

Exploratory Flux Predictions

Nicholas Kaufman

May 3, 2015

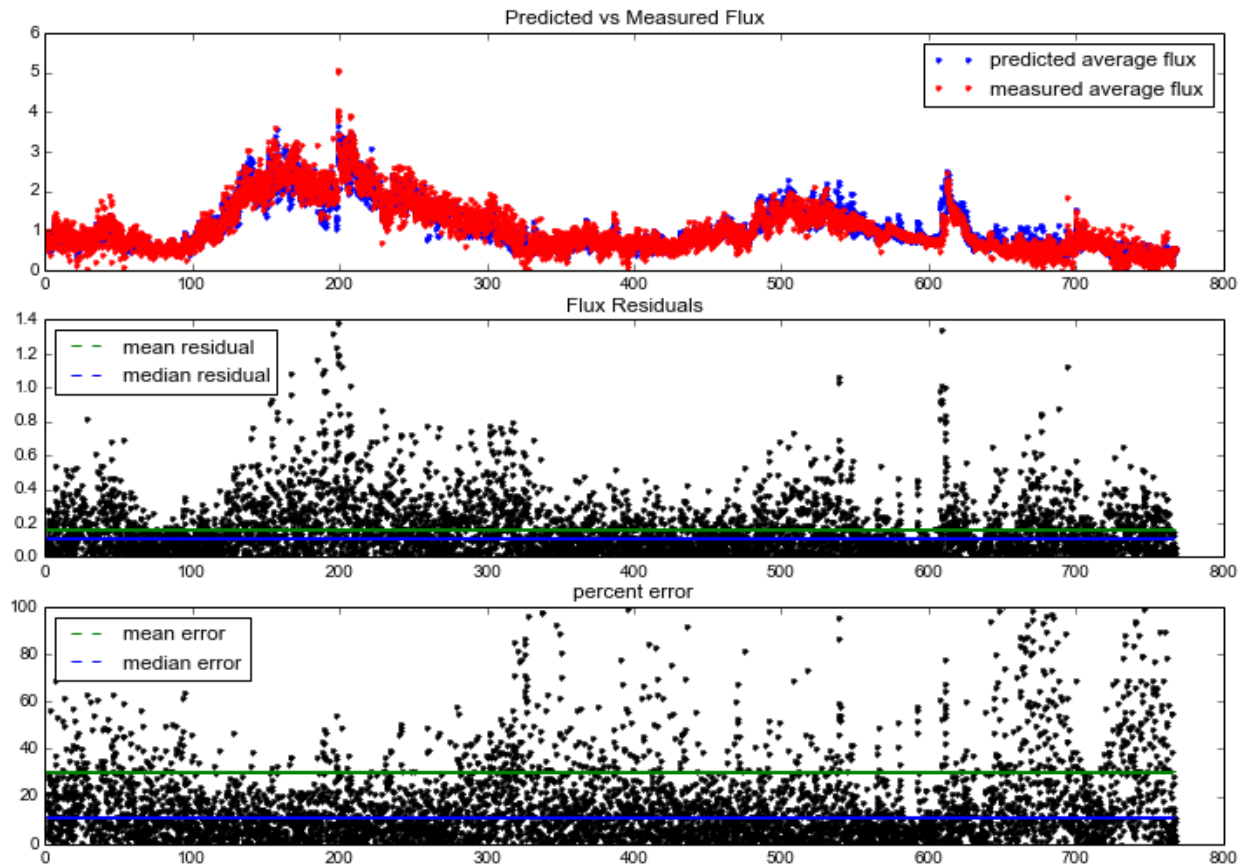
This report will detail the results of predicting average flux emissions using three different models.

Random Forests

The First Model

For the first model, time was eliminated as a training variable. All other factor variables - Air Temperature, three layers of Soil Temperature, Soil Moisture, and Light - were used.

A plot is included below showing some preliminary results.



While training, the R2 score (or the coefficient of determination) is calculated on the training data, and on a cross-validation set. Those values are given below:

R2 score on training data: 0.982474990177

R2 score on 5-folded data: [0.87108384 0.87205808 0.86948272 0.86920592 0.85897424]

Average score across folds: 0.868160959768

R2 obtains a maximum at 1, and is a measure of how well the model is performing on a particular data set, for which we can compare a known baseline.

