



REGRESSION MACHINE LEARNING

PREDICTING FISH WEIGHT

By Dana Nicolas

(NUMERIC) DATA

WHAT'S IN MY DATASET

Vertical Length (cm)

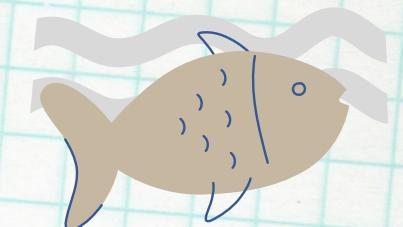
Diagonal Length (cm)

Cross Length (cm)

Width (cm)

(CATEGORICAL) DATA

WHAT'S IN MY DATASET



Species

Perch

Bream

Roach

Pike

Smelt

Parkki

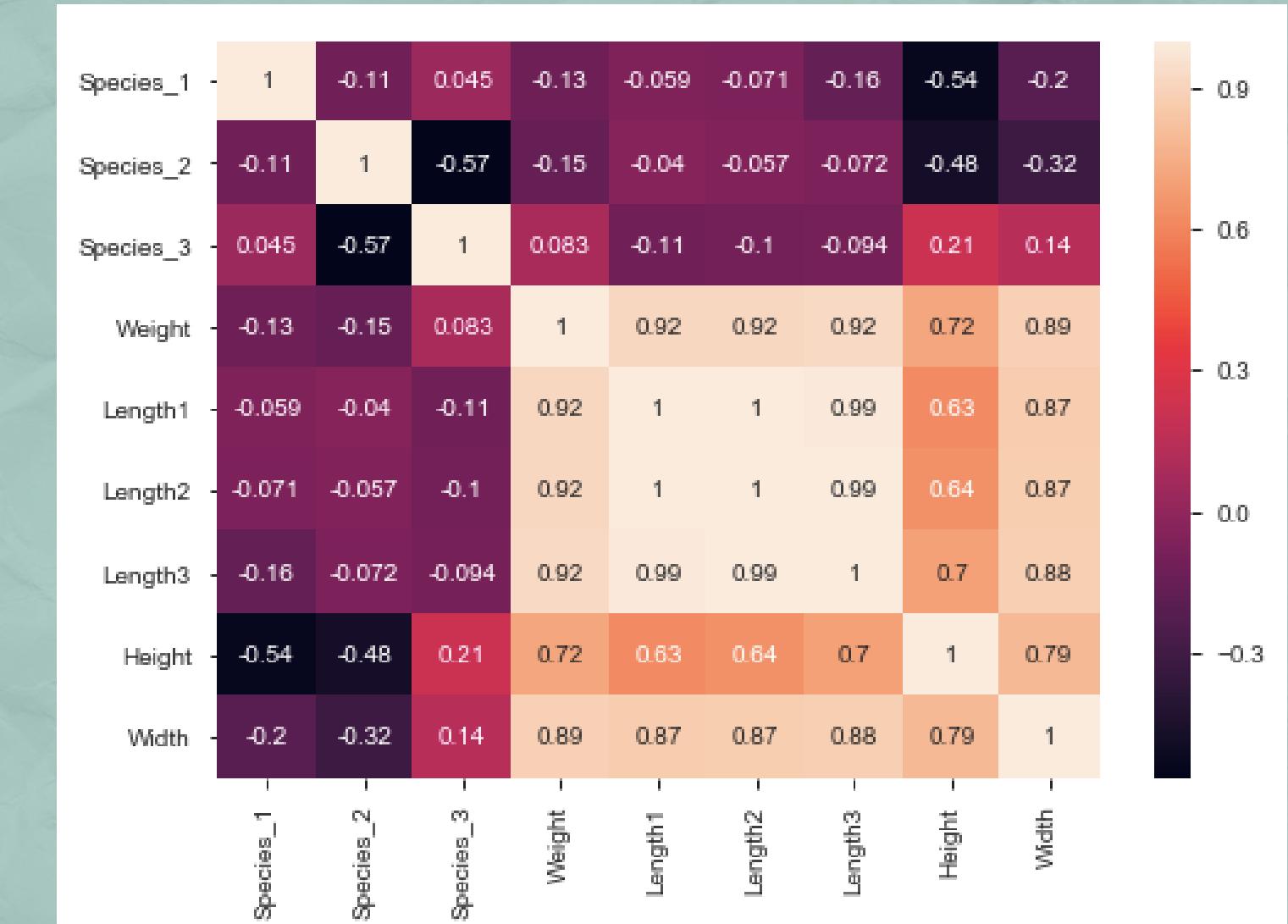
Whitefish

CONVERT CATEGORICAL TO BINARY

```
import category_encoders as ce
binary_encoder = ce.BinaryEncoder(cols=['Species'])
encoded_fishies = binary_encoder.fit_transform(fishies)
```

DATA PREPARATION (EDA)

