03-01 07:28:06| 1.0| 0.0| 5.0|
61.0| 0.0| 0.0| 19.05| 0.0| 0.0| 82.55| 2.5| null|59.233333333333334|
| 1| 2019-03-01 07:08:53| 2019-03-01 08:26:39| 1.0| 0.0| 1.0|
74.4| 0.0| 0.5| 0.0| 0.0| 0.3| 75.2| 0.0| null|77.766666666666667|
| 2| 2019-03-01 08:15:14| 2019-03-01 08:47:36| 1.0| 0.0| 1.0|
2.5| 0.0| 0.5| 0.0| 0.0| 0.3| 5.8| 2.5| null|32.366666666666666|
| 1| 2019-03-01 08:49:40| 2019-03-01 09:59:10| 1.0| 0.0| 1.0|
70.4| 0.0| 0.5| 0.0| 5.76| 0.3| 76.96| 0.0| null|69.5|
| 2| 2019-03-01 07:57:44| 2019-03-01 08:19:31| 2.0| 0.0| 1.0|
-2.5| 0.0| -0.5| 0.0| 0.0| -0.3| -3.3| 0.0| null|21.783333333333335|
+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows

```

## 6. Example – Trip Duration Calculation

```

>>> from pyspark.sql.functions import unix_timestamp, col
>>>
>>> df = df.withColumn("trip_duration_min",
...                   (unix_timestamp("tpep_dropoff_datetime") - unix_timestamp("tpep_pickup_datetime")) / 60)
>>>
>>> df.select("trip_duration_min").describe().show()
+-----+-----+
|summary| trip_duration_min|
+-----+-----+
| count|          7866620|
| mean| 17.729944578056063|
| stddev| 71.28184046473208|
| min| -4139.016666666666|
| max| 27454.783333333333|
+-----+-----+

```

## 7. Convert to Pandas for Plotting (optional)

## 8. Save Cleaned/Analyzed Data to S3 (optional)