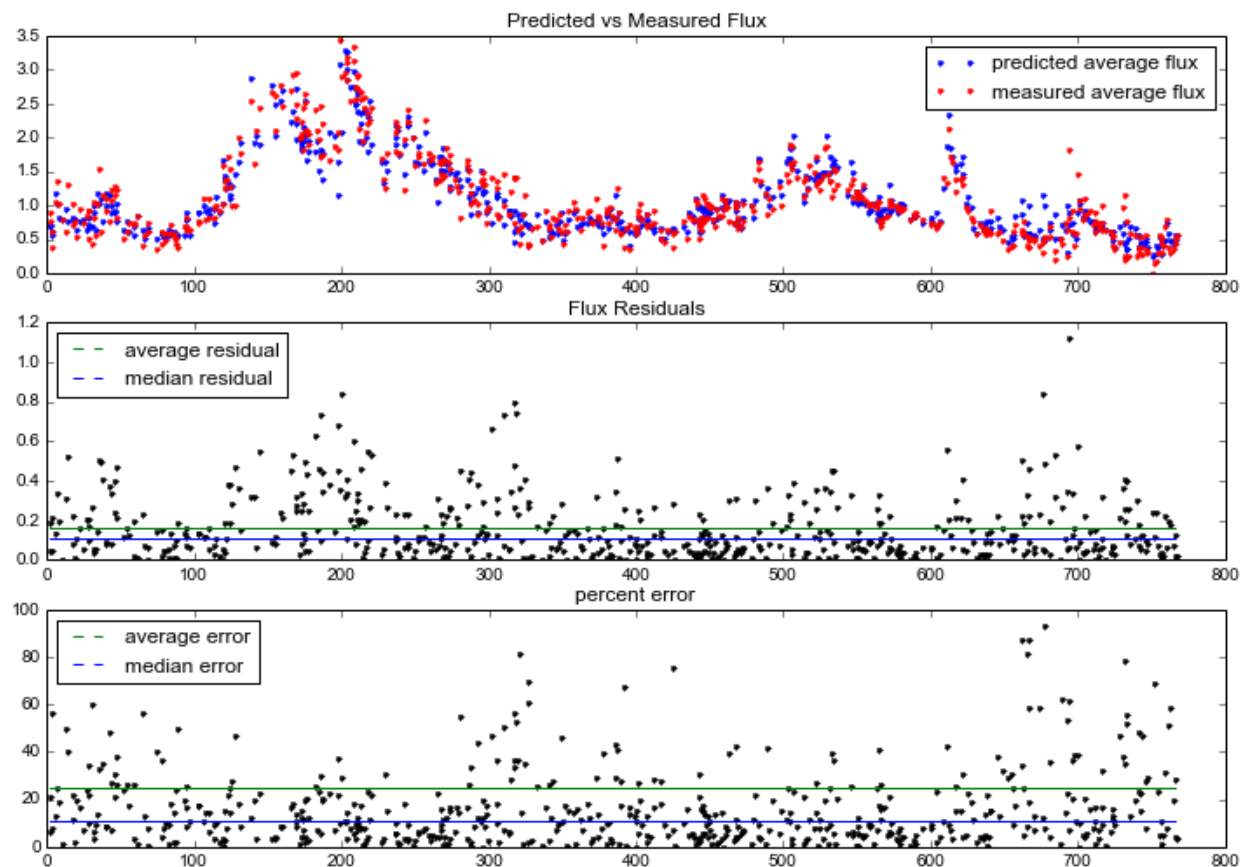
One can also compute the weights the model gives to each of the feature variables. In this case, those are:

Feature Importance Vector: [0.04442861 0.59248679 0.0466869 0.03521874 0.25649003 0.02468893]

The higher the number, the more important the weight. From this we see that, by a large margin, the deepest soil temperature and light are the most important features.