HEATMAP



ROBUST SCALER

SCALED:
LENGTHS, HEIGHT, WEIGHT

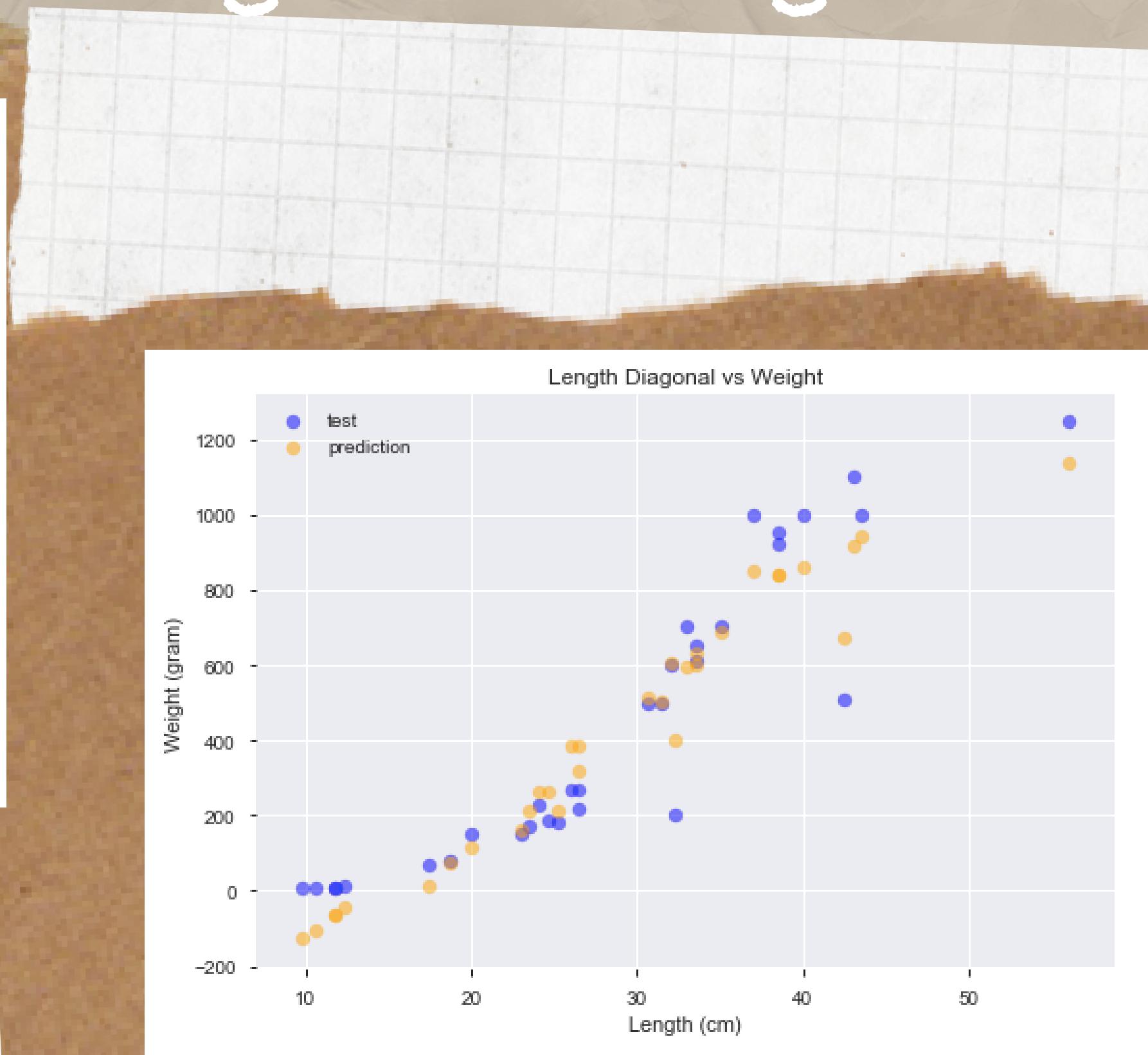
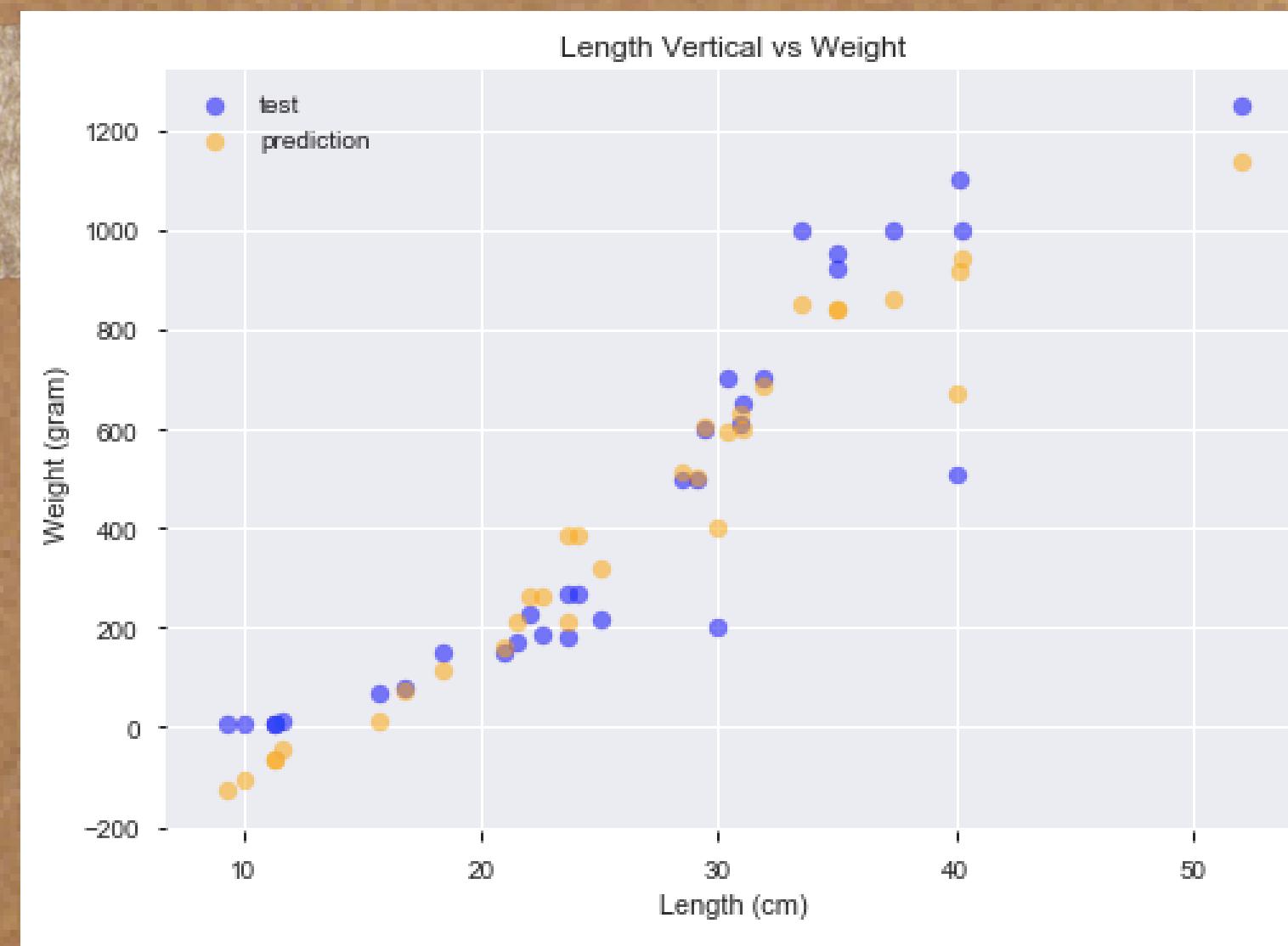
```
from sklearn.preprocessing import RobustScaler  
  
f_columns = ['LengthVer', 'LengthDia', 'LengthCro', 'Height']  
f_transformer = RobustScaler()  
weight_transformer = RobustScaler()  
  
f_transformer = f_transformer.fit(encoded_fishies[f_columns].to_numpy())  
weight_transformer = weight_transformer.fit(encoded_fishies[['Weight']])  
  
train.loc[:, f_columns] = f_transformer.transform(train[f_columns].to_numpy())  
train['Weight'] = weight_transformer.transform(train[['Weight']])  
  
test.loc[:, f_columns] = f_transformer.transform(test[f_columns].to_numpy())  
test['Weight'] = weight_transformer.transform(test[['Weight']])
```

