



```
In [ ]: import pandas as pd
        from sklearn.model_selection import train_test_split
        from sklearn.linear_model import Lasso
        import matplotlib.pyplot as plt
        import kagglehub
        import os
        from sklearn.metrics import r2_score
```

```
In [ ]: path = kagglehub.dataset_download("chicago/chicago-taxi-rides-2016")
        print(os.listdir(path))
```

```
['column_remapping.json', 'chicago_taxi_trips_2016_01.csv', 'chicago_taxi_trips_2016_11.csv', 'data_dictionary.csv', 'chicago_taxi_trips_2016_05.csv', 'chicago_taxi_trips_2016_12.csv', 'chicago_taxi_trips_2016_04.csv', 'chicago_taxi_trips_2016_07.csv', 'chicago_taxi_trips_2016_10.csv', 'chicago_taxi_trips_2016_06.csv', 'chicago_taxi_trips_2016_08.csv', 'chicago_taxi_trips_2016_03.csv', 'chicago_taxi_trips_2016_02.csv', 'chicago_taxi_trips_2016_09.csv']
```

```
In [ ]: csv_path = os.path.join(path, 'chicago_taxi_trips_2016_10.csv')
        df = pd.read_csv(csv_path)
        df.head()
```

```
Out[ ]:   taxi_id  trip_start_timestamp  trip_end_timestamp  trip_seconds  trip_miles  p
0    2059.0    2016-10-2 16:45:00    2016-10-2 17:15:00        2460.0        17.4
1    6308.0    2016-10-26 11:30:00    2016-10-26 12:00:00        1860.0        14.9
2    2595.0    2016-10-21 12:00:00    2016-10-21 12:15:00         420.0         1.0
3    6764.0    2016-10-21 23:15:00    2016-10-21 23:30:00         240.0         0.8
4    3824.0    2016-10-2 00:00:00    2016-10-2 00:15:00         600.0         1.8
```

```
In [ ]: df.shape
```

```
Out[ ]: (1499771, 20)
```

```
In [ ]: df.isnull().sum()
```

Out[]:

0

taxi_id	441
trip_start_timestamp	0
trip_end_timestamp	77
trip_seconds	88
trip_miles	40
pickup_census_tract	1499771
dropoff_census_tract	537186
pickup_community_area	139784
dropoff_community_area	170899
fare	43
tips	43
tolls	43
extras	43
trip_total	43
payment_type	0
company	634747
pickup_latitude	139743
pickup_longitude	139743
dropoff_latitude	167193
dropoff_longitude	167193

dtype: int64

```
In [ ]: numeric_cols = ['trip_seconds', 'trip_miles', 'fare', 'tips', 'tolls', 'extras',  
                        'pickup_latitude', 'pickup_longitude', 'dropoff_latitude', 'dr  
for col in numeric_cols:  
    df[col].fillna(df[col].mean(), inplace=True)
```

/tmp/ipython-input-6-2840052652.py:4: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.

The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
df[col].fillna(df[col].mean(), inplace=True)
```

```
In [ ]: df.isnull().sum()
```

```
Out[ ]:
```

	0
taxi_id	441
trip_start_timestamp	0
trip_end_timestamp	77
trip_seconds	0
trip_miles	0
pickup_census_tract	1499771
dropoff_census_tract	537186
pickup_community_area	139784
dropoff_community_area	170899
fare	0
tips	0
tolls	0
extras	0
trip_total	0
payment_type	0
company	634747
pickup_latitude	0
pickup_longitude	0
dropoff_latitude	0
dropoff_longitude	0

dtype: int64

```
In [ ]: df = df.drop(["trip_start_timestamp"], axis=1)
```

```
In [ ]: df.isnull().sum()
```

```
Out[ ]:
```

	0
trip_seconds	0
trip_miles	0
fare	0
tips	0
tolls	0
extras	0
trip_total	0
pickup_latitude	0
pickup_longitude	0
dropoff_latitude	0
dropoff_longitude	0

dtype: int64

```
In [ ]: x = df.drop("fare", axis=1).values
y = df["fare"].values
```

```
In [ ]: x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=
```

```
In [ ]: model=Lasso()
model.fit(x_train,y_train)
prediction=model.predict(x_test)
print(prediction)
print(y_test)

[ 8.24846699  9.26288814 32.37617242 ...  8.75372855 11.46376676
12.6947243 ]
[ 8.25  9.5  32.25 ...  8.75 11.75 13. ]
```

```
In [ ]: r2 = r2_score(y_test, prediction)
print("R2 Score:", r2)
```

R² Score: 0.9994963028095779

```
In [ ]: plt.scatter(y_test, prediction, color='blue')
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'r--')
plt.xlabel('Actual Fare')
plt.ylabel('Predicted Fare')
plt.title('Actual vs. Predicted Fare')
plt.show()
```

Actual vs. Predicted Fare