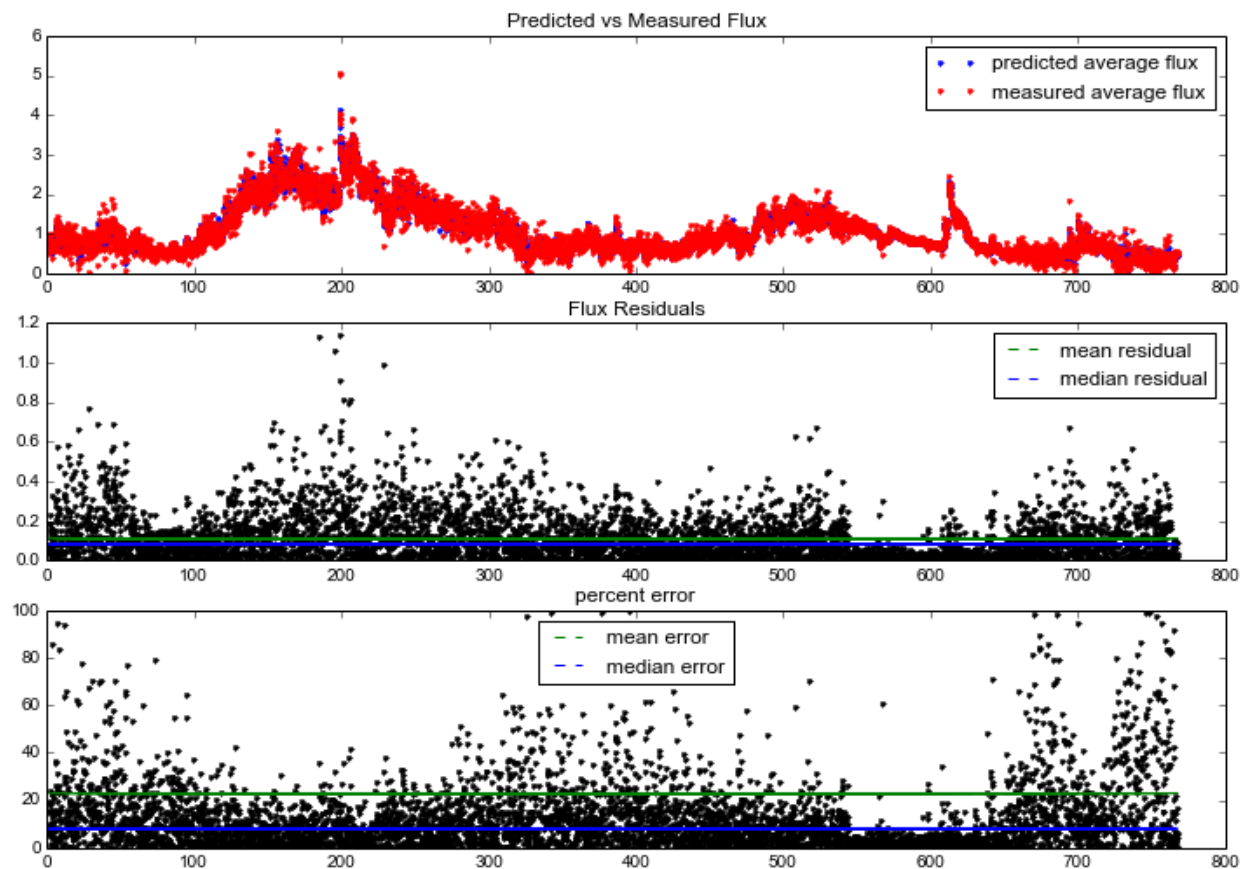
It is interesting to note that, in the predictions plot above, we can see two time distinct time periods in which the models performance is weakest. It is also worthy of mention that since this data is over the span of two years, those areas appear to correspond to roughly the same yearly time range. We can conclude that perhaps it would be useful to include the time as a variable.

In order to present a perhaps cleaner picture, we include the same plots but using only every tenth data point.



The Second Model

Here we include time as a variable, according to the discussion following the previous plot. When training the model with this additional factor variable, the predictions improve.



We notice that the predictions are significantly more accurate here. We include the calculated output as before.

