

# PREDICTING NFL ARRESTS

JUAN SALAZAR

SAMUEL SACHNOFF

MICHAEL SILIN

# PREPARING THE DATA

BEFORE

	season	week_num	day_of_week	gametime_local	home_team	away_team	home_score	away_score	OT_flag	arrests	division_game
0	2011	1	Sunday	1:15:00 PM	Arizona	Carolina	28	21	1	5	n
1	2011	4	Sunday	1:05:00 PM	Arizona	New York Giants	27	31	1	6	n
2	2011	7	Sunday	1:05:00 PM	Arizona	Pittsburgh	20	32	1	9	n
3	2011	9	Sunday	2:15:00 PM	Arizona	St. Louis	19	13	0	6	y
4	2011	13	Sunday	2:15:00 PM	Arizona	Dallas	19	13	0	3	n
5	2011	14	Sunday	2:05:00 PM	Arizona	San Francisco	21	19	1	4	y
6	2011	15	Sunday	2:15:00 PM	Arizona	Cleveland	20	17	0	1	n
7	2011	17	Sunday	2:15:00 PM	Arizona	Seattle	23	20	0	4	y
8	2012	1	Sunday	1:25:00 PM	Arizona	Seattle	20	16	1	0	y
9	2012	3	Sunday	1:05:00 PM	Arizona	Philadelphia	27	6	1	12	n
10	2012	4	Sunday	1:05:00 PM	Arizona	Miami	24	21	0	4	n



# AFTER

	season	week_num	day_of_week	gametime_local	home_team	away_team	home_score	away_score	OT_flag	arrests	division_game
0	2011	1	0	0	0	4	28	21	1	5	0
1	2011	4	0	1	0	20	27	31	1	6	0
2	2011	7	0	1	0	24	20	32	1	9	0
3	2011	9	0	2	0	28	19	13	1	6	1
4	2011	13	0	2	0	8	19	13	1	3	0
5	2011	14	0	3	0	26	21	19	1	4	1
6	2011	15	0	2	0	7	20	17	1	1	0
7	2011	17	0	2	0	27	23	20	1	4	1
8	2012	1	0	4	0	27	20	16	1	0	1
9	2012	3	0	1	0	23	27	6	1	12	0
10	2012	4	0	1	0	16	24	21	1	4	0
11	2012	6	0	1	0	3	19	16	1	1	0
12	2012	8	1	5	0	26	3	24	1	3	1

- MOST OF THE GAMES WERE PLAYED ON SUNDAY
- WE USED CUSTOM FUNCTIONS AND `DATAFRAME.TRANSFORM()` TO CONVERT THE DATA TO VECTORS
- RANDOMLY SHUFFLE THE INDICES BY USING `NP.RANDOM.SHUFFLE()`
- SPLIT THE SHUFFLED DATA FRAME INTO 80% TRAINING AND 20% TESTING SETS



