Before you turn in the homework, make sure everything runs as expected. To do so, select **Kernel**→**Restart & Run All** in the toolbar above. Remember to submit both on **DataHub** and **Gradescope**.

Please fill in your name and include a list of your collaborators below.

```
In [1]:    1  NAME = "Junsheng Pei"
           2  COLLABORATORS = ""
```

# Project 2: NYC Taxi Rides

## Extras

Put all of your extra work in here. Feel free to save figures to use when completing Part 4.

```
In [2]:    1  import os
           2  import pandas as pd
           3  import numpy as np
           4  from pathlib import Path
           5  from sqlalchemy import create_engine
           6  from utils import timeit
           7  import matplotlib.pyplot as plt
           8  import seaborn as sns
```

```
In [3]:    1  !ls -lh /srv/db/taxi_2016_student_small.sqlite
```

```
-rw-r--r-- 1 root root 2.1G Nov  7 04:44 /srv/db/taxi_2016_student_small.sqlite
```

```
In [4]:    1  DB_URI = "sqlite:////srv/db/taxi_2016_student_small.sqlite"
           2  TABLE_NAME = "taxi"
           3
           4  sql_engine = create_engine(DB_URI)
           5  with timeit():
           6      print(f"Table {TABLE_NAME} has {sql_engine.execute(f'SELECT COUNT(*) FROM {TABLE_NAME}').first()[0]} rows!")
```

```
Table taxi has 15000000 rows!
1.15 s elapsed
```

## Data Selection:

we chose data from April instead of January, and then we export the data

```
In [5]:    1  query = f"""
           2          SELECT *
           3          FROM (
           4          SELECT *
           5          FROM (
           6          SELECT *
           7          FROM (
           8          SELECT *
           9          FROM taxi
          10          WHERE tpep_pickup_datetime
          11              BETWEEN '2016-04-01' AND '2016-04-30'
          12              AND record_id % 100 == 0
          13          ORDER BY tpep_pickup_datetime
          14          )
          15          WHERE (julianday(tpep_dropoff_datetime) - julianday(tpep_pickup_datetime)) < 0.5
          16          )
          17          )
          18          WHERE passenger_count > 0
          19          """
          20
          21  with timeit(): # this query should take less than a second
          22      cleaned_df = pd.read_sql(query, sql_engine)
          23  cleaned_df['tpep_pickup_datetime'] = pd.to_datetime(cleaned_df['tpep_pickup_datetime'])
          24  cleaned_df['tpep_dropoff_datetime'] = pd.to_datetime(cleaned_df['tpep_dropoff_datetime'])
          25  cleaned_df['duration'] = cleaned_df["tpep_dropoff_datetime"]-cleaned_df["tpep_pickup_datetime"]
          26  cleaned_df['duration'] = cleaned_df['duration'].dt.total_seconds()
```

```
2.62 s elapsed
```

```
In [6]:    1  from sklearn.model_selection import train_test_split
           2  train_df, val_df = train_test_split(cleaned_df, test_size=0.2, random_state=42)
           3  data_file = Path("./", "cleaned_data_2016.hdf") # Path of hdf file
           4  train_df.to_hdf(data_file, "train") # Train data of hdf file
           5  val_df.to_hdf(data_file, "val") # Val data of hdf file
```

## Better features

We test the the following features: 'tip_amount','trip_distance','ifdaytime','ifweekday','total_amount','fare_amount.