LINEAR REGRESSION

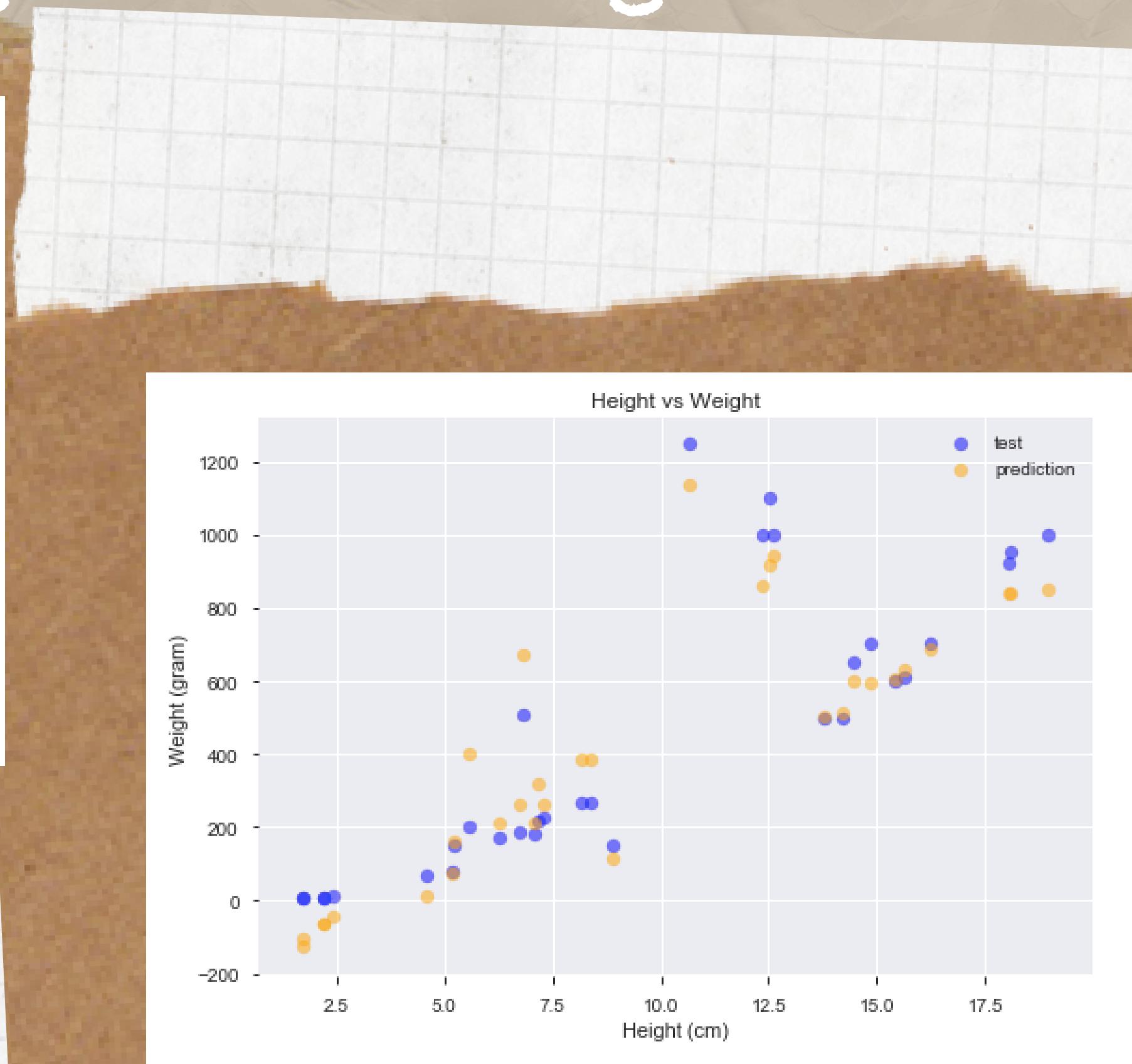
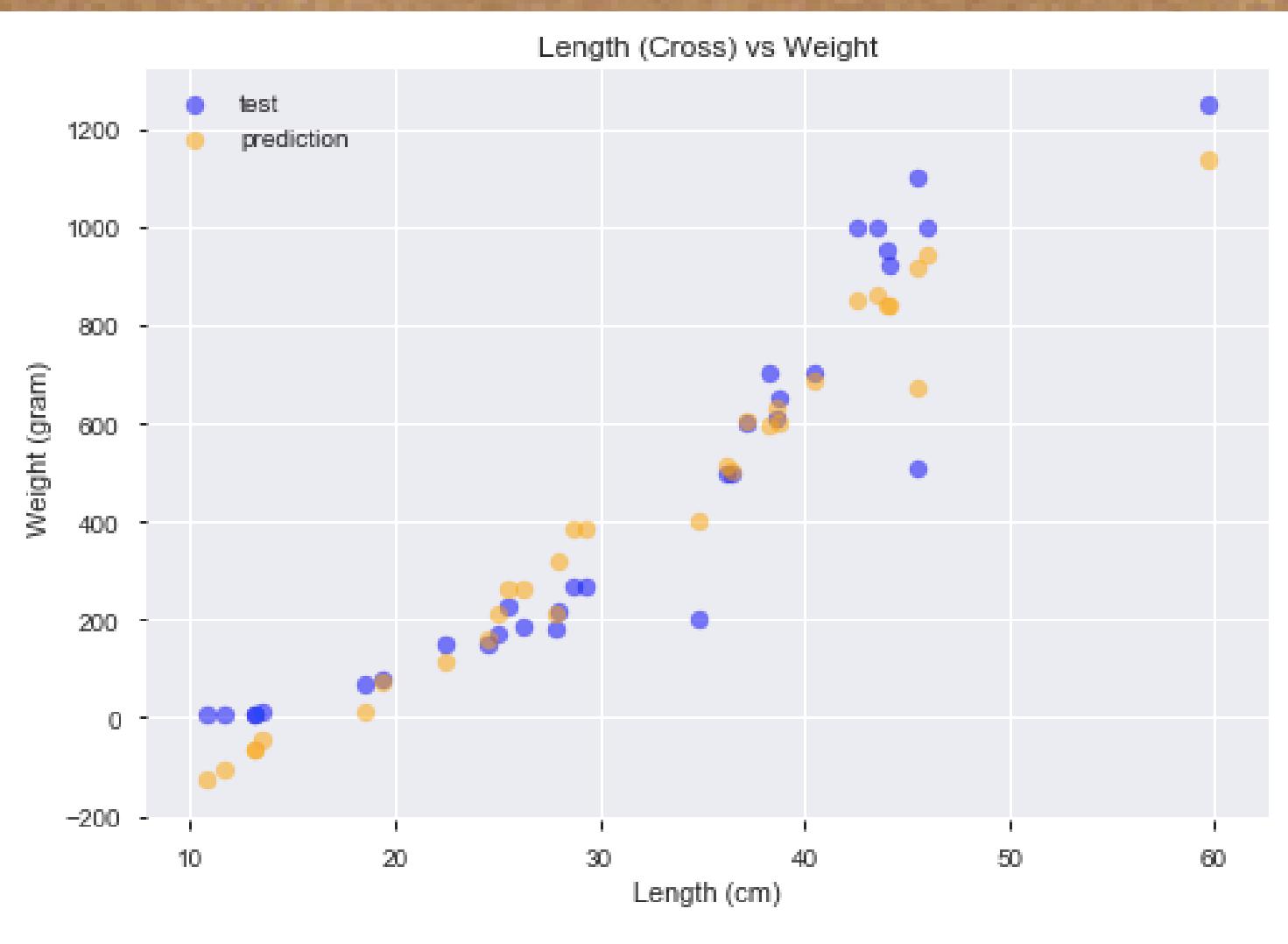
```
lr_model = LinearRegression()  
lr_model.fit(x_train, y_train)  
y_pred = lr_model.predict(x_test)  
y_pred
```