- NORMALIZED TRAINING AND TESTING DATA AND MULTIPLY NEGATIVE CORRELATIONS BY -1 WHETHER IN CROSS VALIDATION OR IN ACTUAL MODEL PREDICTIONS
- WE USED A 10-FOLD CROSS VALIDATION TO CREATE OUR MODEL
- WE THEN TRAINED OUR MODEL AND PREDICTED THE TESTING DATA
- USED A POLYNOMIAL DEGREE OF 5
- WE THEN RAN OUR PROGRAM 10 TIMES TO DETERMINE HOW ACCURATE OUR MODEL WAS
- AVG ERROR: +/- 3.54 ARRESTS

```

def crossValidate(df1):
    errors= []
    #get the different x columns to take correlations of
    x_columns= np.array([])
    x_columns= np.append(df1.columns[0],x_columns)
    x_columns= np.append(df1.columns[2:],x_columns)
    y_data= df1['arrests']
    x_data= df1[x_columns]
    kf = KFold(n_splits=10,shuffle= True)
    sampleIndices= [5,12]
    model= make_pipeline(PolynomialFeatures(5),linear_model.LinearRegression())
    for train_index, test_index in kf.split(x_data):
        train_x= x_data.iloc[train_index]
        test_x= x_data.iloc[test_index]
        train_y= y_data.iloc[train_index]
        test_y= y_data.iloc[test_index]
