```python
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import datetime as dt


df = pd.read_csv("uber.csv")
df.head()
```

```
---------------------------------------------------------------------
-----
FileNotFoundError                          Traceback (most recent call
last)
Cell In[3], line 1
----> 1 df = pd.read_csv("uber.csv")
      2 df.head()

File ~\anaconda3\Lib\site-packages\pandas\util\_decorators.py:211, in
deprecate_kwarg.<locals>._deprecate_kwarg.<locals>.wrapper(*args,
**kwargs)
    209     else:
    210         kwargs[new_arg_name] = new_arg_value
--> 211 return func(*args, **kwargs)

File ~\anaconda3\Lib\site-packages\pandas\util\_decorators.py:331, in
deprecate_nonkeyword_arguments.<locals>.decorate.<locals>.wrapper(*arg
s, **kwargs)
    325 if len(args) > num_allow_args:
    326     warnings.warn(
    327
msg.format(arguments=_format_argument_list(allow_args)),
    328         FutureWarning,
    329         stacklevel=find_stack_level(),
    330     )
--> 331 return func(*args, **kwargs)

File ~\anaconda3\Lib\site-packages\pandas\io\parsers\readers.py:950,
in read_csv(filepath_or_buffer, sep, delimiter, header, names,
index_col, usecols, squeeze, prefix, mangle_dupe_cols, dtype, engine,
converters, true_values, false_values, skipinitialspace, skiprows,
skipfooter, nrows, na_values, keep_default_na, na_filter, verbose,
