

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

## Part 4: Feature Engineering and Model Fitting

In this final part of the project, you will finally build a regression model that attempts to predict the duration of a taxi ride from all other available information.

You will build this model using a processing pipeline and submit your results to Kaggle. We will first walk you through a generic example using the data we saved from Part 1. Please carefully follow these steps as you will need to repeat this for your final model. After, we give you free reign and let you decide how you want to define your final model.

In [2]:

```
1 import os
2 import pandas as pd
3 import numpy as np
4 import sklearn.linear_model as lm
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 from pathlib import Path
8 from sqlalchemy import create_engine
9 from sklearn.model_selection import cross_val_score, train_test_split, GridSearchCV
10
11 sns.set(style="whitegrid", palette="muted")
12
13 plt.rcParams['figure.figsize'] = (12, 9)
14 plt.rcParams['font.size'] = 12
15
16 %matplotlib inline
```

## Training and Validation

The following code loads the training and validation data from part 1 into a Pandas DataFrame.

In [3]:

```
1 # Run this cell to load the data.
2 data_file = Path("./", "cleaned_data.hdf")
3 train_df = pd.read_hdf(data_file, "train")
4 val_df = pd.read_hdf(data_file, "val")
```

In [4]:

```
1 train_df.head()
```

Out[4]:

	record_id	VendorID	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_distance	pickup_longitude	pickup_latitude	RatecodeID	store_and_fwd_flag	...	dropoff_latitude	payment_type	fare_amount	ex
	13242	5711100	2016-01-17 17:48:41	2016-01-17 17:55:53	1	1.00	-74.006470	40.738766	1	N	...	40.735664	1	6.5	
	12723	4989400	2016-01-17 01:18:39	2016-01-17 01:21:15	1	0.40	-73.989365	40.763000	1	N	...	40.766121	2	4.0	
	8508	2436400	2016-01-12 09:07:00	2016-01-12 09:41:17	1	11.40	-73.984108	40.774509	1	N	...	40.770458	1	37.0	
	21304	10899100	2016-01-29 09:07:54	2016-01-29 09:18:25	1	1.42	-74.002907	40.760262	1	N	...	40.742764	1	8.5	
	3817	1319400	2016-01-06 11:44:54	2016-01-06 11:49:55	1	0.80	-73.969742	40.760273	1	N	...	40.751129	2	5.0	

5 rows × 21 columns

## Testing

Here we load our testing data on which we will evaluate your model.

In [5]:

```
1 test_df = pd.read_csv("./proj2_test_data.csv")
2 test_df['tpep_pickup_datetime'] = pd.to_datetime(test_df['tpep_pickup_datetime'])
3 test_df.head()
```

Out[5]:

	record_id	VendorID	tpep_pickup_datetime	passenger_count	trip_distance	pickup_longitude	pickup_latitude	RatecodeID	store_and_fwd_flag	dropoff_longitude	dropoff_latitude	payment_type	fare_amount	extra	mta_tax
0	10000	1	2016-01-02 01:45:37	1	1.20	-73.982224	40.768620	1	N	-73.983765	40.779598	1	6.0	0.5	0.5
1	19000	2	2016-01-02 03:05:16	1	10.90	-73.999977	40.738121	1	N	-73.888657	40.824364	1	31.5	0.5	0.5
2	21000	1	2016-01-02 03:24:36	1	1.80	-73.986618	40.747379	1	N	-73.978508	40.729622	1	8.5	0.5	0.5
3	23000	2	2016-01-02 03:47:38	1	5.95	-74.002922	40.744572	1	N	-73.942413	40.786419	1	20.5	0.5	0.5
4	27000	1	2016-01-02 04:36:44	1	1.60	-73.986366	40.759464	1	N	-73.963081	40.760353	2	8.0	0.5	0.5