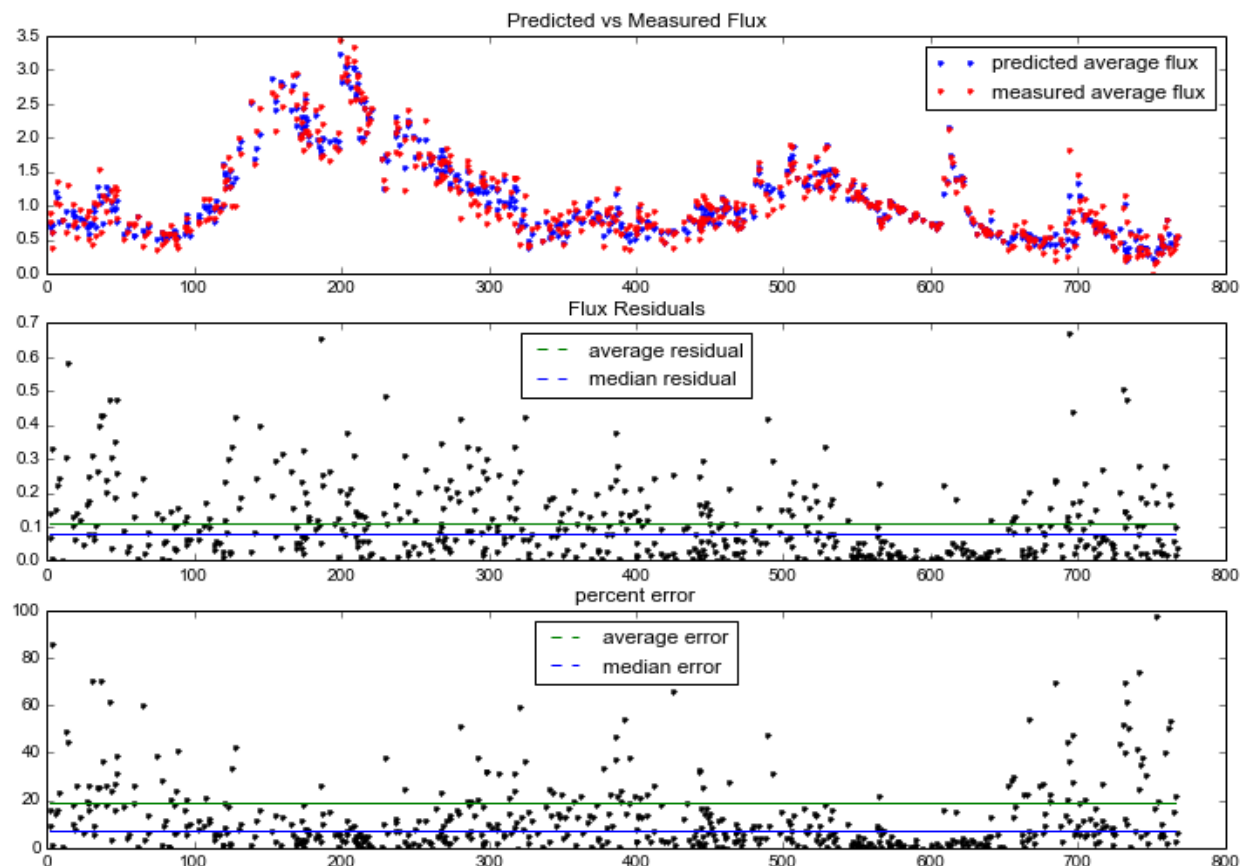
R2 score on training data: 0.990912535056 R2 score on 5-folded data: [0.9345659 0.93356225 0.93235294 0.9353663 0.92977682] Average score across folds: 0.933124842167

R2 score on testing data:

Feature Importance Vector: [0.02230462 0.50997696 0.01256326 0.01180546 0.05988339 0.00664382 0.37682248]

From this, we see that light is now no longer an important feature. Time as taken it's place by a large margin as the penultimate feature.