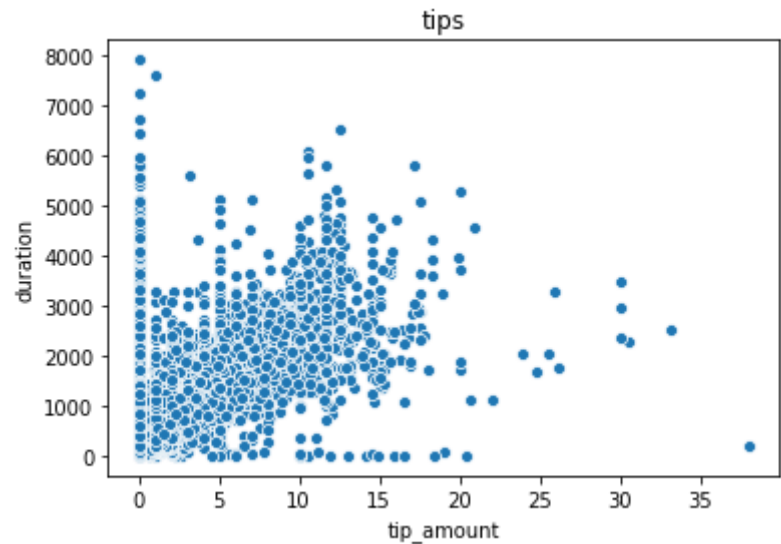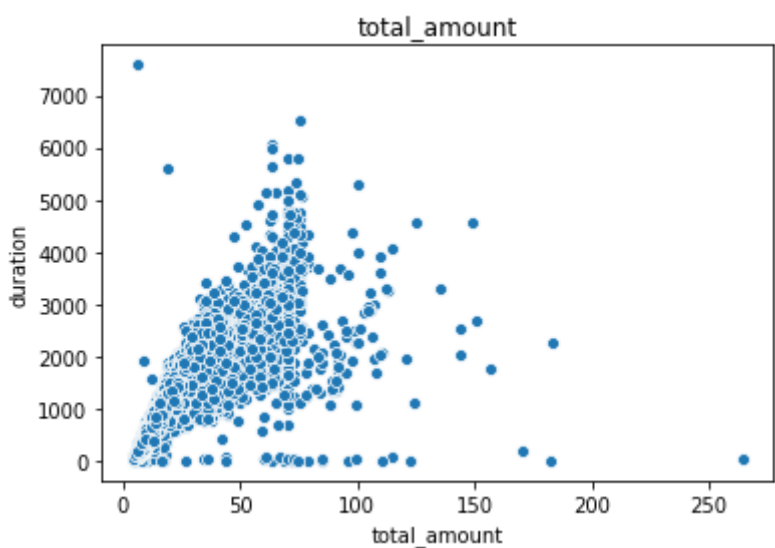And we found those feature can help us predict duration

```
In [7]:    1  cleaned_df = cleaned_df[cleaned_df['duration'] < 10000]
```
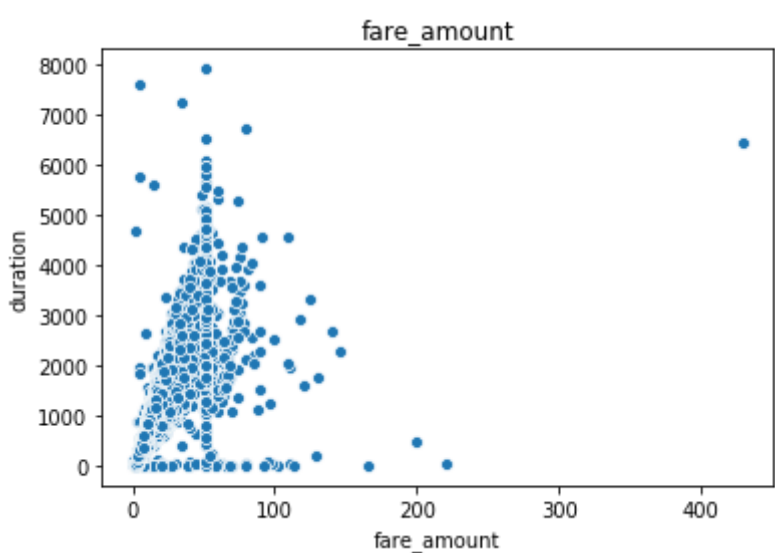
```
In [8]:    1  cleaned_df_feature = cleaned_df[cleaned_df['tip_amount'] < 40]
           2  sns.scatterplot('tip_amount','duration', data = cleaned_df_feature)
           3  plt.title('tips')
           4  plt.show()
```
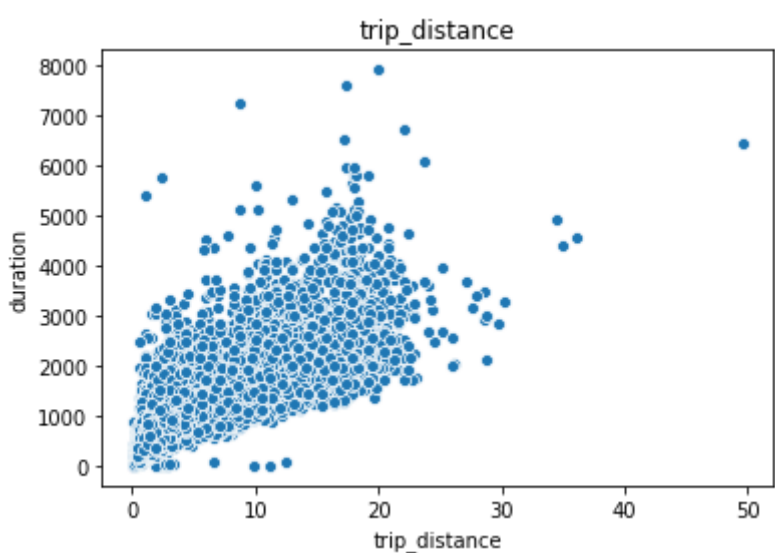
```
In [9]:   1  cleaned_df_feature = cleaned_df[cleaned_df['tip_amount'] > 0]
          2  sns.scatterplot('total_amount','duration', data = cleaned_df_feature)
          3  plt.title('total_amount')
          4  plt.show()
```
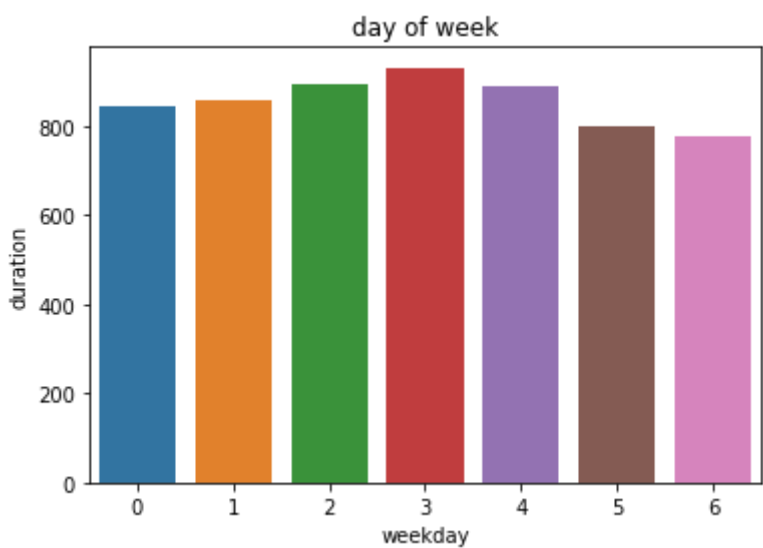


```
In [10]:  1  cleaned_df_feature = cleaned_df[cleaned_df['fare_amount'] > 0]
          2  sns.scatterplot('fare_amount','duration', data = cleaned_df_feature)