In [6]:

```
1 test_df.describe()
```

Out[6]:

	record_id	VendorID	passenger_count	trip_distance	pickup_longitude	pickup_latitude	RatecodeID	dropoff_longitude	dropoff_latitude	payment_type	fare_amount	extra	mta_tax	tip_amount
count	1.377400e+04	13774.000000	13774.000000	13774.000000	13774.000000	13774.000000	13774.000000	13774.000000	13774.000000	13774.000000	13774.000000	13774.000000	13774.000000	13774.000000
mean	3.465950e+07	1.536082	1.663642	2.954688	-72.953619	40.187999	1.043778	-73.055577	40.245056	1.340061	12.836930	0.333091	0.497985	1.805420
std	2.015133e+07	0.498714	1.311739	3.704427	8.628431	4.753186	0.877637	8.191366	4.512564	0.490019	10.707619	0.429590	0.036632	2.416784
min	1.000000e+04	1.000000	0.000000	0.000000	-77.039436	0.000000	1.000000	-77.039436	0.000000	1.000000	-93.300000	-0.500000	-0.500000	0.000000
25%	1.719975e+07	1.000000	1.000000	1.000000	-73.992058	40.735166	1.000000	-73.991318	40.734002	1.000000	6.500000	0.000000	0.500000	0.000000
50%	3.457400e+07	2.000000	1.000000	1.700000	-73.981846	40.752432	1.000000	-73.979897	40.753263	1.000000	9.500000	0.000000	0.500000	1.350000
75%	5.216875e+07	2.000000	2.000000	3.157500	-73.967119	40.767264	1.000000	-73.962749	40.768455	2.000000	14.500000	0.500000	0.500000	2.360000
max	6.940400e+07	2.000000	6.000000	104.800000	0.000000	40.868210	99.000000	0.000000	41.540859	4.000000	156.040000	4.500000	1.740000	40.000000

## Modeling

We've finally gotten to a point where we can specify a simple model. Remember that we will be fitting our model on the training set we created in part 1. We will use our validation set to evaluate how well our model might perform on future data.

### Reusable Pipeline

Throughout this assignment, you should notice that your data flows through a single processing pipeline several times. From a software engineering perspective, this should be sufficient motivation to abstract parts of our code into reusable functions/methods. We will now encapsulate our entire pipeline into a single function `process_data_gm`. `gm` is shorthand for "guided model".

```
In [7]: 1 # Copied from part 2
2 def haversine(lat1, lng1, lat2, lng2):
3     """
4     Compute haversine distance
5     """
6     lat1, lng1, lat2, lng2 = map(np.radians, (lat1, lng1, lat2, lng2))
7     average_earth_radius = 6371
8     lat = lat2 - lat1
9     lng = lng2 - lng1
10    d = np.sin(lat * 0.5) ** 2 + np.cos(lat1) * np.cos(lat2) * np.sin(lng * 0.5) ** 2
11    h = 2 * average_earth_radius * np.arcsin(np.sqrt(d))
12    return h
13
14 # Copied from part 2
15 def manhattan_distance(lat1, lng1, lat2, lng2):
16     """
17     Compute Manhattan distance
18     """
19     a = haversine(lat1, lng1, lat1, lng2)
20     b = haversine(lat1, lng1, lat2, lng1)
21     return a + b
22
23 # Copied from part 2
24 def bearing(lat1, lng1, lat2, lng2):
25     """
26     Compute the bearing, or angle, from (lat1, lng1) to (lat2, lng2).
27     A bearing of 0 refers to a NORTH orientation.
28     """
29     lng_delta_rad = np.radians(lng2 - lng1)
30     lat1, lng1, lat2, lng2 = map(np.radians, (lat1, lng1, lat2, lng2))
31     y = np.sin(lng_delta_rad) * np.cos(lat2)
32     x = np.cos(lat1) * np.sin(lat2) - np.sin(lat1) * np.cos(lat2) * np.cos(lng_delta_rad)
33     return np.degrees(np.arctan2(y, x))
34
35 # Copied from part 2
36 def add_time_columns(df):
37     """
38     Add temporal features to df
39     """
40     df.is_copy = False # propogate write to original dataframe
41     df.loc[:, 'month'] = df['tpep_pickup_datetime'].dt.month
42     df.loc[:, 'week_of_year'] = df['tpep_pickup_datetime'].dt.weekofyear
43     df.loc[:, 'day_of_month'] = df['tpep_pickup_datetime'].dt.day
44     df.loc[:, 'day_of_week'] = df['tpep_pickup_datetime'].dt.dayofweek
45     df.loc[:, 'hour'] = df['tpep_pickup_datetime'].dt.hour
46     df.loc[:, 'week_hour'] = df['tpep_pickup_datetime'].dt.weekday * 24 + df['hour']
47     return df
48
49 # Copied from part 2
50 def add_distance_columns(df):
51     """
52     Add distance features to df
53     """
54     df.is_copy = False # propogate write to original dataframe
55     df.loc[:, 'manhattan'] = manhattan_distance(lat1=df['pickup_latitude'],
56                                                  lng1=df['pickup_longitude'],
57                                                  lat2=df['dropoff_latitude'],
58                                                  lng2=df['dropoff_longitude'])
59
60     df.loc[:, 'bearing'] = bearing(lat1=df['pickup_latitude'],
61                                   lng1=df['pickup_longitude'],
62                                   lat2=df['dropoff_latitude'],
63                                   lng2=df['dropoff_longitude'])
64     df.loc[:, 'haversine'] = haversine(lat1=df['pickup_latitude'],
65                                       lng1=df['pickup_longitude'],
66                                       lat2=df['dropoff_latitude'],
67                                       lng2=df['dropoff_longitude'])
68     return df
69
70 def select_columns(data, *columns):
71     return data.loc[:, columns]
```

```
In [8]: 1 def process_data_gml(data, test=False):
2     X = (
3         data
4
5         # Transform data
6         .pipe(add_time_columns)
7         .pipe(add_distance_columns)
8
9         .pipe(select_columns,
10              'pickup_longitude',
11              'pickup_latitude',
12              'dropoff_longitude',
13              'dropoff_latitude',
14              'manhattan',
15              )
16     )
17     if test:
18         y = None
19     else:
20         y = data['duration']
21
22     return X, y
```