        #normalize the data
        for col in x_columns:
            if train_x[col].corr(train_y)<0:
                train_x[col]=train_x.loc[:,col]*-1
                test_x[col]= test_x.loc[:,col]*-1
        train_x= (train_x-train_x.mean())/train_x.std()
        test_x= (test_x-test_x.mean())/test_x.std()
        model.fit(train_x,train_y)
        #get the error
        y_fit= model.predict(test_x)
        error= (test_y-y_fit).mean()
        errors.append(error)
    errors= np.array(errors)
    print(errors.mean())

    #Plot graph of errors
    erroSeries = pd.Series(errors, index=np.arange(1,len(errors)+1))
    for i in range(len(erroSeries)):
        if(erroSeries.iloc[i] < 0):
            erroSeries.iloc[i] *= -1
    erroSeries.plot()

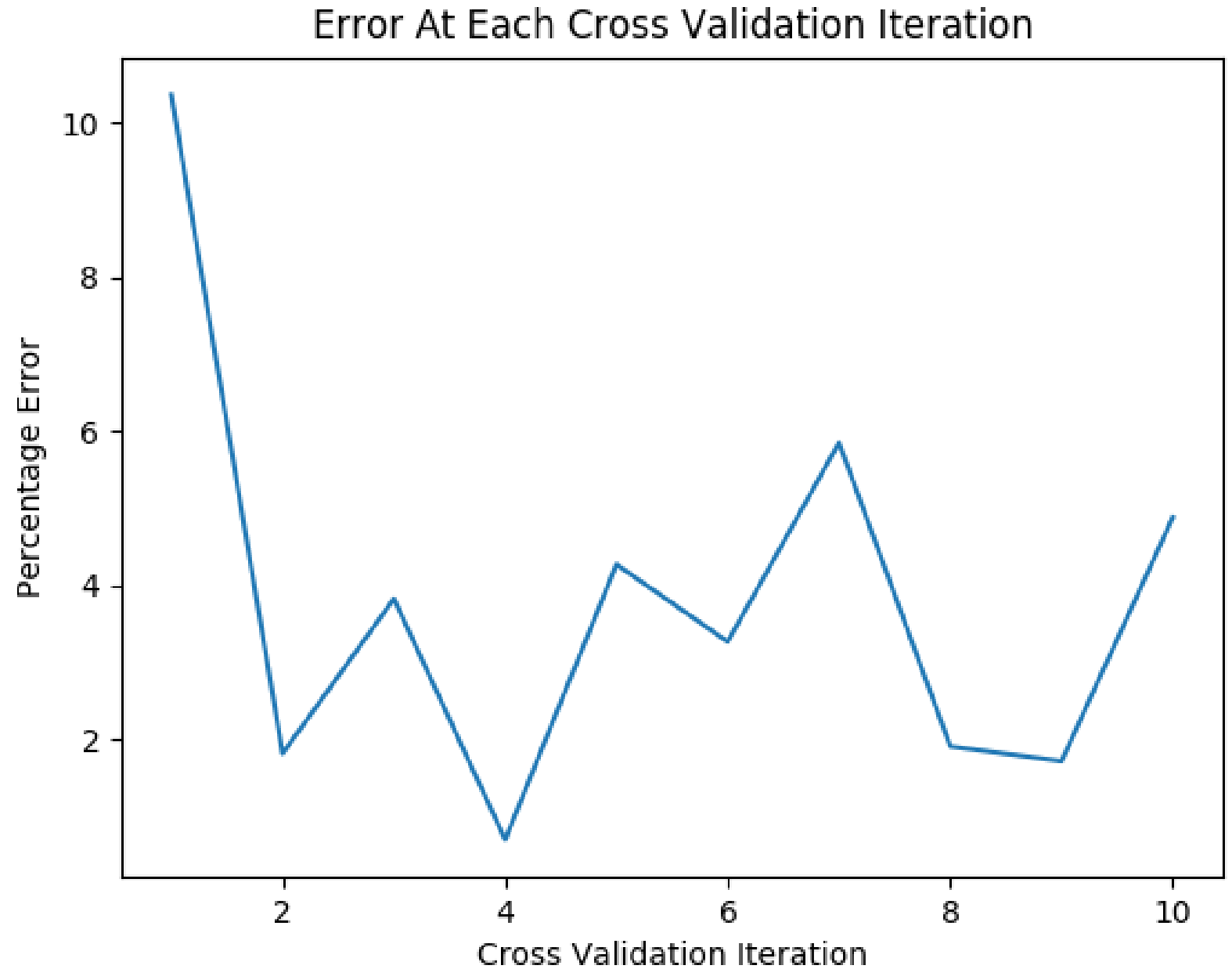
    import matplotlib.pyplot as plt
    plt.ylabel("Percentage Error")
    plt.xlabel("Cross Validation Iteration")
    plt.title("Error At Each Cross Validation Iteration")

    return model

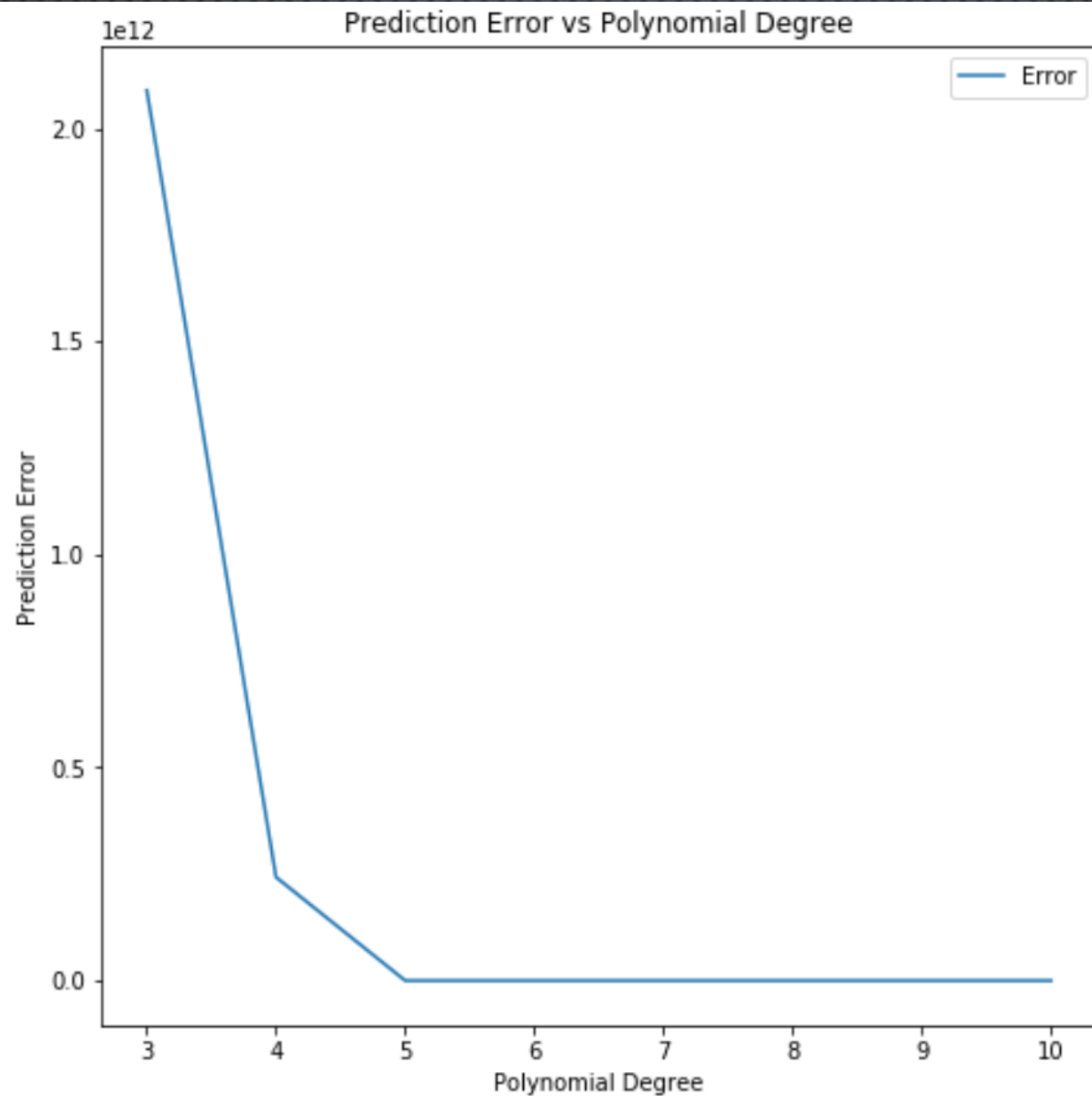
```