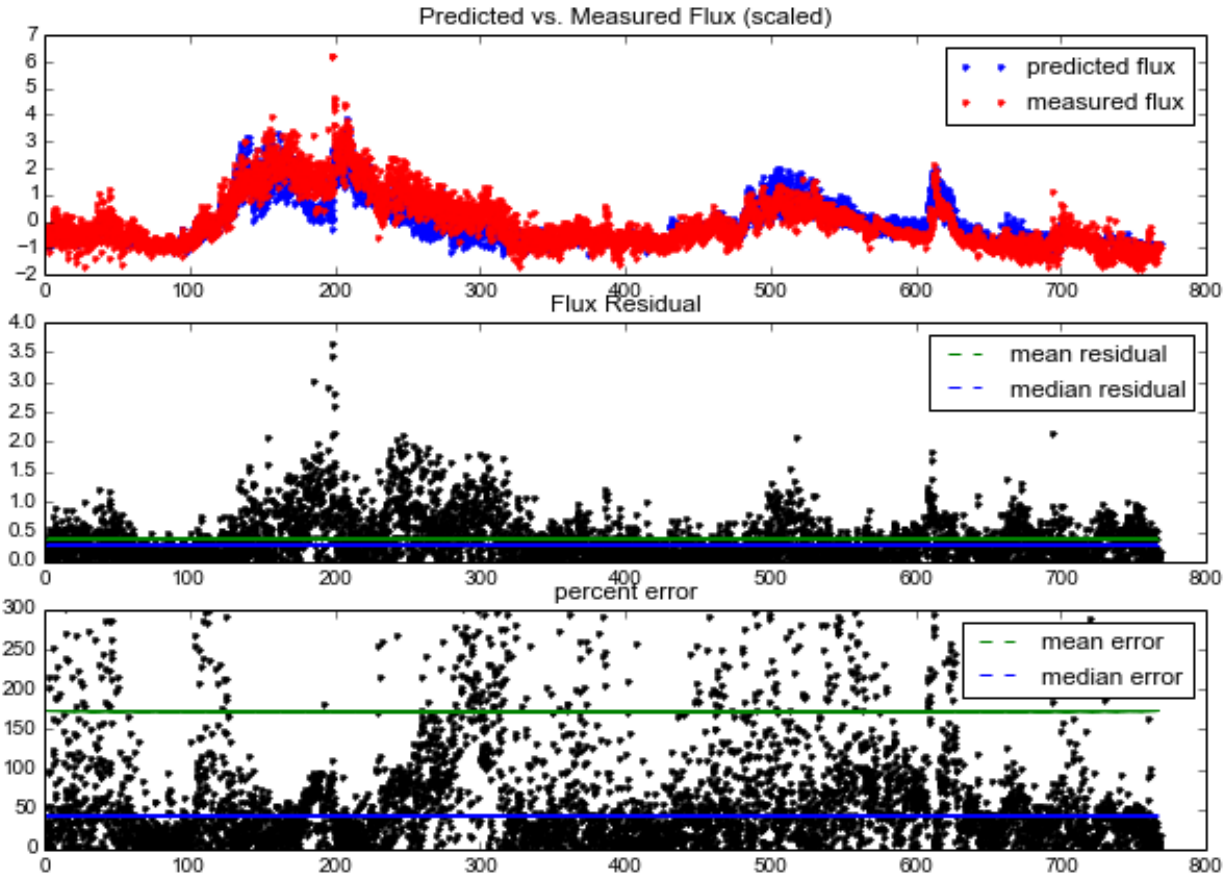
Again, we plot the results using a tenth of the data for a neater picture.



1 Support Vector Machine

1.1 The Third Model

This model's performance is quite inferior to the previous models. The results are included because they highlight the importance of time as a variable. Note that it is impossible to use time as a variable in SVM models because the model requires the data to be scaled to fit the Standard Distribution, per feature. This is obviously infeasible to do with time.



Here we can very clearly see the two regions in time in which the prediction most notably breaks down. I believe that given more feature variables, we can highlight which attributes are most significant in predicting the average flux. Further, I believe that it would be possible to use similar methodologies to predict the amount of flux for varying gasses based on a larger range of dependent variables.

2 Source Code

```

1  #Flux prediction
2  #data courtesy of Dr. Rodrigo Vargas
3  #Last modified: 5/2/15
4
5  import numpy as np
6  import matplotlib.pyplot as plt
7  import pandas as pd
8  from sklearn.cross_validation import KFold, cross_val_score, train_test_split
9  from sklearn import ensemble
10 from sklearn.svm import SVR
11
12 #read in the data
13 raw_vars = pd.read_csv('vargas_raw_data.csv')
14 raw_data = np.array(raw_vars.values)
15
16 bad_ind = []
17 for i in enumerate(raw_vars['average_flux']):