We will use our pipeline defined above to pre-process our training and test data in exactly the same way. Our functions make this relatively easy to do!

```
In [9]: 1 # Train
2 X_train, y_train = process_data_gml(train_df)
3 X_val, y_val = process_data_gml(val_df)
4 guided_model_1 = lm.LinearRegression(fit_intercept=True)
5 guided_model_1.fit(X_train, y_train)
6
7 # Predict
8 y_train_pred = guided_model_1.predict(X_train)
9 y_val_pred = guided_model_1.predict(X_val)
```

```
/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.__getattr__(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)
```

Here, `y_val` are the correct durations for each ride, and `y_val_pred` are the predicted durations based on the 7 features above (`vendorID`, `passenger_count`, `pickup_longitude`, `pickup_latitude`, `dropoff_longitude`, `dropoff_latitude`, `manhattan`).

```
In [10]: 1 assert 600 <= np.median(y_train_pred) <= 700
2         assert 600 <= np.median(y_val_pred) <= 700
```

The resulting model really is a linear model just like we saw in class, i.e. the predictions are simply generated by the product  $\Phi\theta$ . For example, the line of code below generates a prediction for  $x_1$  by computing  $\phi_1^T\theta$ . Here `guided_model_1.coef_` is  $\theta$  and `X_train.iloc[0, :]` is  $\phi_1$ .

Note that unlike in class, here the dummy intercept term is not included in  $\Phi$ .

```
In [11]: 1 X_train.iloc[0, :].dot(guided_model_1.coef_) + guided_model_1.intercept_

Out[11]: 558.751330511368
```

We see that this prediction is exactly the same (except for possible floating point error) as generated by the `predict` function, which simply computes the product  $\Phi\theta$ , yielding predictions for every input.

```
In [12]: 1 y_train_pred[0]

Out[12]: 558.75133051135344
```

In this assignment, we will use Mean Absolute Error (MAE), a.k.a. mean L1 loss, to measure the quality of our models. As a reminder, this quantity is defined as:

$$MAE = \frac{1}{n} \sum_i |y_i - \hat{y}_i|$$

Why may we want to use the MAE as a metric, as opposed to Mean Squared Error (MSE)? Using our domain knowledge that most rides are short in duration (median is roughly 600 seconds), we know that MSE is susceptible to outliers. Given that some of the outliers in our dataset are quite extreme, it is probably better to optimize for the majority of rides rather than for the outliers. You may want to remove some of these outliers later on.

```
In [13]: 1 def mae(actual, predicted):
2     """
3     Calculates MAE from actual and predicted values
4     Input:
5         actual (1D array-like): vector of actual values
6         predicted (1D array-like): vector of predicted/fitted values
7     Output:
8         a float, the MAE
9     """
10
11     mae = np.mean(np.abs(actual - predicted))
12     return mae
```

```
In [14]: 1 assert 200 <= mae(y_val_pred, y_val) <= 300
2 print("Validation Error: ", mae(y_val_pred, y_val))
```