skip_blank_lines, parse_dates, infer_datetime_format, keep_date_col,
date_parser, dayfirst, cache_dates, iterator, chunksize, compression,
thousands, decimal, lineterminator, quotechar, quoting, doublequote,
escapechar, comment, encoding, encoding_errors, dialect,
error_bad_lines, warn_bad_lines, on_bad_lines, delim_whitespace,
low_memory, memory_map, float_precision, storage_options)
    935 kwds_defaults = _refine_defaults_read(
    936     dialect,
    937     delimiter,
```

```
    (...)
    946        defaults={"delimiter": ","},
    947 )
    948 kwds.update(kwds_defaults)
--> 950 return _read(filepath_or_buffer, kwds)

File ~\anaconda3\Lib\site-packages\pandas\io\parsers\readers.py:605,
in _read(filepath_or_buffer, kwds)
    602 _validate_names(kwds.get("names", None))
    604 # Create the parser.
--> 605 parser = TextFileReader(filepath_or_buffer, **kwds)
    607 if chunksize or iterator:
    608        return parser

File ~\anaconda3\Lib\site-packages\pandas\io\parsers\readers.py:1442,
in TextFileReader.__init__(self, f, engine, **kwds)
    1439        self.options["has_index_names"] = kwds["has_index_names"]
    1441 self.handles: IOHandles | None = None
-> 1442 self._engine = self._make_engine(f, self.engine)

File ~\anaconda3\Lib\site-packages\pandas\io\parsers\readers.py:1735,
in TextFileReader._make_engine(self, f, engine)