R2: 0.9373
MSE: 0.032

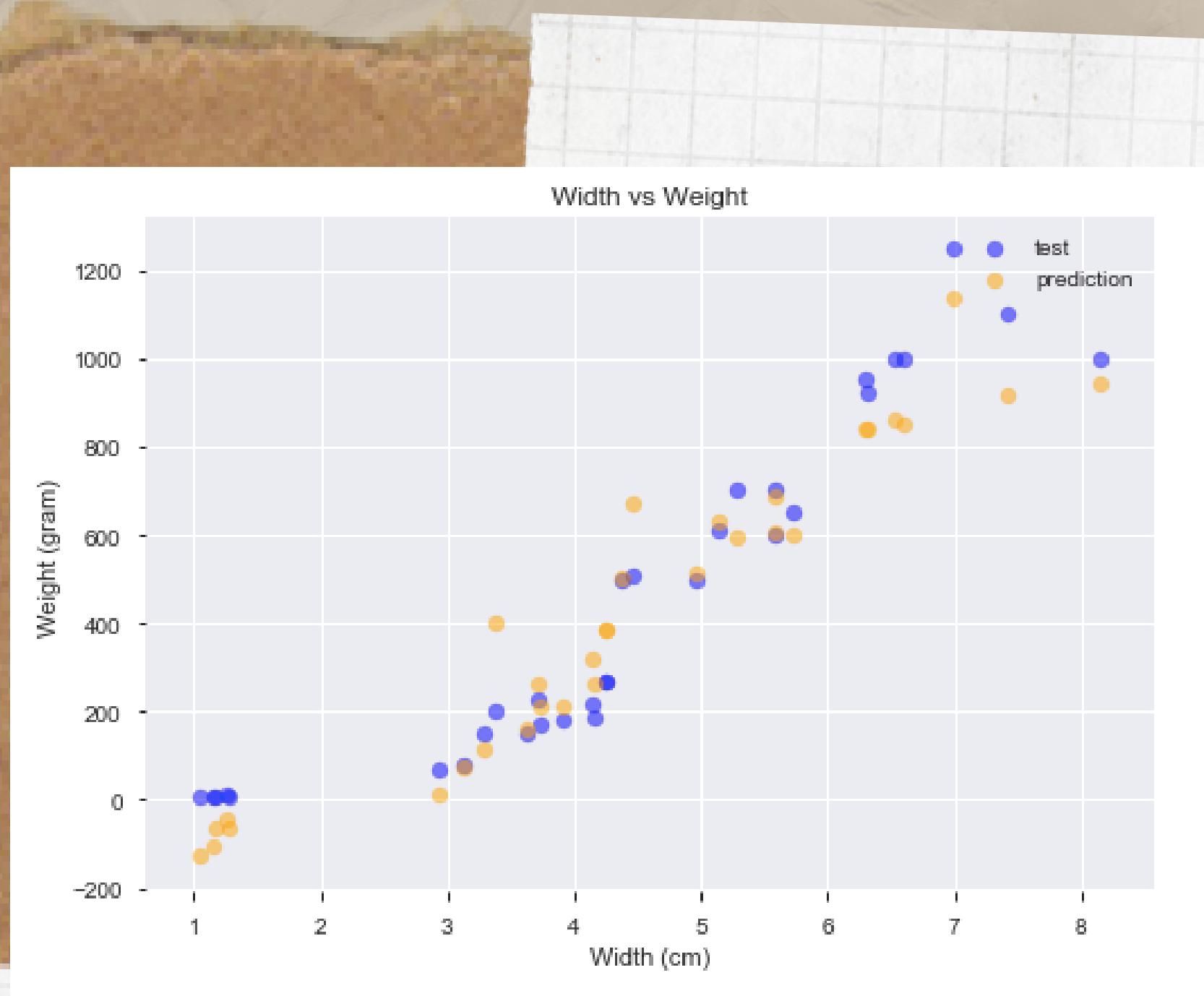
Vertical and Diagonal Length



Cross Length and Height



Width



SUPPORT VECTOR REGRESSION

```
from sklearn.svm import SVR  
  
svr = SVR(kernel='linear')  
svr.fit(x_train, y_train)  
y_pred = svr.predict(x_test)
```

R2: 0.9349
MSE: 0.0330

R2 SCORES

C = [0.0001, 0.001, 0.01,
0.1, 1, 10, 100]

[-0.1062827014720864,
0.6161316444213405,
0.8566696778533788,
0.9116113655775042,
0.9349222046340373,
0.9364409153802398,
0.9361393307666539]

TOL = [0.0001,
0.001, 0.01, 0.1, 1,
10, 100]

[0.9347963924796352,
0.9349222046340373,
0.9357756312722021,
0.9306948621802559,
0.8568955679912522,
-1.0200386582113041,
-1.0200386582113041]

EPS = [0.0001,
0.001, 0.01, 0.1, 1,
10, 100]

[0.9254868652460827,
0.9258435358742665,
0.9285374838240378,
0.9349222046340373,
0.3486499222581494,
-1.0200386582113037,
-1.0200386582113161]

MEAN SQUARE ERROR

C = [0.0001, 0.001, 0.01,
0.1, 1, 10, 100]

[0.5601881897341542,
0.19437935612813315,
0.0725781516699914,
0.044757338286430026,
0.03295343254435603,
0.03218440323339159,
0.0323371165846096]