Validation Error: 266.136130855

Side note: scikit-learn also has tools to compute mean absolute error ( `sklearn.metrics.mean_absolute_error` ). In fact, most metrics that we have discussed in this class can be found as part of the [sklearn.metrics module](https://scikit-learn.org/stable/modules/classes.html#sklearn-metrics-metrics) (https://scikit-learn.org/stable/modules/classes.html#sklearn-metrics-metrics). Some of these may come in handy as part of your feature engineering!

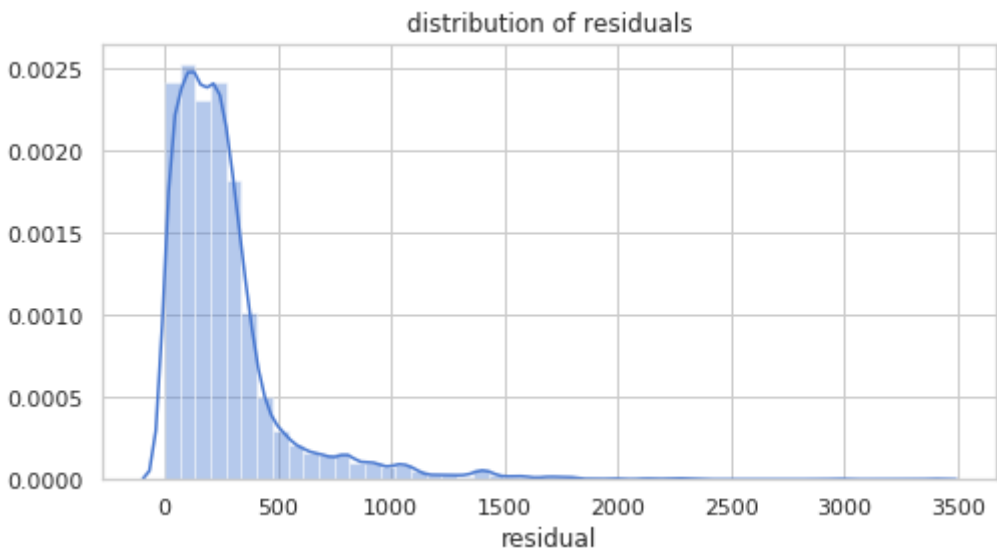
## Visualizing Error

You should be getting between 200 and 300 MAE, which means your model was off by roughly 3-5 minutes on trips of average length 12 minutes. This is fairly decent performance given that our basic model uses only using the pickup/dropoff latitude and manhattan distance of the trip. 3-5 minutes may seem like a lot for a trip of 12 minutes, but keep in mind that this is the *average* error. This metric is susceptible to extreme outliers, which exist in our dataset.

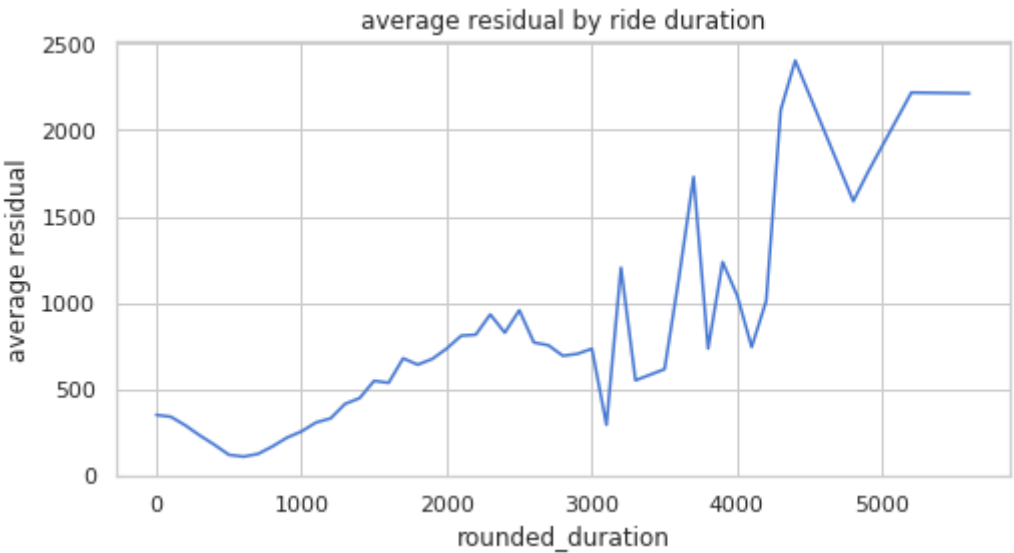
Now we will visualize the residual for the validation set. We will plot the following:

1. Distribution of residuals
2. Average residual grouping by ride duration

```
In [15]: 1 # Distribution of residuals
2 plt.figure(figsize=(8,4))
3 sns.distplot(np.abs(y_val - y_val_pred))
4 plt.xlabel('residual')
5 plt.title('distribution of residuals');
```



```
In [16]: 1 # Average residual grouping by ride duration
2 val_residual = X_val.copy()
3 val_residual['duration'] = y_val
4 val_residual['rounded_duration'] = np.around(y_val, -2)
5 val_residual['residual'] = np.abs(y_val - y_val_pred)
6 tmp = val_residual.groupby('rounded_duration').mean()
7 plt.figure(figsize=(8,4))
8 tmp['residual'].plot()
9 plt.ylabel('average residual')
10 plt.title('average residual by ride duration');
```



In the first visualization, we see that most of the residuals are centered around 250 seconds ~ 4 minutes. There is a minor right tail, suggesting that we are still unable to accurately fit some outliers in our data. The second visualization also suggests this, as we see the average residual increasing as a somewhat linear function of duration. But given that our average ride duration is roughly 600-700 seconds, it seems that we are indeed optimizing for the average ride because the residuals are smallest around 600-700.

Keep this in mind when creating your final model! Visualizing the error is a powerful tool and may help diagnose shortcomings of your model. Let's go ahead and submit to kaggle, although your error on the test set may be higher than 300.

## Submission to Kaggle

The following code will write your predictions on the test dataset to a CSV, which you can submit to Kaggle. You may need to modify it to suit your needs, but we recommend you make a copy and preserve the original function.



Remember that if you've performed transformations or featurization on the training data, you must also perform the same transformations on the test data in order to make predictions. For example, if you've created features for the columns `pickup_datetime` or `pickup_latitude` on the training data, you must also extract the same features in order to use scikit-learn's `.predict(...)` method.

```
In [17]: 1 from datetime import datetime
2 def generate_submission(test, predictions, force=False):
3     if force:
4         if not os.path.isdir("submissions"):
5             os.mkdir("submissions")
6         submission_df = pd.DataFrame({
7             "id": test_df.index.values,
8             "duration": predictions,
9         },
10         columns=['id', 'duration'])
11
12     timestamp = datetime.isoformat(datetime.now()).split(".")[0]
13
14     submission_df.to_csv(f'submissions/submission_{timestamp}.csv', index=False)
15
16     print(f'Created a CSV file: submission_{timestamp}.csv')
17     print('You may now upload this CSV file to Kaggle for scoring.')
```

```
In [18]: 1 X_test, _ = process_data_gml(test_df, True)

/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.__getattr__(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)
```

```
In [19]: 1 assert list(X_train.columns) == list(X_test.columns), "Different columns or different column ordering"
2 submission_predictions = (guided_model_1
3                             .fit(X_train, y_train)
4                             .predict(X_test))
5 submission_predictions = submission_predictions.astype(int)
6 submission_predictions[submission_predictions < 0] = 0
7 generate_submission(test_df, submission_predictions, True)
```

Created a CSV file: submission\_2018-12-04T22:28:34.csv  
You may now upload this CSV file to Kaggle for scoring.

```
In [20]: 1 # Check your submission
2 assert isinstance(submission_predictions, np.ndarray), "Submission not an array"
3 assert all(submission_predictions >= 0), "Duration must be non-negative"
4 assert issubclass(submission_predictions.dtype.type, np.integer), "Seconds must be integers"
```

## Your Turn!

Now it's your turn! Draw upon everything you have learned this semester to find the best features to help your model accurately predict the duration of a taxi ride.