          3  plt.title('fare_amount')
          4  plt.show()
```
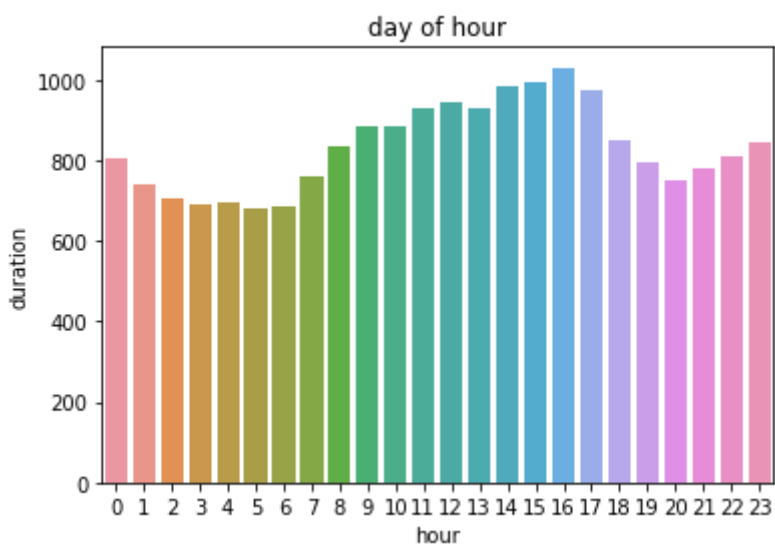


```
In [11]:  1  cleaned_df_feature = cleaned_df[cleaned_df['trip_distance'] > 0]
          2  sns.scatterplot('trip_distance','duration', data = cleaned_df_feature)
          3  plt.title('trip_distance')
          4  plt.show()
```



```
In [12]:  1  cleaned_df['weekday'] =  cleaned_df['tpep_pickup_datetime'].dt.dayofweek
          2  groupByweekday = cleaned_df.groupby('weekday')['duration'].mean()
          3  sns.barplot(x =groupByweekday.index, y = groupByweekday)
          4  plt.title('day of week')
          5  plt.show()
```



```
In [13]:  1  cleaned_df['hour'] =  cleaned_df['tpep_pickup_datetime'].dt.hour
          2  groupByhour = cleaned_df.groupby('hour')['duration'].mean()
          3  sns.barplot(x =groupByhour.index, y = groupByhour)
          4  plt.title('day of hour')
          5  plt.show()
```