    1733        if "b" not in mode:
    1734            mode += "b"
-> 1735 self.handles = get_handle(
    1736        f,
    1737        mode,
    1738        encoding=self.options.get("encoding", None),
    1739        compression=self.options.get("compression", None),
    1740        memory_map=self.options.get("memory_map", False),
    1741        is_text=is_text,
    1742        errors=self.options.get("encoding_errors", "strict"),
    1743        storage_options=self.options.get("storage_options", None),
    1744 )
    1745 assert self.handles is not None
    1746 f = self.handles.handle

File ~\anaconda3\Lib\site-packages\pandas\io\common.py:856, in
get_handle(path_or_buf, mode, encoding, compression, memory_map,
is_text, errors, storage_options)
    851 elif isinstance(handle, str):
    852        # Check whether the filename is to be opened in binary
mode.
    853        # Binary mode does not support 'encoding' and 'newline'.
    854        if ioargs.encoding and "b" not in ioargs.mode:
    855            # Encoding
--> 856            handle = open(
    857                handle,
    858                ioargs.mode,
    859                encoding=ioargs.encoding,
```

```
    860             errors=errors,
    861             newline="",
    862         )
    863     else:
    864         # Binary mode
    865         handle = open(handle, ioargs.mode)

FileNotFoundError: [Errno 2] No such file or directory: 'uber.csv'


df.drop(columns=['Unnamed: 0','key'],inplace=True)

df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200000 entries, 0 to 199999
Data columns (total 7 columns):
 #   Column             Non-Null Count   Dtype
---  ------             --------------   -----
 0   fare_amount        200000 non-null  float64
 1   pickup_datetime    200000 non-null  object
 2   pickup_longitude   200000 non-null  float64
 3   pickup_latitude    200000 non-null  float64
 4   dropoff_longitude  199999 non-null  float64
 5   dropoff_latitude   199999 non-null  float64
 6   passenger_count    200000 non-null  int64
dtypes: float64(5), int64(1), object(1)
memory usage: 10.7+ MB

df.dropna(how='any',inplace=True)


df.isnull().sum()

fare_amount          0
pickup_datetime      0
pickup_longitude     0
pickup_latitude      0
dropoff_longitude    0
dropoff_latitude     0
passenger_count      0
dtype: int64

for col in df.select_dtypes(exclude=['object']):
    plt.figure()
    sns.boxplot(data=df,x=col)
```
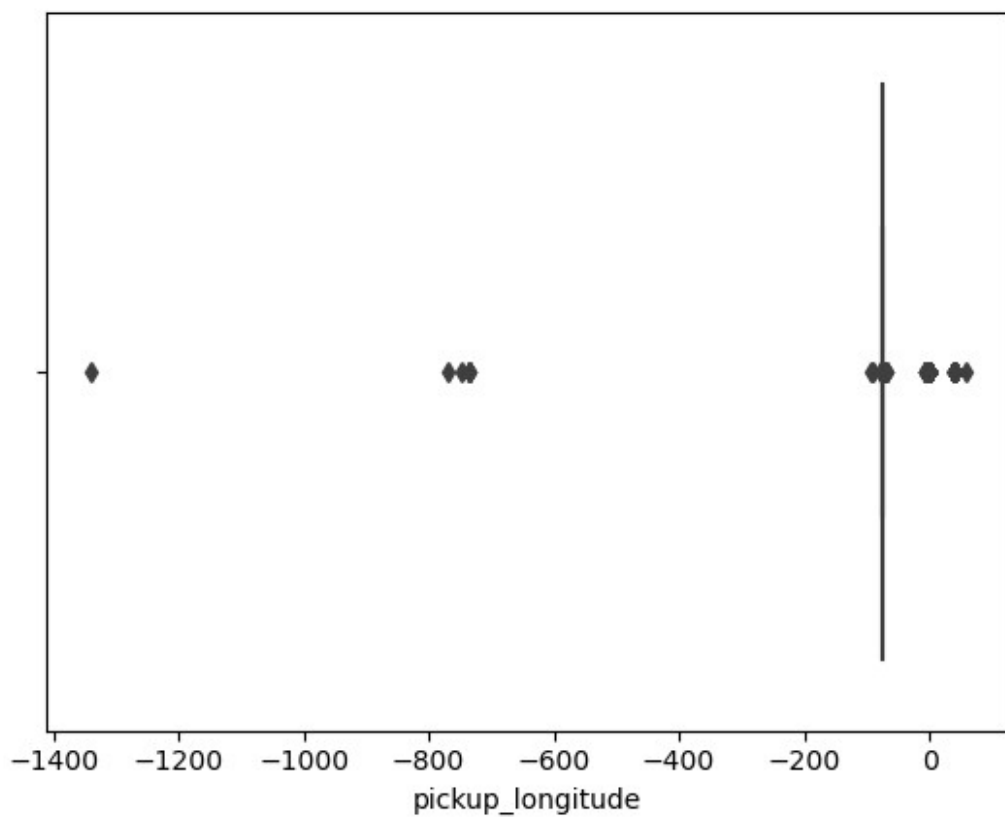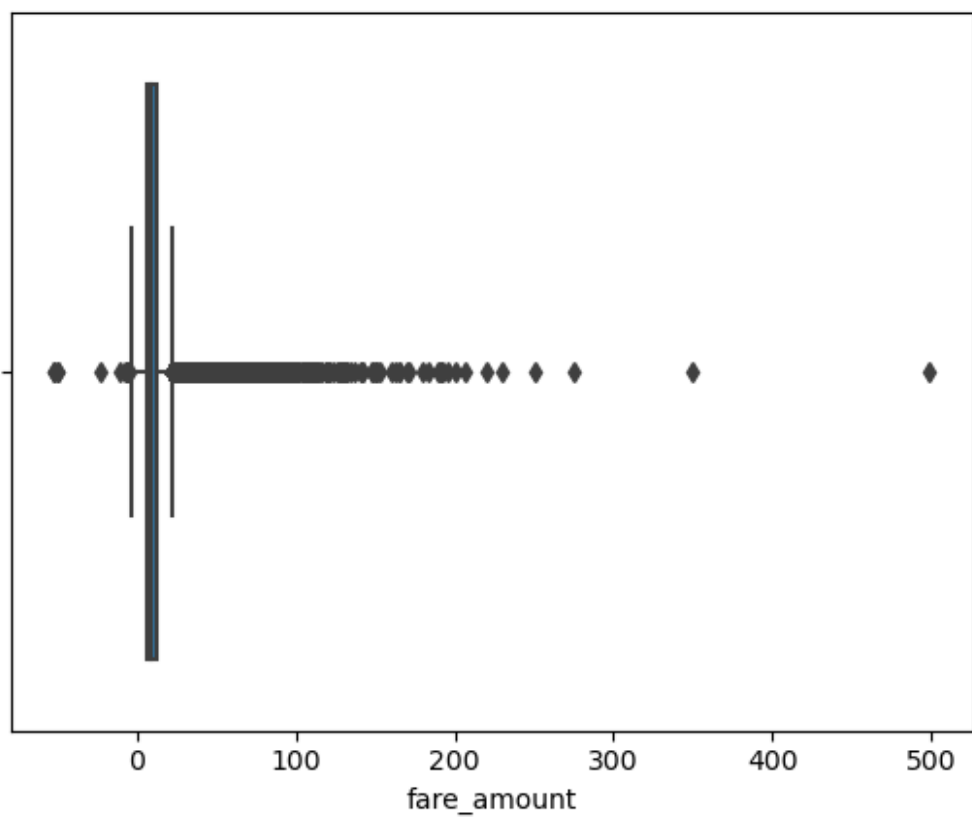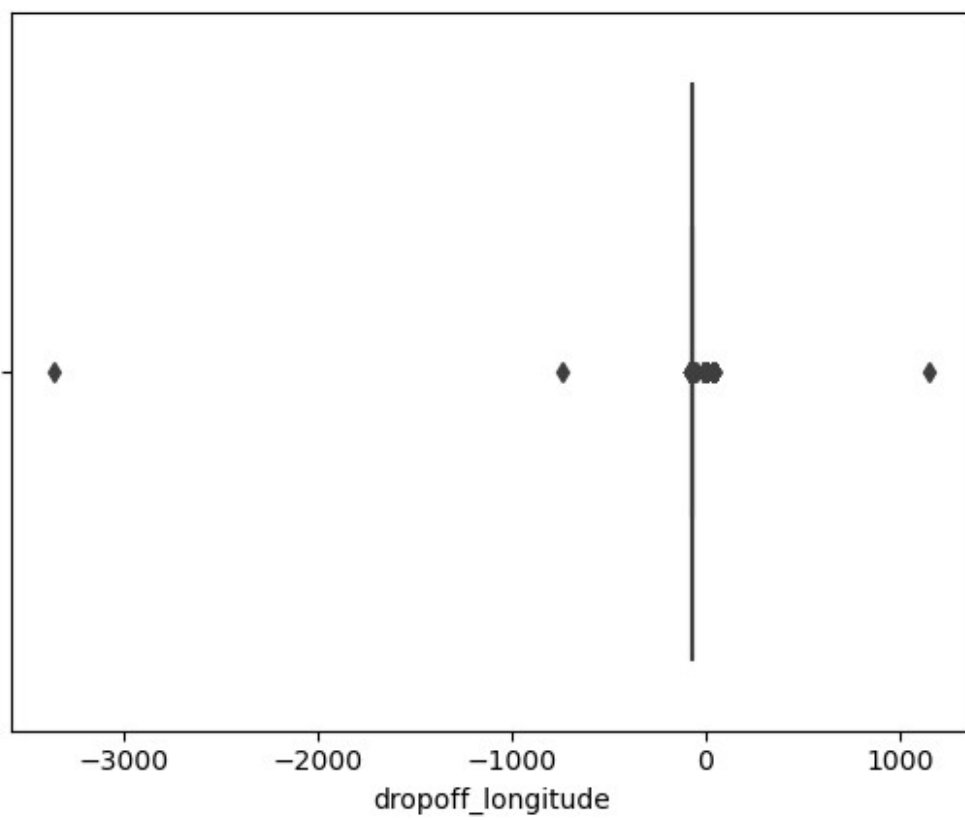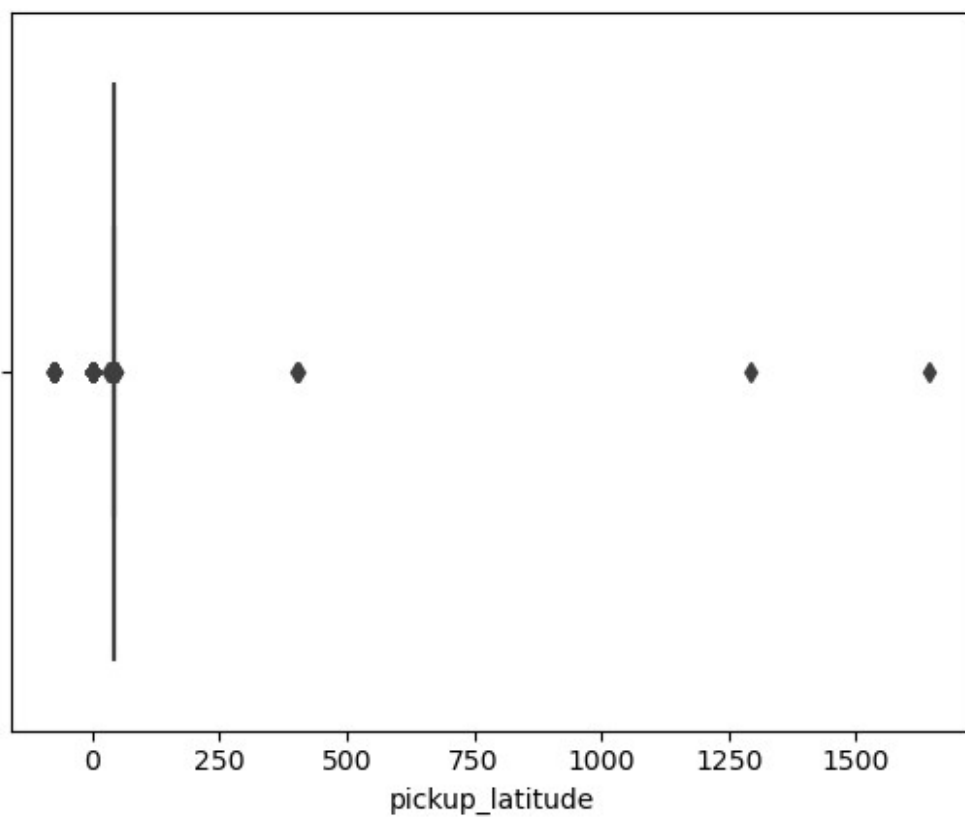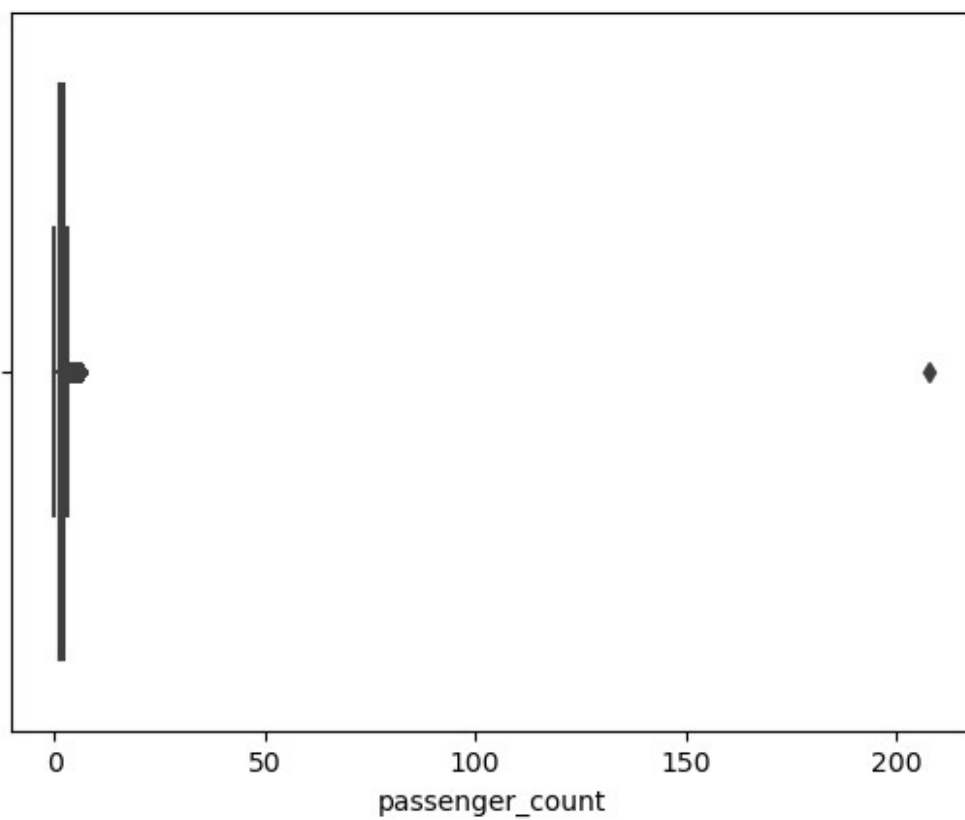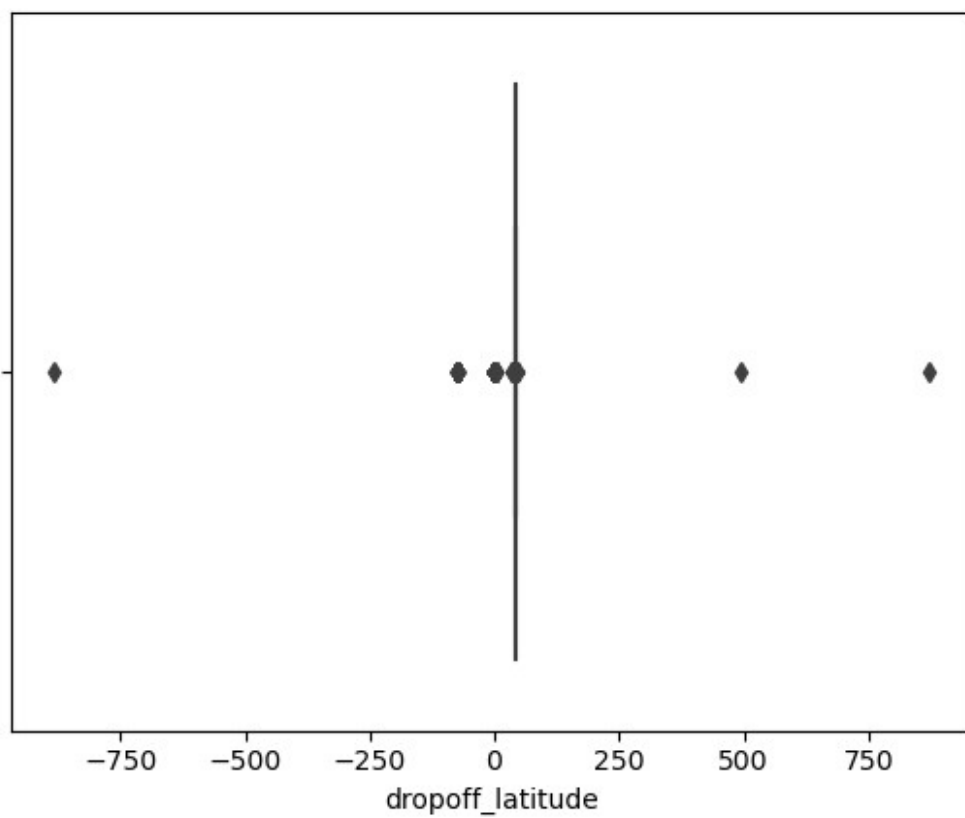
fare_amount

