

```
In [274]: 1 import pandas as pd
2 import numpy as np
3 from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier
4 from sklearn.model_selection import train_test_split
5 from sklearn.neighbors import KNeighborsRegressor
6 from sklearn.impute import SimpleImputer
7 from sklearn.model_selection import train_test_split, GridSearchCV
8 from sklearn.linear_model import LogisticRegression, SGDClassifier
9 from sklearn.ensemble import VotingClassifier, ExtraTreesClassifier, RandomForestClassifier
10 from sklearn.svm import SVC
11 from sklearn.metrics import mean_squared_error
12 from sklearn.neural_network import MLPClassifier
13
14 from sklearn.linear_model import LinearRegression
15 from sklearn.preprocessing import FunctionTransformer
16 from sklearn.preprocessing import StandardScaler, OneHotEncoder
17 from sklearn.pipeline import Pipeline
18 from sklearn.compose import ColumnTransformer
19
20 from sklearn.tree import DecisionTreeRegressor
21 from sklearn.model_selection import train_test_split
22 from sklearn.neighbors import KNeighborsRegressor
23 from sklearn.model_selection import train_test_split
24 pd.set_option('display.max_columns', 500)
```

## Part 1, extract the mean time grouby (origin,dest,weekday,timeslot)

```
In [275]: 1 data=pd.read_csv('https://raw.githubusercontent.com/Xiao-Wang-UCSD/dataset/master/20170101-20170102.csv')
2
```

```
In [276]: 1 data['Origin Name'].unique()
```

```
Out[276]: array(['The Palace Of Fine Arts, 3601 Lyon St, San Francisco, CA',
                  'Embarcadero, San Francisco, CA',
                  'Oracle Park, 24 Willie Mays Plaza, San Francisco, CA',
                  "Fisherman's Wharf, 286-298 Jefferson St, San Francisco, CA",
                  '2nd Street and Stevenson Street, San Francisco, CA',
                  'Powell BART Station, Market St and Powell St, San Francisco, CA'],
              dtype=object)
```

```
In [518]: 1
2 def cleaning2(data,name):
3     col=data['Daily Mean Travel Time']
4     #extract the percentile of col in col
5     func=lambda a:percentileofscore(col,a)
6     percent=col.apply(func)
7     new=data[name].dropna()
8     #apply the percentile of col to name
9     func2=lambda b:np.percentile(new, b)
10    value=percent.apply(func2)
11    data[name]=data[name].fillna(value)
12    return data
13
14
15
```

```
In [ ]: 1 <option value="The Palace Of Fine Arts, 3601 Lyon St, San Francisco, CA"
2 <option value="Fisherman's Wharf, 286-298 Jefferson St, San Francisco, C
3 <option value="Oracle Park, 24 Willie Mays Plaza, San Francisco, CA">Ora
4 <option value="Embarcadero, San Francisco, CA">Embarcadero Bart</option>
5 <option value="2nd Street and Stevenson Street, San Francisco, CA">Montg
6 <option value="Powell BART Station, Market St and Powell St, San Francis
```

```
In [586]: 1 data=pd.read_csv('https://raw.githubusercontent.com/Xiao-Wang-UCSD/dat
2 data['Destination Name']=data['Destination Name'].replace({'2nd Street
3
4 column_name=['AM Mean Travel Time', 'AM Range - Lower',
5             'AM Range - Upper', 'PM Mean Travel Time',
6             'PM Range - Lower Bound Travel Time',
7             'PM Range - Upper Bound Travel Time', 'Midday Mean Travel Time'
8             'Midday Range - Lower', 'Midday Range - Upper',
9             'Evening Mean Travel Time', 'Evening Range - Lower',
10            'Evening Range - Upper', 'Early Morning Mean Travel Time',
11            'Early Morning Range - Lower', 'Early Morning Range - Upper']
12 for i in column_name:
13     data=cleaning2(data,i)
14
```

```
In [ ]: 1
```

```
In [576]: 1 data1=data[['Origin Name','Destination Name','AM Mean Travel Time','PM
2 name=list(data1.groupby(['Origin Name','Destination Name','weekday']).
3 frame=data1.groupby(['Origin Name','Destination Name','weekday']).mean
4
```

```
In [577]: 1 frame=frame.reset_index()
2 frame['weekday']=frame['weekday'].replace({0:'Mon',1:'Tues',2:'Wed',3:
```

```
In [578]: 1 # temp = frame.reset_index()
2 # temp = temp.groupby('Origin Name')[['Destination Name','AM Mean Trav
3 # temp
4
```