You may use whatever method you prefer in order to create features. You may use features that we created and features that you discovered yourself from any of the 2 datasets. However, we want to make it fair to students who are seeing these techniques for the first time. As such, you are only allowed regression models and their regularized forms. This means no random forest, k-nearest-neighbors, neural nets, etc.

Here are some ideas to improve your model:

- **Data selection:** January 2016 was an odd month for taxi rides due to the blizzard. Would it help to select training data differently?
- **Data cleaning:** Try cleaning your data in different ways. In particular, consider how to handle outliers.
- **Better features:** Explore the 2 datasets and find what features are most helpful. Utilize external datasets to improve your accuracy.
- **Regularization:** Try different forms of regularization to avoid fitting to the training set. Recall that `Ridge` and `Lasso` are the names of the classes in `sklearn.linear_model` that combine `LinearRegression` with regularization techniques.
- **Model selection:** You can adjust parameters of your model (e.g., the regularization parameter) to achieve higher accuracy. [GridSearchCV](http://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html) ([http://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.GridSearchCV.html](http://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html)) may be helpful.
- **Validation:** Recall that you should use cross-validation to do feature and model selection properly! Otherwise, you will likely overfit to your training data.

There's many things you could try that could help your model. We have only suggested a few. Be creative and innovative! Please use `proj2_extras.ipynb` for all of your extraneous work. Note that you will be submitting `proj2_extras.ipynb` and we will be grading it. Please properly comment and format this notebook!

Once you are satisfied with your results, answer the questions in the Deliverables section. You may want to read this section in advance so you have an idea of what we're looking for.

## Deliverables

### Feature/Model Selection Process

Let's first look at selection of better features. In this following cell, describe the process of choosing good features to improve your model. You should use at least 3-4 sentences each to address the follow questions. Backup your responses with graphs supporting your claim (you can save figures and load them, no need to add the plotting code here). Use these questions to concisely summarize all of your extra work!

#### Question 1a

How did you find better features for your model?

```
In [21]: 1 q1a_answer = r"""
2 I testet more columns in the given dataset by considering each coloumn as a feature and visualizing the
3 relationship between such feature and duration.
4
5 If the values of the feature are discrete, I use barplot. And if the duration distincts from each value, I add
6 this feature to my feature matirx.
7
8 If the values of the feature are continuous, I use scatter plot(we can also find outliers with scatter plot).
9 And if the duration have a linear relationship with such feature, I add this feature to my feature matirx.
10
11 If the feature helps improve the accuray, it's a good feature.
12
13 """
14 # YOUR CODE HERE
15 #raise NotImplementedError()
```

#### Question 1b

What did you try that worked / didn't work?

```
In [22]: 1 q1b_answer = r"""
2
3 I found that 'tip_amount', 'haversine', 'trip_distance', 'ifdaytime', 'ifweekday', 'total_amount' and
4 'fare_amount' worked.
5
6 However, record_id', 'VendorID', passenger_count', 'RatecodeID', 'store_and_fwd_flag', 'payment_type',
7 'improvement_surcharge', 'total_amount', 'month', 'week_of_year' and 'day_of_month' didn't work.
8
9 """
10 # YOUR CODE HERE
11 #raise NotImplementedError()
```

#### Question 1c

What was surprising in your search for good features?

```
In [23]: 1 g1c_answer = r"""
2 I found that 'total_amount' is extremely helpful, it is proportional to the duration.
3 And I found replacing the outliers with average value can make the feature more helpful, like pickup_latitude,
4 pickup_longtitude have many values as 0.
5
6 """
7 # YOUR CODE HERE
8 #raise NotImplementedError()
```

Question 2

Just as in the guided model above, you should encapsulate as much of your workflow into functions as possible. Define `process_data_fm` and `final_model` in the cell below. In order to calculate your final model's MAE, we will run the code in the cell after that.