pickup_longitude

pickup_latitude

dropoff_longitude

```python
df = df[
    (df.pickup_latitude > -90) & (df.pickup_latitude < 90) &
    (df.dropoff_latitude > -90) & (df.dropoff_latitude < 90) &
    (df.pickup_longitude > -180) & (df.pickup_longitude < 180) &
    (df.dropoff_longitude > -180) & (df.dropoff_longitude < 180) &
    (df.fare_amount > 0) & (df.passenger_count > 0)  &
(df.passenger_count < 50)
]

from math import cos, asin, sqrt, pi
import numpy as np

def distance(lat_1,lon_1,lat_2,lon_2):
#       lat1 = row.pickup_latitude
#       lon1 = row.pickup_longitude
#       lat2 = row.dropoff_latitude
#       lon2 = row.dropoff_longitude
    lon_1, lon_2, lat_1, lat_2 = map(np.radians, [lon_1, lon_2, lat_1,
lat_2])   #Degrees to Radians


    diff_lon = lon_2 - lon_1
    diff_lat = lat_2 - lat_1


    km = 2 * 6371 * np.arcsin(np.sqrt(np.sin(diff_lat/2.0)**2 +
np.cos(lat_1) * np.cos(lat_2) * np.sin(diff_lon/2.0)**2))

    return km


temp =
distance(df['pickup_latitude'],df['pickup_longitude'],df['dropoff_lati
tude'],df['dropoff_longitude'])
temp.head()

0    1.683323
1    2.457590
2    5.036377
3    1.661683
4    4.475450
dtype: float64

df_new = df.copy()
df_new['Distance'] = temp
df = df_new
df.head()

   fare_amount          pickup_datetime  pickup_longitude
pickup_latitude  \
0          7.5  2015-05-07 19:52:06 UTC        -73.999817
```

```
        40.738354
1             7.7   2009-07-17 20:04:56 UTC          -73.994355
        40.728225
2            12.9   2009-08-24 21:45:00 UTC          -74.005043
        40.740770
3             5.3   2009-06-26 08:22:21 UTC          -73.976124
        40.790844
4            16.0   2014-08-28 17:47:00 UTC          -73.925023
        40.744085

   dropoff_longitude   dropoff_latitude   passenger_count   Distance
0         -73.999512          40.723217                 1   1.683323
1         -73.994710          40.750325                 1   2.457590
2         -73.962565          40.772647                 1   5.036377
3         -73.965316          40.803349                 3   1.661683
4         -73.973082          40.761247                 5   4.475450
```
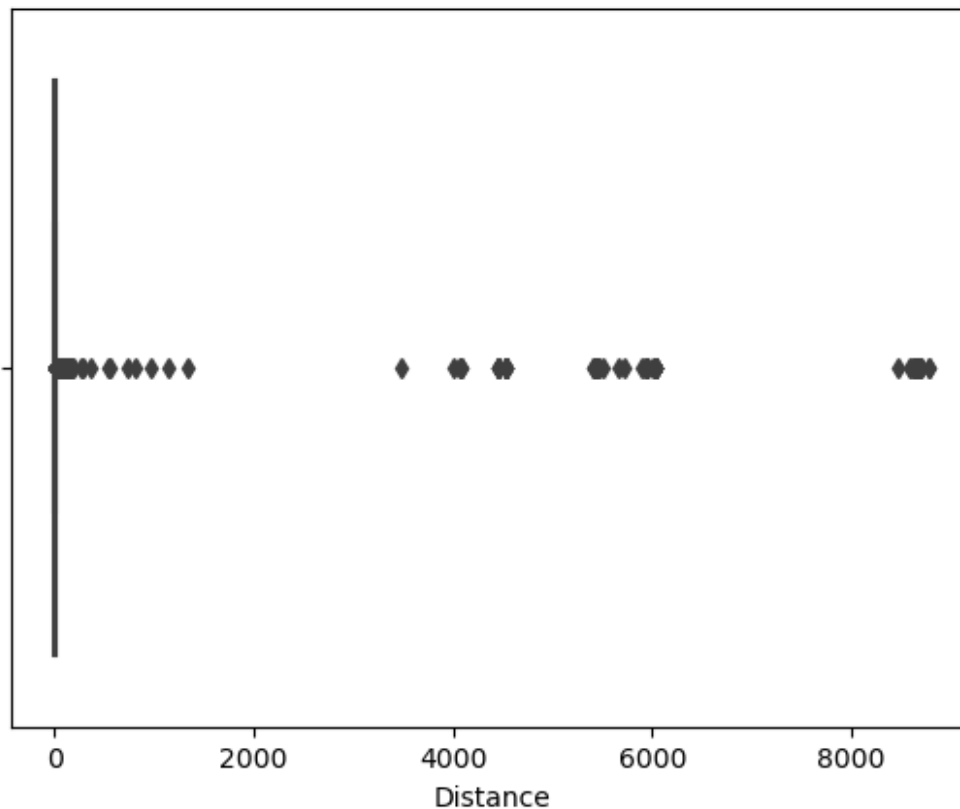
```python
sns.boxplot(data=df,x='Distance')
```

```
<Axes: xlabel='Distance'>
```



```python
df = df[(df['Distance'] < 200) & (df['Distance'] > 0)]
```

```python
df['pickup_datetime'] = pd.to_datetime(df['pickup_datetime'])
```