In [579]:

1	frame							
	Origin Name	St, San F	Weekday	Travel Time	Travel Time	Mean Travel Time	Mean Travel Time	M
	Francisco...	Name						
	2nd Street and Stevenson Street, San Francisco...	Oracle Park, 24 Willie Mays Plaza, San Francis...	Mon	443.320000	586.480000	537.400000	411.080000	3
7								
	2nd Street and Stevenson Street, San Francisco...	Oracle Park, 24 Willie Mays Plaza, San Francis...	Tues	467.769231	657.961538	526.461538	443.615385	3
8								
	2nd Street and Stevenson Street, San Francisco...	Oracle Park, 24 Willie Mays Plaza, San Francis...	Wed	512.153846	696.769231	584.884615	451.423077	3
9								

In [ ]:

1

In [580]:

```
1 a=frame['Origin Name'].unique()  
2 b=frame['Destination Name'].unique()  
3 c=['Mon', 'Tues', 'Wed', 'Thur', 'Fri', 'Sat', 'Sun']
```

```
In [581]: 1 dict={}
          2 for aa in a:
          3     dict[aa]={}
          4     for bb in b:
          5         dict[aa][bb]={}
          6         for cc in c:
          7             dict[aa][bb][cc]=0
          8
          9
         10 dict
         11
```

```
Out[581]: {'2nd Street and Stevenson Street, San Francisco, CA': {"Fisherman's Wharf, 286-298 Jefferson St, San Francisco, CA": {'Mon': 0,
    'Tues': 0,
    'Wed': 0,
    'Thur': 0,
    'Fri': 0,
    'Sat': 0,
    'Sun': 0},
    'Oracle Park, 24 Willie Mays Plaza, San Francisco, CA': {'Mon': 0,
    'Tues': 0,
    'Wed': 0,
    'Thur': 0,
    'Fri': 0,
    'Sat': 0,
    'Sun': 0},
    'The Palace Of Fine Arts, 3601 Lyon St, San Francisco, CA': {'Mon': 0,
    'Tues': 0,
    'Wed': 0,
    'Thur': 0,
    'Fri': 0,
    'Sat': 0,
    'Sun': 0}}
```

```
In [582]: 1
2 for inx in range(frame.shape[0]):
3     dict[frame.iloc[inx,0]][frame.iloc[inx,1]][frame.iloc[inx,2]]=list
4
5 dict
689.6538461538462,
511.38461538461536],
'Wed': [875.3076923076923,
961.7307692307693,
782.8076923076923,
708.9615384615385,
529.7307692307693],
'Thur': [931.9615384615385,
981.7307692307693,
814.7692307692307,
749.5,
512.9230769230769],
'Fri': [860.1923076923077,
950.9230769230769,
860.0384615384615,
753.7692307692307,
543.8461538461538],
'Sat': [668.3846153846154,
875.2692307692307,
822.7307692307693]
```

```
In [ ]: 1
```

```
In [583]: 1 import json
2
3 t=json.dumps(dict,indent=5)
4 with open("mean.json", "w") as outfile:
5     outfile.write(t)
```

## 2 machine learning

```
In [741]: 1 data=pd.read_csv('https://raw.githubusercontent.com/Xiao-Wang-UCSD/data/master/mean.json')
2 data['month']=data['Date'].apply(lambda a:pd.to_datetime(a).month)
3
4 column_name=['AM Mean Travel Time', 'AM Range - Lower',
5             'AM Range - Upper', 'PM Mean Travel Time',
6             'PM Range - Lower Bound Travel Time',
7             'PM Range - Upper Bound Travel Time', 'Midday Mean Travel Time',
8             'Midday Range - Lower', 'Midday Range - Upper',
9             'Evening Mean Travel Time', 'Evening Range - Lower',
10            'Evening Range - Upper', 'Early Morning Mean Travel Time',
11            'Early Morning Range - Lower', 'Early Morning Range - Upper']
12 for i in column_name:
13     data=cleaning2(data,i)
14
```

```
In [742]: 1 data['month']=data['month'].apply(lambda a: a+1)
          2
```

```
In [743]: 1 # data1=data[['Origin Name','Destination Name','AM Mean Travel Time','
          2
          3 # data1=data[['Origin Name','Destination Name','AM Mean Travel Time','
          4
          5
          6 data['Destination Name']=data['Destination Name'].replace({'2nd Street
          7 data['Origin Name']=data['Origin Name'].replace({'2nd Street and Steve
          8 # data.to_csv('out.csv')
          9
         10
         11
```

```
In [744]: 1 neww=pd.DataFrame([])
          2 for i in range(data.shape[0]):
          3     neww=neww.append(data[['Date']].loc[[i,i,i,i,i ] ])
          4 neww
```