```

```

18     if (np.isnan(i[1]) or np.isnan(raw_vars['air_temp_avg'][i[0]]) \
19     or np.isnan(raw_vars['deep_soil_temp'][i[0]]) \
20     or np.isnan(raw_vars['mid_soil_temp'][i[0]]) \
21     or np.isnan(raw_vars['shallow_soil_temp'][i[0]]) \
22     or np.isnan(raw_vars['soil_moisture'][i[0]]) \
23     or np.isnan(raw_vars['light'][i[0]]) or np.isnan(raw_vars['time'][i[0]])) :
24         bad_ind.append(i[0])
25
26 raw_data = np.delete(raw_data, bad_ind, axis=0)
27
28 #we split the data/target into two different types:
29 #type 1: time is not a feature.
30 #need to carry time until you can after splitting for training/testing
31 #type 2: time is a feature.
32
33 data_type1 = raw_data[:,1:7]
34 target_type1 = raw_data[:,[0,7]]
35 data_type2 = raw_data[:,1:]
36 target_type2 = raw_data[:,0]
37
38 #split data up into training/testing for each type
39 X_t1_train, X_t1_test, y_t1_train, y_t1_test = train_test_split(data_type1,\
40     target_type1, test_size=.3, random_state=42)
41
42 X_t2_train, X_t2_test, y_t2_train, y_t2_test = train_test_split(data_type2,\
43     target_type2, test_size=.3, random_state=42)
44 #pull out the time vectors for the type 1 data, and excise it from target
45 time_t1_train = y_t1_train[:,1]
46 y_t1_train = y_t1_train[:,0]
47 time_t1_test = y_t1_test[:,1]
48 y_t1_test = y_t1_test[:,0]
49
50 #I don't know python coding procedure. where do you put all of your functions?
51
52 def scale_data(train_data, test_data, train_target, test_target):
53     from sklearn import preprocessing
54     scaler_X = preprocessing.StandardScaler().fit(train_data)
55     scaler_y = preprocessing.StandardScaler().fit(train_target)
56     train_data = scaler_X.transform(train_data)
57     test_data = scaler_X.transform(test_data)
58     train_target = scaler_y.transform(train_target)
59     test_target = scaler_y.transform(test_target)
60
61     return train_data, test_data, train_target, test_target
62
63 X_svm_train, X_svm_test, y_svm_train, y_svm_test = scale_data \
64     (X_t1_train, X_t1_test, y_t1_train, y_t1_test)
65
66 def train_and_evaluate(clf, X, y):
67     clf.fit(X,y)
68     print "R2_score_on_training_data:", clf.score(X,y)
69
70     kf = KFold(X.shape[0], n_folds = 5, shuffle = True, random_state = 42)
71     score = cross_val_score(clf, X,y,cv=kf)

```

```

72     print "R2_score_on_5-folded_data:", score
73     print "Average_score_across_folds:", np.mean(score)
74
75 rf_t1 = ensemble.RandomForestRegressor(n_estimators=75, random_state=42)
76 rf_t2 = ensemble.RandomForestRegressor(n_estimators=75, random_state=33)
77 svm_t1 = SVR()
78
79 print "training_random_forest_1_model..."
80 train_and_evaluate(rf_t1, X_t1_train, y_t1_train)
81
82 print "training_random_forest_2_model..."
83 train_and_evaluate(rf_t2, X_t2_train, y_t2_train)
84
85 print "training_svm..."
86 train_and_evaluate(svm_t1, X_svm_train, y_svm_train)
87
88 #predict flux values using each model
89 rf_t1_pred = rf_t1.predict(X_t1_test)
90 rf_t2_pred = rf_t2.predict(X_t2_test)
91 svm_pred = svm_t1.predict(X_svm_test)
92
93 #compute plotting vectors, per model
94 rf_t1_residual = np.array(rf_t1_pred - y_t1_test)
95 rf_t1_error = np.array(np.abs(rf_t1_residual) / y_t1_test)
96 rf_t2_residual = np.array(rf_t2_pred - y_t2_test)
97 rf_t2_error = np.array(np.abs(rf_t2_residual) / y_t2_test)
98 svm_residual = np.array(svm_pred - y_svm_test)
99 svm_error = np.array(np.abs(svm_residual) / y_svm_test))
100
101 #Now we generate plots.
102 f, ax = plt.subplots(figsize=(10,7), nrows=3)
103
104 ax[0].plot(time_t1_test, rf_t1_pred, 'b.', label='predicted_flux')
105 ax[0].plot(time_t1_test, y_t1_test, 'r.', label='measured_flux')
106 ax[0].set_title('Predicted_vs_Measure_Flux')
107 ax[0].legend(loc='best')
108 ax[1].plot(time_t1_test, np.abs(rf_t1_residual), 'k.')
109 ax[1].plot(time_t1_test, np.ones(len(time_t1_test)) \
110           *np.mean(np.abs(rf_t1_residual)), 'g—', label='average_residual')
111 ax[1].plot(time_t1_test, np.ones(len(time_t1_test)) \
112           *np.median(np.abs(rf_t1_residual)), 'b—', label='median_residual')
113 ax[1].set_title('Flux_Residual_(abs_diff_of_predicted_and_measured)')
114 ax[1].legend(loc='best')
115 ax[2].plot(time_t1_test, rf_t1_error*100, 'k.')
116 ax[2].plot(time_t1_test, np.ones(len(time_t1_test))*np.mean(rf_t1_error*100)\
117           , 'g—', label="average_error")
118 ax[2].plot(time_t1_test, np.ones(len(time_t1_test))*np.median(rf_t1_error*100)\
119           , 'b—', label="median_error")
120 ax[2].set_title('Percent_Error_in_flux_prediction_values')
121 ax[2].legend(loc='best')
122
123 f.show()
124
125 g, ay = plt.subplots(figsize=(10,7), nrows=3)