# CROSS VALIDATION ERROR

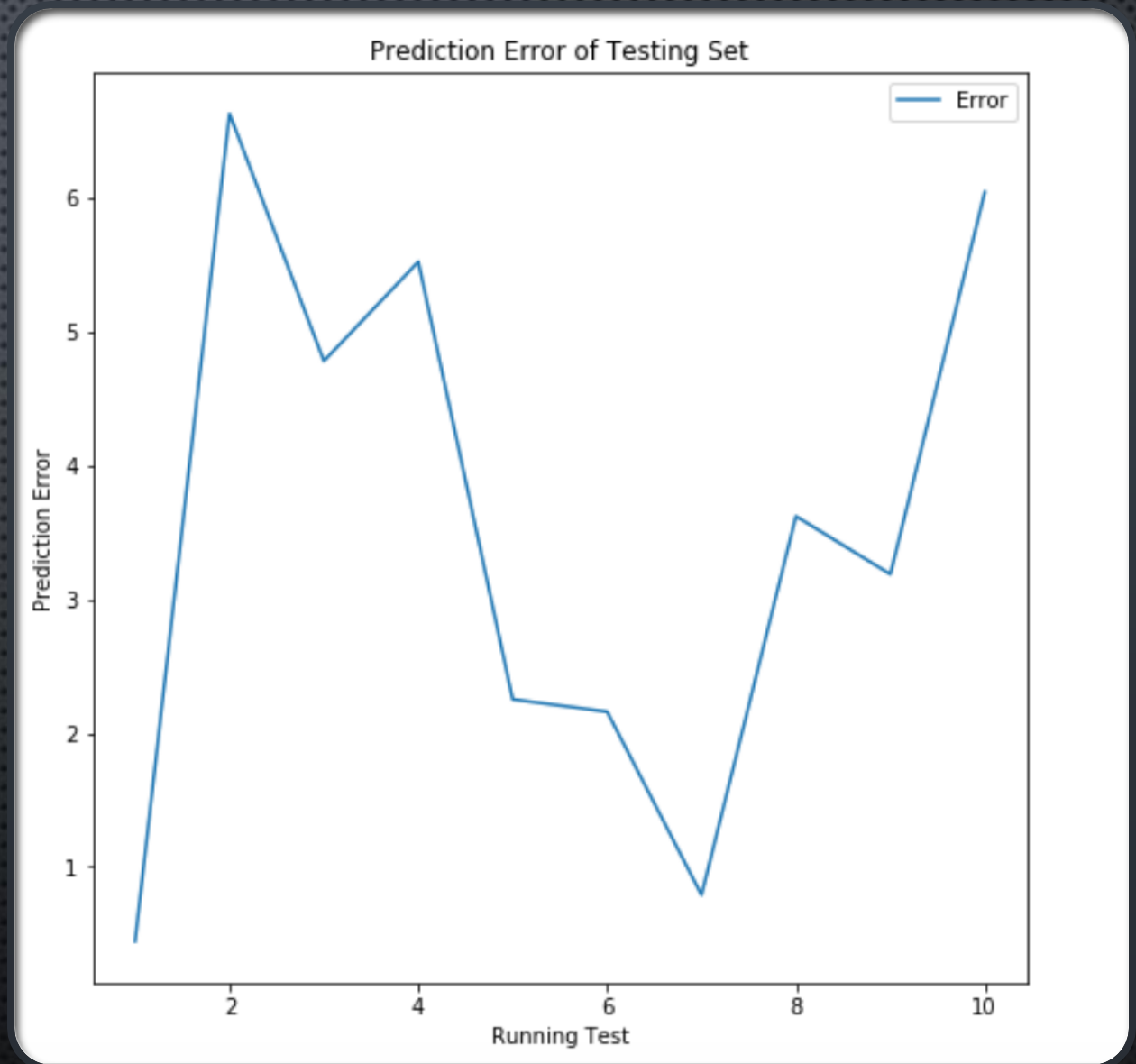


# ERRORS OF DIFFERENT