## Feature Selection and Data Cleaning

```
In [14]:  1  data_file_fm = Path("./", "cleaned_data_2016.hdf")
          2  train_df_fm = pd.read_hdf(data_file_fm, "train")
          3  val_df_fm = pd.read_hdf(data_file_fm, "val")
```

```python
# Copied from part 2
def haversine(lat1, lng1, lat2, lng2):
    """
    Compute haversine distance
    """
    lat1, lng1, lat2, lng2 = map(np.radians, (lat1, lng1, lat2, lng2))
    average_earth_radius = 6371
    lat = lat2 - lat1
    lng = lng2 - lng1
    d = np.sin(lat * 0.5) ** 2 + np.cos(lat1) * np.cos(lat2) * np.sin(lng * 0.5) ** 2
    h = 2 * average_earth_radius * np.arcsin(np.sqrt(d))
    return h

# Copied from part 2
def manhattan_distance(lat1, lng1, lat2, lng2):
    """
    Compute Manhattan distance
    """
    a = haversine(lat1, lng1, lat1, lng2)
    b = haversine(lat1, lng1, lat2, lng1)
    return a + b

# Copied from part 2
def bearing(lat1, lng1, lat2, lng2):
    """
    Compute the bearing, or angle, from (lat1, lng1) to (lat2, lng2).
    A bearing of 0 refers to a NORTH orientation.
    """
    lng_delta_rad = np.radians(lng2 - lng1)
    lat1, lng1, lat2, lng2 = map(np.radians, (lat1, lng1, lat2, lng2))
    y = np.sin(lng_delta_rad) * np.cos(lat2)
    x = np.cos(lat1) * np.sin(lat2) - np.sin(lat1) * np.cos(lat2) * np.cos(lng_delta_rad)
    return np.degrees(np.arctan2(y, x))

# Copied from part 2
def add_time_columns(df):
    """
    Add temporal features to df
    """
    df.is_copy = False # propogate write to original dataframe
    df.loc[:, 'month'] = df['tpep_pickup_datetime'].dt.month
    df.loc[:, 'week_of_year'] = df['tpep_pickup_datetime'].dt.weekofyear
    df.loc[:, 'day_of_month'] = df['tpep_pickup_datetime'].dt.day
    df.loc[:, 'day_of_week'] = df['tpep_pickup_datetime'].dt.dayofweek
    df.loc[:, 'hour'] = df['tpep_pickup_datetime'].dt.hour
    df.loc[:, 'week_hour'] = df['tpep_pickup_datetime'].dt.weekday * 24 + df['hour']
    return df

# Copied from part 2
def add_distance_columns(df):
    """
    Add distance features to df
    """
    df.is_copy = False # propogate write to original dataframe
    df.loc[:, 'manhattan'] = manhattan_distance(lat1=df['pickup_latitude'],
                                                lng1=df['pickup_longitude'],
                                                lat2=df['dropoff_latitude'],
                                                lng2=df['dropoff_longitude'])

    df.loc[:, 'bearing'] = bearing(lat1=df['pickup_latitude'],
                                   lng1=df['pickup_longitude'],
                                   lat2=df['dropoff_latitude'],
                                   lng2=df['dropoff_longitude'])
    df.loc[:, 'haversine'] = haversine(lat1=df['pickup_latitude'],
                                       lng1=df['pickup_longitude'],
                                       lat2=df['dropoff_latitude'],
                                       lng2=df['dropoff_longitude'])
    return df

def select_columns(data, *columns):
    return data.loc[:, columns]

def mae(actual, predicted):
    """
    Calculates MAE from actual and predicted values
    Input:
      actual (1D array-like): vector of actual values
      predicted (1D array-like): vector of predicted/fitted values
    Output:
      a float, the MAE
    """

    mae = np.mean(np.abs(actual - predicted))
    return mae
```

```python
def add_ifdaytime(data):
    data['ifdaytime'] = (data['hour'] >= 8) & (data['hour'] <= 18)
    return data

def add_ifweekday(data):
    data['ifweekday'] = data['day_of_week'] > 4
    return data

def drop_outlier(data, col, _filter):
    return data.loc[data[col][lambda x: _filter(x)].index]

def replace_outlier(data, col, _filter):
    mean = data[col][lambda x : _filter(x)].mean()
    data[col] = data[col].apply(lambda x : x if _filter(x) else mean)
    return data
```

```python
In [17]:    1  def process_data_fm(data, test=False):
            2      # Put your final pipeline here
            3
            4      # data cleaning
            5      if(test):
            6          clean_data = replace_outlier
            7      else:
            8          clean_data = drop_outlier
            9          data = clean_data(data,'duration', lambda x : (x < 8000) & x > 0)
           10
           11      filter_latitude = lambda x : (x >= 40.63) & (x <= 40.85)
           12
           13      data = clean_data(data,'pickup_latitude', filter_latitude )
           14      data = clean_data(data,'dropoff_latitude', filter_latitude )
           15
           16      filter_longitude  = lambda x : (x >= -74.03) & (x <= -73.75)
           17
           18      data = clean_data(data,'pickup_longitude', filter_longitude)
           19      data = clean_data(data,'dropoff_longitude', filter_longitude )
           20
           21      data = clean_data(data,'total_amount', lambda x: (x>0) & (x <= 90))
           22      data = clean_data(data,'fare_amount', lambda x: (x>0) & (x <= 80))
           23      data = clean_data(data,'tip_amount', lambda x :(x>0) & (x <= 20) )
           24
           25      data = clean_data(data, 'trip_distance', lambda x : x < 50)
           26
           27      X = (
           28          data
           29          # Transform data
           30          .pipe(add_time_columns)
           31          .pipe(add_distance_columns)
           32          .pipe(add_ifdaytime)
           33          .pipe(add_ifweekday)
           34          .pipe(select_columns,
           35                  'pickup_longitude',
           36                  'pickup_latitude',
           37                  'dropoff_longitude',
           38                  'dropoff_latitude',
           39                  'manhattan',
           40                  'tip_amount',
           41                  'haversine',
           42                  'trip_distance',
           43                   'ifdaytime',
           44                  'ifweekday',
           45                  'total_amount',
           46                  'fare_amount',
           47                  )
           48      )
           49      if test:
           50          y = None
           51      else:
           52          y = data['duration']
           53
           54      return X, y
           55
           56  # YOUR CODE HERE
           57  #raise NotImplementedError()
```

