

Abon, Benedict Aldous A.

CPE311-CPE22S3

```
import pandas as pd
```

1. Create DataFrame; Examine the first 5 rows

```
data = pd.read_csv('2019_Yellow_Taxi_Trip_Data.csv')
data.head(5)

{"summary": {"name": "data", "rows": 10000, "fields": [{"column": "vendorid", "properties": {"dtype": "number", "std": 0, "min": 1, "max": 2, "num_unique_values": 2, "samples": [1, 2], "semantic_type": "\\", "description": "\n"}, "semantic_type": "object", "dtype": "object", "num_unique_values": 4281, "samples": ["2019-10-23T16:44:33.000", "2019-10-23T16:26:46.000"], "description": "\n"}, {"column": "tpep_pickup_datetime", "properties": {"dtype": "object", "num_unique_values": 4281, "samples": ["2019-10-23T16:44:33.000", "2019-10-23T16:26:46.000"], "semantic_type": "\\", "description": "\n"}, "semantic_type": "date", "dtype": "date", "num_unique_values": 4281, "samples": ["2019-10-23T16:44:33.000", "2019-10-23T16:26:46.000"], "description": "\n"}, {"column": "tpep_dropoff_datetime", "properties": {"dtype": "object", "num_unique_values": 5079, "samples": ["2019-10-23T17:45:43.000", "2019-10-23T17:49:28.000"], "semantic_type": "\\", "description": "\n"}, "semantic_type": "date", "dtype": "date", "num_unique_values": 5079, "samples": ["2019-10-23T17:45:43.000", "2019-10-23T17:49:28.000"], "description": "\n"}, {"column": "passenger_count", "properties": {"dtype": "number", "std": 1, "min": 0, "max": 6, "num_unique_values": 7, "samples": [1, 2], "semantic_type": "\\", "description": "\n"}, "semantic_type": "integer", "dtype": "integer", "min": 0, "max": 6, "num_unique_values": 7, "samples": [1, 2], "description": "\n"}, {"column": "trip_distance", "properties": {"dtype": "number", "std": 4.148063386748824, "min": 0.0, "max": 38.11, "num_unique_values": 1243, "samples": [5.41, 10.29], "semantic_type": "\\", "description": "\n"}, "semantic_type": "float", "dtype": "float", "min": 0.0, "max": 38.11, "num_unique_values": 1243, "samples": [5.41, 10.29], "description": "\n"}], "semantic_type": "\\", "description": "\n"}, {"column": "ratecodeid", "properties": {"dtype": "number", "std": 0, "min": 1, "max": 5, "num_unique_values": 5, "samples": [2, 4], "semantic_type": "\\", "description": "\n"}, "semantic_type": "integer", "dtype": "integer", "min": 1, "max": 5, "num_unique_values": 5, "samples": [2, 4], "description": "\n"}, {"column": "store_and_fwd_flag", "properties": {"dtype": "category", "num_unique_values": 2, "samples": ["Y", "N"], "semantic_type": "\\", "description": "\n"}, "semantic_type": "category", "dtype": "category", "num_unique_values": 2, "samples": ["Y", "N"], "description": "\n"}, {"column": "pulocationid", "properties": {"dtype": "number", "std": 0, "min": 1, "max": 2, "num_unique_values": 2, "samples": [1, 2], "semantic_type": "\\", "description": "\n"}, "semantic_type": "integer", "dtype": "integer", "min": 1, "max": 2, "num_unique_values": 2, "samples": [1, 2], "description": "\n"}], "semantic_type": "\\", "description": "\n"}}
```



```

\"total_amount\", \n      \"properties\": {\n          \"dtype\":\n\"number\", \n          \"std\": 19.209255488783793, \n          \"min\": -\n65.92, \n          \"max\": 671.8, \n          \"num_unique_values\": 1097,\n          \"samples\": [\n              14.16, \n              85.7\n          ], \n          \"semantic_type\": \"\", \n          \"description\": \"\"\n      }, \n      { \n          \"column\": \"congestion_surcharge\", \n          \"properties\": {\n              \"dtype\": \"number\", \n              \"std\": 0.7209463477704644, \n              \"min\": -2.5, \n              \"max\": 2.75,\n              \"num_unique_values\": 4, \n              \"samples\": [\n                  0.0, \n                  2.75\n              ], \n              \"semantic_type\": \"\", \n              \"description\": \"\"\n          }\n      }\n  },\n  \"type\": \"dataframe\", \"variable_name\": \"data\"}

```

- Find the dimensions (number of rows and number of columns) in the data.

```
data.shape
```

```
(10000, 18)
```

- Calculate summary statistics for the fare_amount, tip_amount, tolls_amount, and total_amount columns.

```

# fare amount stats
fa_mean = data['fare_amount'].mean()
fa_median = data['fare_amount'].median()
fa_mode = data['fare_amount'].mode()
fa_std = data['fare_amount'].std()
# print stats
print(f'Fare amount mean: {fa_mean}')
print(f'Fare amount median: {fa_median}')
print(f'Fare amount mode: {fa_mode[0]}')
print(f'Fare amount standard deviation: {fa_std}')

Fare amount mean: 15.106313
Fare amount median: 10.0
Fare amount mode: 52.0
Fare amount standard deviation: 13.954761757679899

# tip_amount stats
ta_mean = data['tip_amount'].mean()
ta_median = data['tip_amount'].median()
ta_mode = data['tip_amount'].mode()
ta_std = data['tip_amount'].std()
# print stats
print(f'Tip amount mean:{ta_mean}')
print(f'Tip amount median: {ta_median}')
print(f'Tip amount mode: {ta_mode[0]}')
print(f'Tip amount standard deviation: {ta_std}')

Tip amount:2.6344939999999997
Tip amount median: 2.0

```

```

Tip amount mode: 0.0
Tip amount standard deviation: 3.409800031906561

# tolls_amount stats
taa_mean = data['tolls_amount'].mean()
taa_median = data['tolls_amount'].median()
taa_mode = data['tolls_amount'].mode()
taa_std = data['tolls_amount'].std()
# print stats
print(f'Tolls amount mean: {taa_mean}')
print(f'Tolls amount median: {taa_median}')
print(f'Tolls amount mode: {taa_mode[0]}')
print(f'Tolls amount standard deviation: {taa_std}')

Tolls amount: 0.623447
Tolls amount median: 0.0
Tolls amount: 0.0
Tolls amount standard deviation: 6.437507127609063

# total_amount stats
taat_mean = data['total_amount'].mean()
taat_median = data['total_amount'].median()
taat_mode = data['total_amount'].mode()
taat_std = data['total_amount'].std()
# print stats
print(f'Total amount mean: {taat_mean}')
print(f'Total amount median: {taat_median}')
print(f'Total amount mode: {taat_mode[0]}')
print(f'Total amount standard deviation: {taat_std}')

Total amount mean: 22.564659
Total amount median: 16.3
Total amount mode: 11.8
Total amount standard deviation: 19.209255488783793

```

1. Isolate the fare_amount, tip_amount, tolls_amount, and total_amount for the longest trip by distance (trip_distance).

```

data.sort_values(by='trip_distance', ascending=False).head(1)

{"repr_error": "0", "type": "dataframe"}

```