```



```

126
127 ay[0].plot(X_t2_test[:, -1:], rf_t2_pred, 'b.', label='predicted_flux')
128 ay[0].plot(X_t2_test[:, -1:], y_t2_test, 'r.', label='measured_flux')
129 ay[0].set_title('Predicted_vs_Measure_Flux')
130 ay[0].legend(loc='best')
131 ay[1].plot(X_t2_test[:, -1:], np.abs(rf_t2_residual), 'k.')
132 ay[1].plot(X_t2_test[:, -1:], np.ones(len(X_t2_test[:, -1:])) \
133         *np.mean(np.abs(rf_t2_residual)), 'g—', label='average_residual')
134 ay[1].plot(X_t2_test[:, -1:], np.ones(len(X_t2_test[:, -1:])) \
135         *np.median(np.abs(rf_t2_residual)), 'b—', label='median_residual')
136 ay[1].set_title('Flux_Residual_(abs_diff_of_predicted_and_measured)')
137 ay[1].legend(loc='best')
138 ay[2].plot(X_t2_test[:, -1:], rf_t2_error*100, 'k.')
139 ay[2].plot(X_t2_test[:, -1:], np.ones(len(X_t2_test[:, -1:])) \
140         *np.mean(rf_t2_error*100), 'g—', label="average_error")
141 ay[2].plot(X_t2_test[:, -1:], np.ones(len(X_t2_test[:, -1:])) \
142         *np.median(rf_t2_error*100), 'b—', label="median_error")
143 ay[2].set_title('Percent_Error_in_flux_prediction_values')
144 ay[2].legend(loc='best')
145
146 g.show()
147
148 h, az = plt.subplots(figsize=(10,7), nrows=3)
149
150 ay[0].plot(time_t1_test, svm_pred, 'b.', label='predicted_flux')
151 ay[0].plot(time_t1_test, y_svm_test, 'r.', label='measured_flux')
152 ay[0].set_title('Predicted_vs_Measure_Flux')
153 ay[0].legend(loc='best')
154 ay[1].plot(time_t1_test, np.abs(svm_residual), 'k.')
155 ay[1].plot(time_t1_test, np.ones(len(time_t1_test)) \
156         *np.mean(np.abs(svm_residual)), 'g—', label='average_residual')
157 ay[1].plot(time_t1_test, np.ones(len(time_t1_test)) \
158         *np.median(np.abs(svm_residual)), 'b—', label='median_residual')
159 ay[1].set_title('Flux_Residual_(abs_diff_of_predicted_and_measured)')
160 ay[1].legend(loc='best')
161 ay[2].plot(time_t1_test, svm_error*100, 'k.')
162 ay[2].plot(time_t1_test, np.ones(len(time_t1_test))*np.mean(svm_error*100) \
163         , 'g—', label="average_error")
164 ay[2].plot(time_t1_test, np.ones(len(time_t1_test))*np.median(svm_error*100) \
165         , 'b—', label="median_error")
166 ay[2].set_title('Percent_Error_in_flux_prediction_values')
167 ay[2].legend(loc='best')
168
169 h.show()

```