## Parameter Section for Ridge and Lasso and Model Selection

```python
In [18]:    1  import sklearn.linear_model as lm
```

```python
In [19]:    1  # Parameter Section for Ridge
            2  _lambdas = [0.01,0.1,0.2,0.5,1,2,3,5,10]
            3  X_train_fm, y_train_fm = process_data_fm(train_df_fm)
            4  X_val_fm, y_val_fm = process_data_fm(val_df_fm)
            5
            6  #final_model = lm.LinearRegression(fit_intercept=True)
            7  for _lambda in _lambdas:
            8      final_model = lm.Ridge(alpha = _lambda, fit_intercept=True)
            9
           10      # Define your final model here, feel free to try other forms of regression
           11      final_model.fit(X_train_fm, y_train_fm)
           12
           13      y_val_pred_fm = final_model.predict(X_val_fm)
           14
           15      print("for lambda " + str(_lambda) + ", validation accuracy: " + str(mae(y_val_pred_fm,y_val_fm)))
```

```
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generic.py:4388: FutureWarning: Attribute 'is_copy' is deprecated and will be removed in a future ver
sion.
  object.__getattribute__(self, name)
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generic.py:4389: FutureWarning: Attribute 'is_copy' is deprecated and will be removed in a future ver
sion.
  return object.__setattr__(self, name, value)

for lambda 0.01, validation accuracy: 122.773109937
for lambda 0.1, validation accuracy: 122.715747715
for lambda 0.2, validation accuracy: 122.655868278
for lambda 0.5, validation accuracy: 122.502909548
for lambda 1, validation accuracy: 122.325242267
for lambda 2, validation accuracy: 122.109814835
for lambda 3, validation accuracy: 122.027044081
for lambda 5, validation accuracy: 122.060878093
for lambda 10, validation accuracy: 122.361729584
```