POLY DEGREES

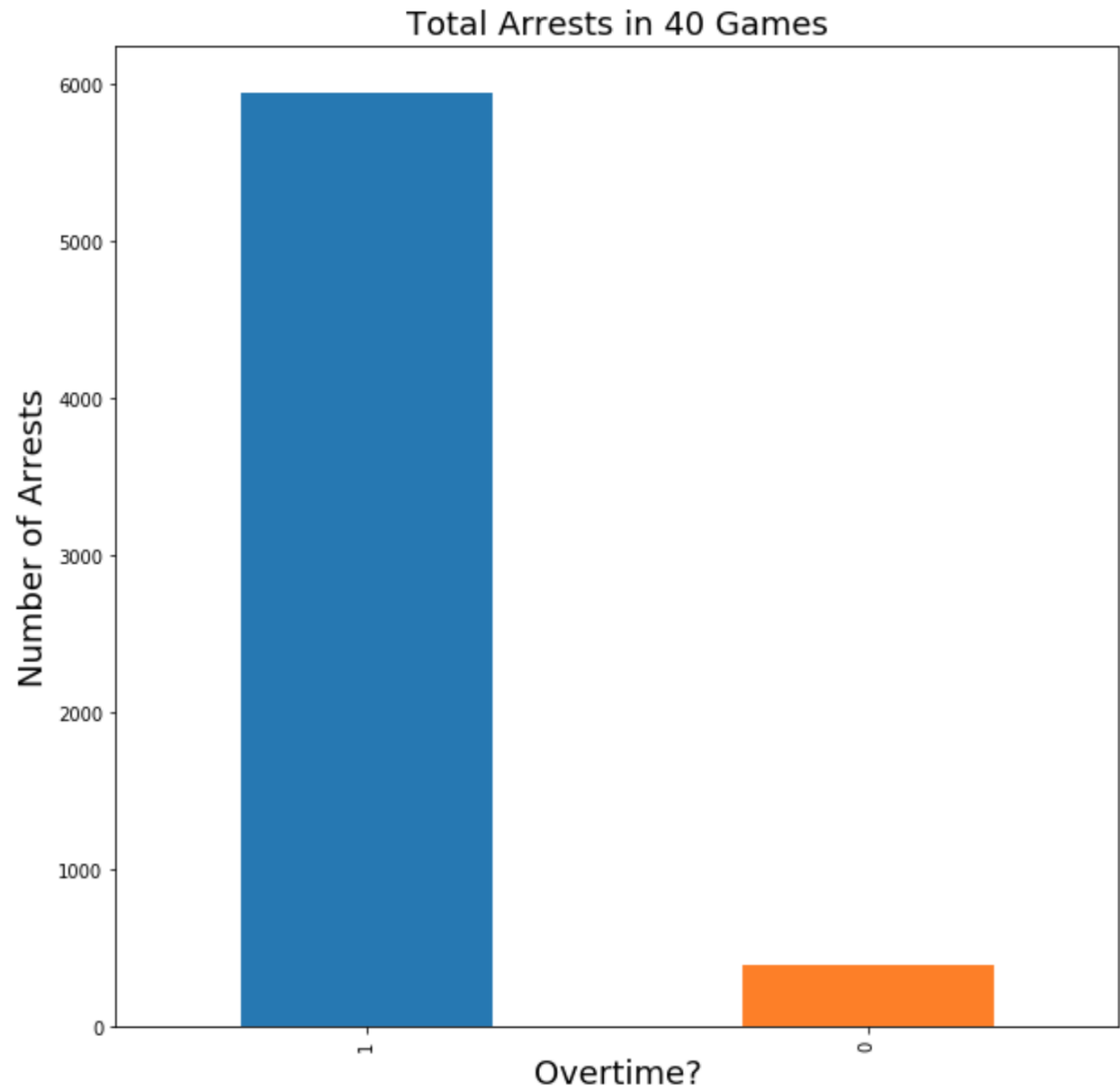




#PREDICTION  
ERRORS AFTER  