C:\Users\prajw\AppData\Local\Temp\ipykernel_26664\200932308.py:1:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation:
https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#
returning-a-view-versus-a-copy
  df['pickup_datetime'] = pd.to_datetime(df['pickup_datetime'])

```python
df['week_day'] = df['pickup_datetime'].dt.day_name()
df['Year'] = df['pickup_datetime'].dt.year
df['Month'] = df['pickup_datetime'].dt.month
df['Hour'] = df['pickup_datetime'].dt.hour
```

C:\Users\prajw\AppData\Local\Temp\ipykernel_26664\2592915223.py:1:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation:
https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#
returning-a-view-versus-a-copy
  df['week_day'] = df['pickup_datetime'].dt.day_name()
C:\Users\prajw\AppData\Local\Temp\ipykernel_26664\2592915223.py:2:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation:
https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#
returning-a-view-versus-a-copy
  df['Year'] = df['pickup_datetime'].dt.year
C:\Users\prajw\AppData\Local\Temp\ipykernel_26664\2592915223.py:3:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation:
https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#
returning-a-view-versus-a-copy
  df['Month'] = df['pickup_datetime'].dt.month
C:\Users\prajw\AppData\Local\Temp\ipykernel_26664\2592915223.py:4:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation:
```

```
https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#
returning-a-view-versus-a-copy
  df['Hour'] = df['pickup_datetime'].dt.hour

df.drop(columns=['pickup_datetime','pickup_latitude','pickup_longitude
','dropoff_latitude','dropoff_longitude'],inplace=True)

C:\Users\prajw\AppData\Local\Temp\ipykernel_26664\3782303944.py:1:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation:
https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#
returning-a-view-versus-a-copy

df.drop(columns=['pickup_datetime','pickup_latitude','pickup_longitude
','dropoff_latitude','dropoff_longitude'],inplace=True)

df.head()

    fare_amount  passenger_count  Distance  week_day   Year  Month  Hour
0          7.5                1  1.683323  Thursday   2015      5    19
1          7.7                1  2.457590    Friday   2009      7    20
2         12.9                1  5.036377    Monday   2009      8    21
3          5.3                3  1.661683    Friday   2009      6     8
4         16.0                5  4.475450  Thursday   2014      8    17

temp = df.copy()

def convert_week_day(day):
    if day in ['Monday','Tuesday','Wednesday','Thursday']:
        return 0 # Weekday
    return 1 # Weekend

def convert_hour(hour):
    if 5 <= hour <= 12:
        return 1
    elif 12 < hour <= 17:
        return 2
    elif 17 < hour < 24:
        return 3
    return 0

df['week_day'] = temp['week_day'].apply(convert_week_day)
df['Hour'] = temp['Hour'].apply(convert_hour)
df.head()


C:\Users\prajw\AppData\Local\Temp\ipykernel_26664\3260682206.py:17:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
```