0 2019-06-09

0 2019-06-09

0 2019-06-09

0 2019-06-09

1 2019-04-28

1 2019-04-28

1 2019-04-28

1 2019-04-28

1 2019-04-28

2 2019-05-18

2 2019-05-18

2 2019-05-18

2 2019-05-18

```
In [745]: 1 new=pd.DataFrame([])
2         for i in range(data1.shape[0]):
3             print(i)
4             each_row=data1.iloc[i,:].to_frame().transpose()
5             interval=list(each_row.iloc[0,2:7])
6             nkm=data1.loc[[i,i,i,i,i]]
7             nkm['result']=interval
8             nkm=nkm.drop(['AM Mean Travel Time','PM Mean Travel Time','Midday
9             nkm['interval']=[0,1,2,3,4]
10            nkm=nkm.reset_index(drop=True)
11            new=new.append(nkm)
12
13 new.head()
14 # for i in range(data1.shape[0]):
15
```

```
0
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
```

In [746]:

```

1 new=new.reset_index(drop=True)
2 new

```

0	3601 Lyon St, San Fra...	Francisco, CA	6	7	1588.000000	0
1	The Palace Of Fine Arts, 3601 Lyon St, San Fra...	Embarcadero, San Francisco, CA	6	7	1252.000000	1
2	The Palace Of Fine Arts, 3601 Lyon St, San Fra...	Embarcadero, San Francisco, CA	6	7	1398.500000	2
3	The Palace Of Fine Arts, 3601 Lyon St, San Fra...	Embarcadero, San Francisco, CA	6	7	1149.500000	3
4	The Palace Of Fine Arts, 3601 Lyon St, San Fra...	Embarcadero, San Francisco, CA	6	7	1228.000000	4
5	Embarcadero, San Francisco, CA	Fisherman's Wharf, 286-298 Jefferson St, San F...	6	5	579.000000	0
6	Embarcadero, San Francisco, CA	Fisherman's Wharf, 286-298 Jefferson St, San F...	6	5	664.000000	1
7	Embarcadero, San Francisco, CA	Fisherman's Wharf, 286-298 Jefferson St, San F...	6	5	654.000000	2
8	Embarcadero, San Francisco, CA	Fisherman's Wharf, 286-298 Jefferson St, San F...	6	5	608.000000	3

In [747]:

```

1 import math
2 from sklearn.ensemble import RandomForestRegressor
3 from sklearn.ensemble import GradientBoostingRegressor
4
5 def train(train_df):
6     train_df = train_df.copy()
7     X_train, X_test, y_train, y_test = train_test_split(train_df.drop(
8
9     col_c3 = ['Destination Name','Origin Name','weekday','month','inte
10    c3_transformer = Pipeline(steps=[('onehot', OneHotEncoder())])
11
12    preproc = ColumnTransformer(transformers=[
13                                ('c3', c3_transformer, co
14
15    pl = Pipeline(steps=[('preprocessor', preproc), ('regressor', Deci
16    pl.fit(X_train,y_train)
17    pred = pl.predict(X_test)
18    rms = math.sqrt(mean_squared_error(y_test, pred))
19    preds = pl.predict(train_df.drop(['result'],axis=1))
20    train_df['pred'] = preds
21    return rms,train_df

```

In [753]:

```

1 # ddddd=train(new)[1]
2 # neww=neww.reset_index(drop=True)
3 # ddddd['date']=neww['Date']
4 # ddddd.to_csv('outpututtt.csv')
5 import statsmodels
6

```

In [ ]:

1