RUNNING THE  
PROGRAM 10  
DIFFERENT  
TIMES

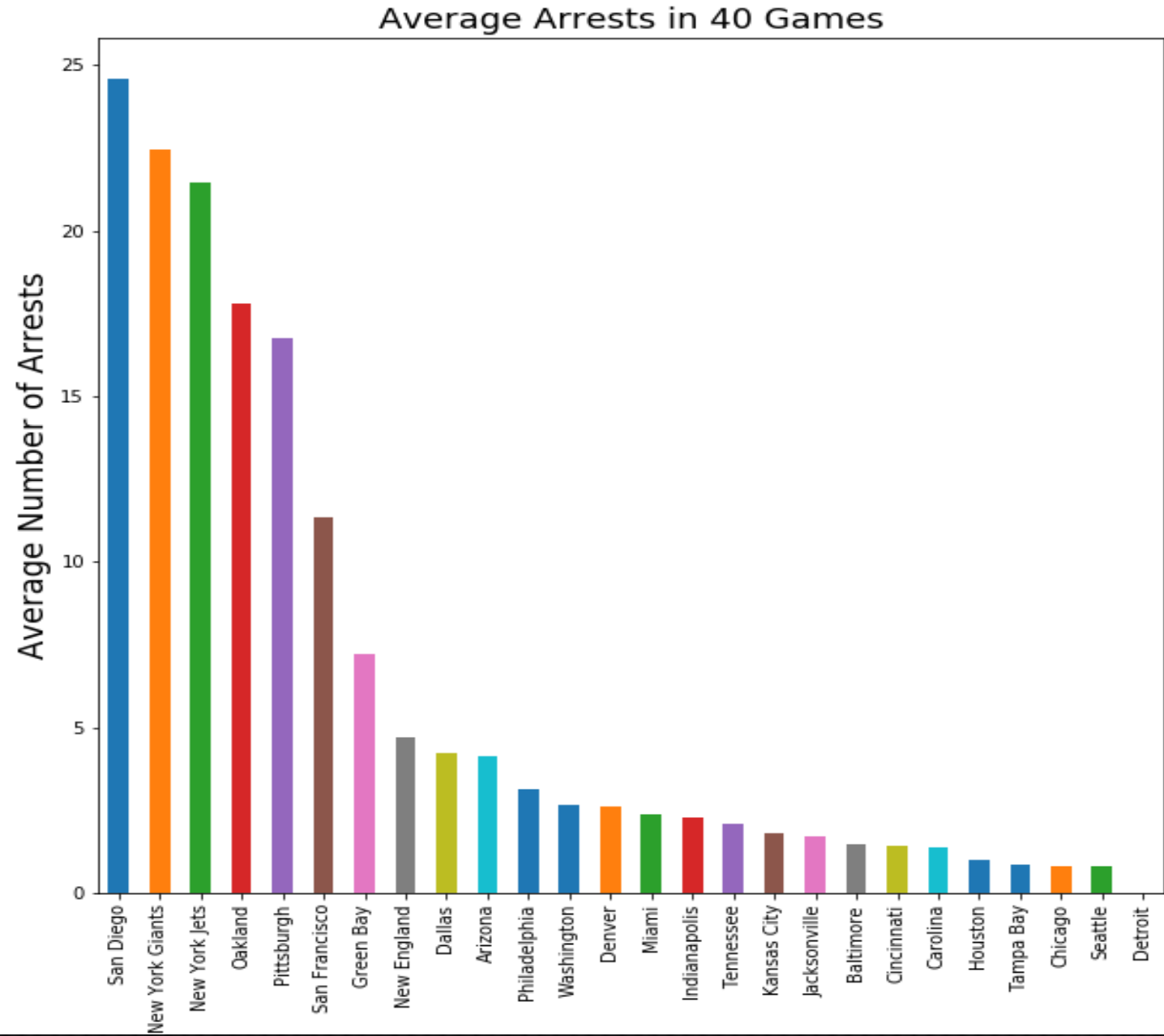