```
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation:
https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#
returning-a-view-versus-a-copy
  df['week_day'] = temp['week_day'].apply(convert_week_day)
C:\Users\prajw\AppData\Local\Temp\ipykernel_26664\3260682206.py:18:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation:
https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#
returning-a-view-versus-a-copy
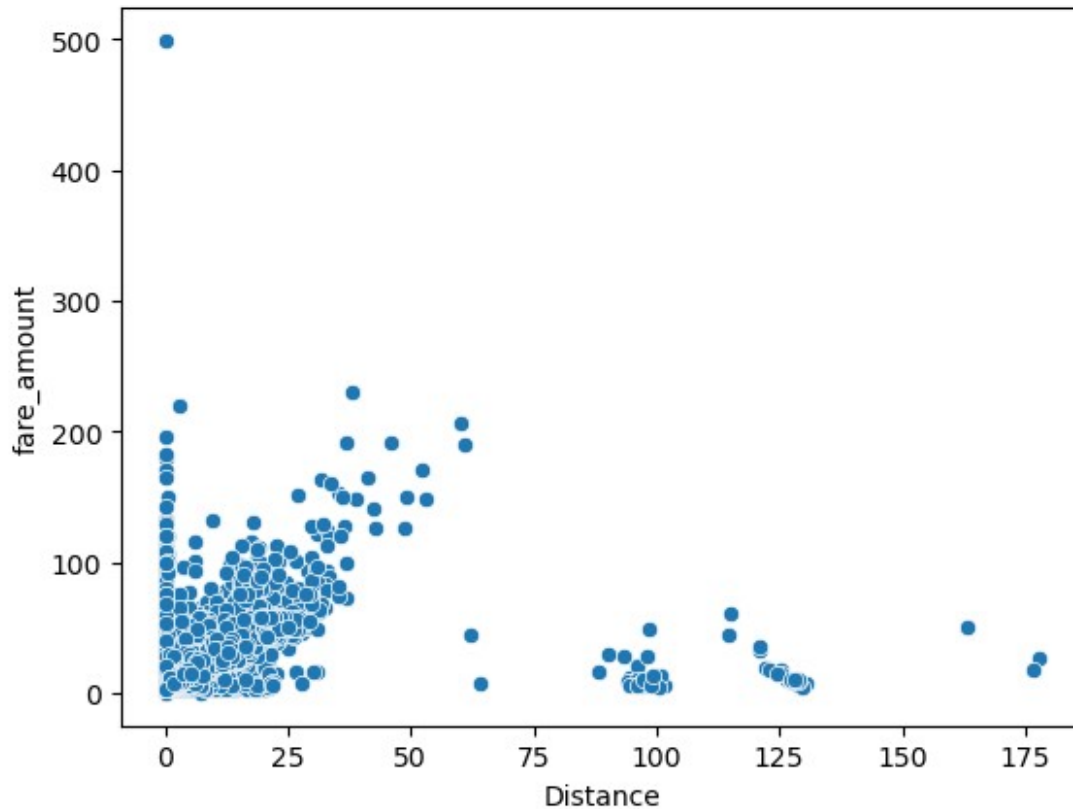  df['Hour'] = temp['Hour'].apply(convert_hour)
```

|   | fare_amount | passenger_count | Distance | week_day | Year | Month | Hour |
|---|---|---|---|---|---|---|---|
| 0 | 7.5 | 1 | 1.683323 | 0 | 2015 | 5 | 3 |
| 1 | 7.7 | 1 | 2.457590 | 1 | 2009 | 7 | 3 |
| 2 | 12.9 | 1 | 5.036377 | 0 | 2009 | 8 | 3 |
| 3 | 5.3 | 3 | 1.661683 | 1 | 2009 | 6 | 1 |
| 4 | 16.0 | 5 | 4.475450 | 0 | 2014 | 8 | 2 |

```
df.corr()
```

|  | fare_amount | passenger_count | Distance | week_day | Year |
|---|---|---|---|---|---|
| fare_amount | 1.000000 | 0.011884 | 0.778667 | 0.002305 | 0.120430 |
| passenger_count | 0.011884 | 1.000000 | 0.005112 | 0.035882 | 0.005339 |
| Distance | 0.778667 | 0.005112 | 1.000000 | 0.014518 | 0.018617 |
| week_day | 0.002305 | 0.035882 | 0.014518 | 1.000000 | 0.006910 |
| Year | 0.120430 | 0.005339 | 0.018617 | 0.006910 | 1.000000 |
| Month | 0.024120 | 0.008818 | 0.007373 | -0.007328 | -0.115182 |
| Hour | -0.021078 | 0.013572 | -0.022691 | -0.078129 | 0.001131 |

|  | Month | Hour |
|---|---|---|
| fare_amount | 0.024120 | -0.021078 |
| passenger_count | 0.008818 | 0.013572 |
| Distance | 0.007373 | -0.022691 |
| week_day | -0.007328 | -0.078129 |
| Year | -0.115182 | 0.001131 |
| Month | 1.000000 | -0.005410 |
| Hour | -0.005410 | 1.000000 |