In [20]:
```python
# Parameter Section for Ridge
_lambdas = [0.01,0.1,0.2,0.5,1,2,3,5,10]
X_train_fm, y_train_fm = process_data_fm(train_df_fm)
X_val_fm, y_val_fm = process_data_fm(val_df_fm)

#final_model = lm.LinearRegression(fit_intercept=True)
for _lambda in _lambdas:
    final_model = lm.Lasso(alpha = _lambda, fit_intercept=True)

    # Define your final model here, feel free to try other forms of regression
    final_model.fit(X_train_fm, y_train_fm)

    y_val_pred_fm = final_model.predict(X_val_fm)

    print("for lambda " + str(_lambda) + ", validation accuracy: " + str(mae(y_val_pred_fm,y_val_fm)))
```

```
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generic.py:4388: FutureWarning: Attribute 'is_copy' is deprecated and will be removed in a future ver
sion.
  object.__getattribute__(self, name)
/srv/conda/envs/data100/lib/python3.6/site-packages/pandas/core/generic.py:4389: FutureWarning: Attribute 'is_copy' is deprecated and will be removed in a future ver
sion.
  return object.__setattr__(self, name, value)
/srv/conda/envs/data100/lib/python3.6/site-packages/sklearn/linear_model/coordinate_descent.py:491: ConvergenceWarning: Objective did not converge. You might want to
increase the number of iterations. Fitting data with very small alpha may cause precision problems.
  ConvergenceWarning)

for lambda 0.01, validation accuracy: 122.720904993
for lambda 0.1, validation accuracy: 122.359093687
for lambda 0.2, validation accuracy: 122.2259669
for lambda 0.5, validation accuracy: 122.37614613
for lambda 1, validation accuracy: 124.429678243
for lambda 2, validation accuracy: 124.532659775
for lambda 3, validation accuracy: 124.788272559
for lambda 5, validation accuracy: 125.625109785
for lambda 10, validation accuracy: 128.680014266
```

So we find using Ridge regression with lambda = 3 is the best

## Submission

You're almost done!

Before submitting this assignment, ensure that you have:

1. Restarted the Kernel (in the menubar, select Kernel→Restart & Run All)
2. Validated the notebook by clicking the "Validate" button.

Then,

1. **Submit** the assignment via the Assignments tab in **Datahub**
2. **Upload and tag** the manually reviewed portions of the assignment on **Gradescope**

In [20]:
```python
# Parameter Section for Ridge
_lambdas = [0.01,0.1,0.2,0.5,1,2,3,5,10]
X_train_fm, y_train_fm = process_data_fm(train_df_fm)
X_val_fm, y_val_fm = process_data_fm(val_df_fm)

#final_model = lm.LinearRegression(fit_intercept=True)
for _lambda in _lambdas:
    final_model = lm.Lasso(alpha = _lambda, fit_intercept=True)

    # Define your final model here, feel free to try other forms of regression
    final_model.fit(X_train_fm, y_train_fm)

    y_val_pred_fm = final_model.predict(X_val_fm)

    print("for lambda " + str(_lambda) + ", validation accuracy: " + str(mae(y_val_pred_fm,y_val_fm)))
```