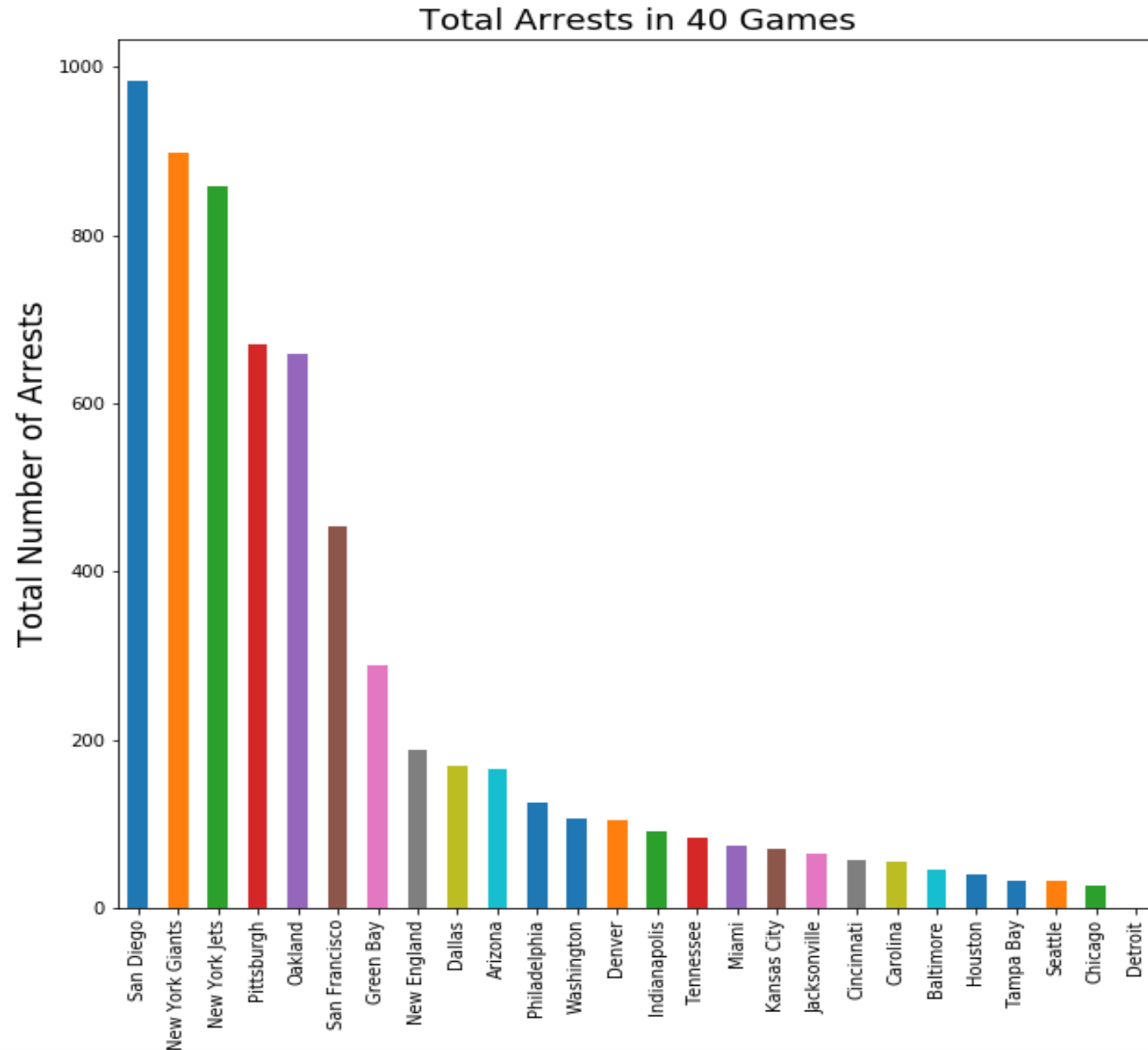
Total Arrests grouped  
by Overtime flag