```
sns.scatterplot(y=df['fare_amount'],x=df['Distance'])
```

```
<Axes: xlabel='Distance', ylabel='fare_amount'>
```



```python
from sklearn.preprocessing import StandardScaler
x = df[['Distance']].values
y = df['fare_amount'].values.reshape(-1,1)


from sklearn.model_selection import train_test_split
x_train, x_test, y_train,y_test =
train_test_split(x,y,random_state=10)

std_x = StandardScaler()
x_train = std_x.fit_transform(x_train)

x_test = std_x.transform(x_test)

std_y = StandardScaler()
y_train = std_y.fit_transform(y_train)

y_test = std_y.transform(y_test)
```

```python
from sklearn.metrics import mean_squared_error,r2_score,
mean_absolute_error
def fit_predict(model):
    model.fit(x_train,y_train.ravel())
    y_pred = model.predict(x_test)
    r_squared = r2_score(y_test,y_pred)
    RMSE = mean_squared_error(y_test, y_pred,squared=False)
    MAE = mean_absolute_error(y_test,y_pred)
    print('R-squared: ', r_squared)
    print('RMSE: ', RMSE)
    print("MAE:  ",MAE)

from sklearn.linear_model import LinearRegression

fit_predict(LinearRegression())
```

```
R-squared:  0.604116792084117
RMSE:  0.6290054895695945
MAE:    0.2755232959095982
```

```python
from sklearn.ensemble import RandomForestRegressor
fit_predict(RandomForestRegressor())
```

```
R-squared:  0.6526968425632438
RMSE:  0.589149157739537
MAE:    0.2919005686906083
```

```python
model = RandomForestRegressor()
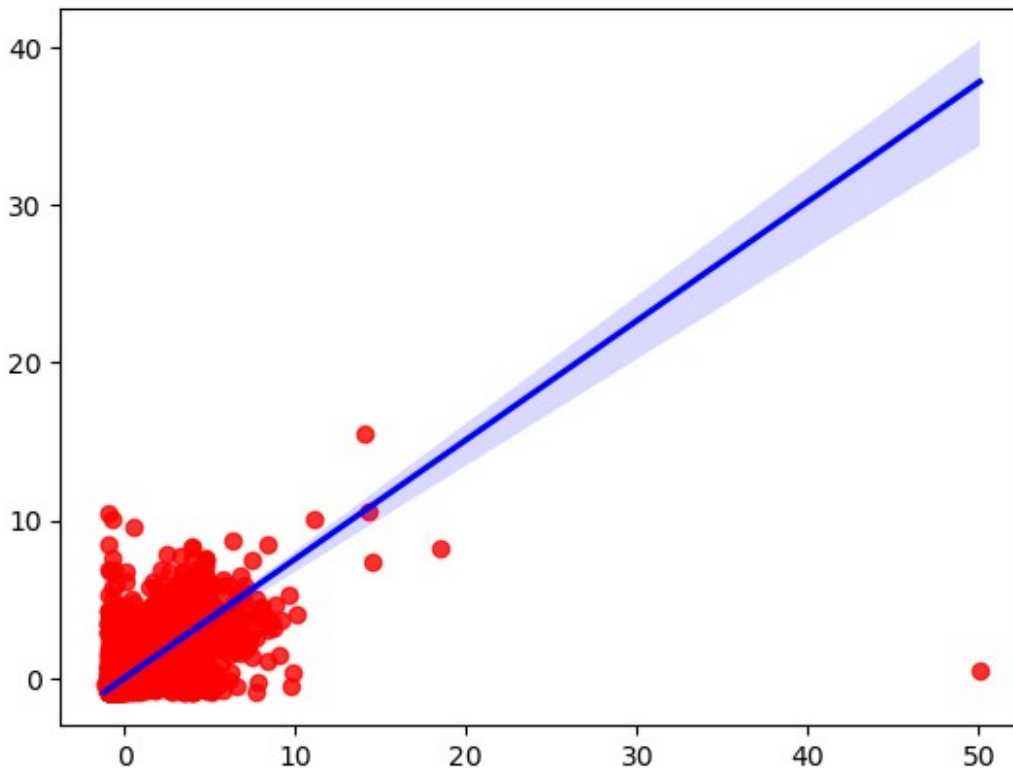model.fit(x_train, y_train)
```

```
C:\Users\prajw\anaconda3\Lib\site-packages\sklearn\base.py:1151:
DataConversionWarning: A column-vector y was passed when a 1d array
was expected. Please change the shape of y to (n_samples,), for
example using ravel().
  return fit_method(estimator, *args, **kwargs)
```

```
RandomForestRegressor()
```

```python
y_pred = model.predict(x_test)

sns.regplot(x=y_test, y=y_pred, color="red", line_kws={"color" :
"blue"})
plt.show()
```

```
model = LinearRegression()
model.fit(x_train, y_train)
scaler = StandardScaler()

---------------------------------------------------------------------
-----
NameError                                 Traceback (most recent call
last)
Cell In[1], line 1
----> 1 model = LinearRegression()
      2 model.fit(x_train, y_train)
      3 scaler = StandardScaler()

NameError: name 'LinearRegression' is not defined

y_pred = model.predict(x_test)

sns.regplot(x=y_test, y=y_pred, color="red", line_kws={"color" :
"blue"})
plt.show()
```