**Note:** You *MUST* name the model you wish to be evaluated on `final_model` . This is what we will be using to generate your predictions. We will take the state of `final_model` right after executing the cell below and run the following code:

```
# Load in test_df, solutions
X_test, _ = process_data_fm(test_df, True)
submission_predictions = final_model.predict(X_test)
# Generate score for autograding
```

We encourage you to conduct all of your exploratory work in `proj2_extras.ipynb` , which will be graded for 10 points.

```
In [24]: 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")
```

```
In [25]: 1 def add_ifdaytime(data):
2     data['ifdaytime'] = (data['hour'] >= 8) & (data['hour'] <= 18)
3     return data
4
5 def add_ifweekday(data):
6     data['ifweekday'] = data['day_of_week'] > 4
7     return data
8
9 def drop_outlier(data, col, _filter):
10     return data.loc[data[col][lambda x: _filter(x)].index]
11
12 def replace_outlier(data, col, _filter):
13     mean = data[col][lambda x: _filter(x)].mean()
14     data[col] = data[col].apply(lambda x: x if _filter(x) else mean)
15     return data
16
17
18
```

```
In [26]: 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()
```

In [27]:

```
1 X_train_fm, y_train_fm = process_data_fm(train_df_fm)
2 X_val_fm, y_val_fm = process_data_fm(val_df_fm)
3
4 #final_model = lm.LinearRegression(fit_intercept=True)
5 final_model = lm.Ridge(alpha = 3, fit_intercept=True)
6
7 # Define your final model here, feel free to try other forms of regression
8 final_model.fit(X_train_fm, y_train_fm)
9
10 y_train_pred_fm = final_model.predict(X_train_fm)
11 y_val_pred_fm = final_model.predict(X_val_fm)
12
13 print(mae(y_train_pred_fm,y_train_fm))
14 print(mae(y_val_pred_fm,y_val_fm))
```

122.552163418  
122.027044081

/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 version.  
object.\_\_getattr\_\_(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 version.  
return object.\_\_setattr\_\_(self, name, value)

In [28]:

```
1 # Feel free to change this cell
2 X_test, _ = process_data_fm(test_df, True)
3 final_predictions = final_model.predict(X_test)
4 final_predictions = final_predictions.astype(int)
5 generate_submission(test_df, final_predictions, True) # Change to true to generate prediction
```

Created a CSV file: submission\_2018-12-04T22:28:35.csv  
You may now upload this CSV file to Kaggle for scoring.

/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 version.  
object.\_\_getattr\_\_(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 version.  
return object.\_\_setattr\_\_(self, name, value)

### Question 3

The following hidden cells will test your model on the test set. Please do not delete any of them if you want credit!

In [29]:

```
1 # NO TOUCH
```

In [30]:

```
1 # NOH
```

In [31]:

```
1 # STAHP
```

In [32]:

```
1 # NO MOLESTE
```

In [33]:

```
1 # VA-T'EN
```

In [34]:

```
1 # NEIN
```

In [35]:

```
1 # PLSNO
```

In [36]:

```
1 # THIS SPACE IS NOT YOURS
```

In [37]:

```
1 # TAWDEETAW
```

In [38]:

```
1 # MAU LEN
```

In [39]:

```
1 # ALMOST
```

In [40]:

```
1 # TO
```

In [41]:

```
1 # THE
```

In [42]:

```
1 # END
```

In [43]:

```
1 # Hmph
```

In [44]:

```
1 # Good riddance
```

In [45]:

```
1 generate_submission(test_df, submission_predictions, True)
```

Created a CSV file: submission\_2018-12-04T22:28:35.csv  
You may now upload this CSV file to Kaggle for scoring.

This should be the format of your CSV file.  
Unix-users can verify it running `!head submission_{datetime}.csv` in a jupyter notebook cell.

id,duration  
id3004672,965.3950873305439  
id3505355,1375.0665915134596  
id1217141,963.2285454171943  
id2150126,1134.7680929570924  
id1598245,878.5495792656438  
id0668992,831.6700312449248  
id1765014,993.1692116960185  
id0898117,1091.1171629594755  
id3905224,887.9037911118357

Kaggle link: <https://www.kaggle.com/t/f8b3c6acc3a045cab152060a5bc79670> (<https://www.kaggle.com/t/f8b3c6acc3a045cab152060a5bc79670>)

### 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**