## Average Arrests grouped By Team

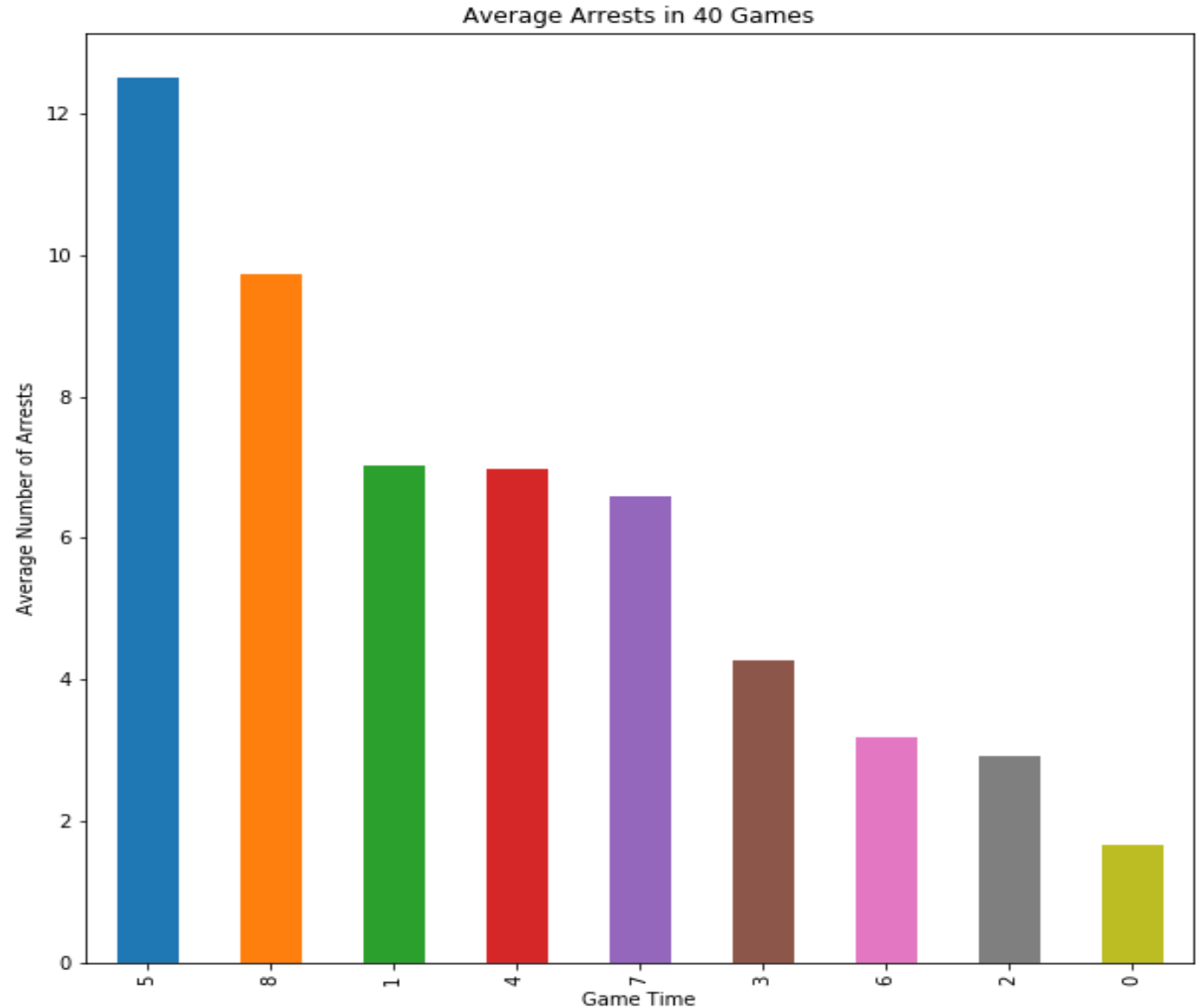




TOTAL  
ARRESTS  
GROUPED BY  
TEAM



# Average Arrests grouped by Game Time



# OBSTACLES

- SKLEARN K-FOLD CROSS VALIDATION DID NOT AUTOMATICALLY NORMALIZE AND MULT NEGATIVE CORRELATIONS BY  $-1$
- WAS GETTING CHAINED ASSIGNMENT WARNING MESSAGE WHEN DOING CROSS VALIDATION WHEN MULT NEGATIVE CORRELATIONS BY  $-1$



# WHAT WE LEARNED!!!

- CROSS VALIDATION
- USING SKLEARN FOR LINEAR REGRESSION
- GRAPHING DATA
- FULL PICTURE OF CLASS
- SAN DIEGO HAS THE HIGHEST NUMBER OF ARRESTS (ALMOST 1000 IN 40 GAMES)