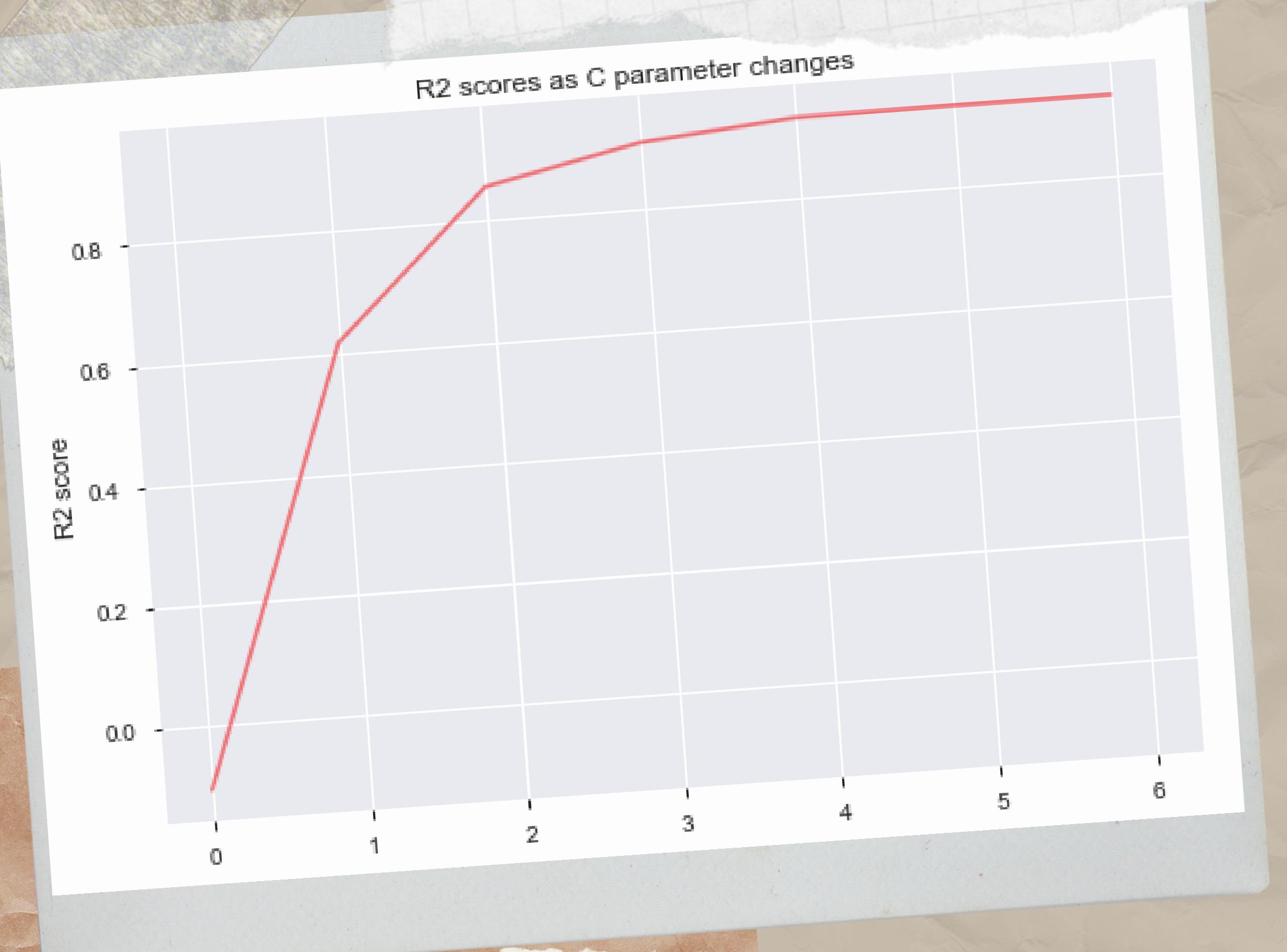
TOL = [0.0001,
0.001, 0.01, 0.1, 1,
10, 100]

[0.03301714002430427,
0.03295343254435603,
0.03252128303784359,
0.03509403124794151,
0.072463767718,
1.0228866433784265,
1.0228866433784325]

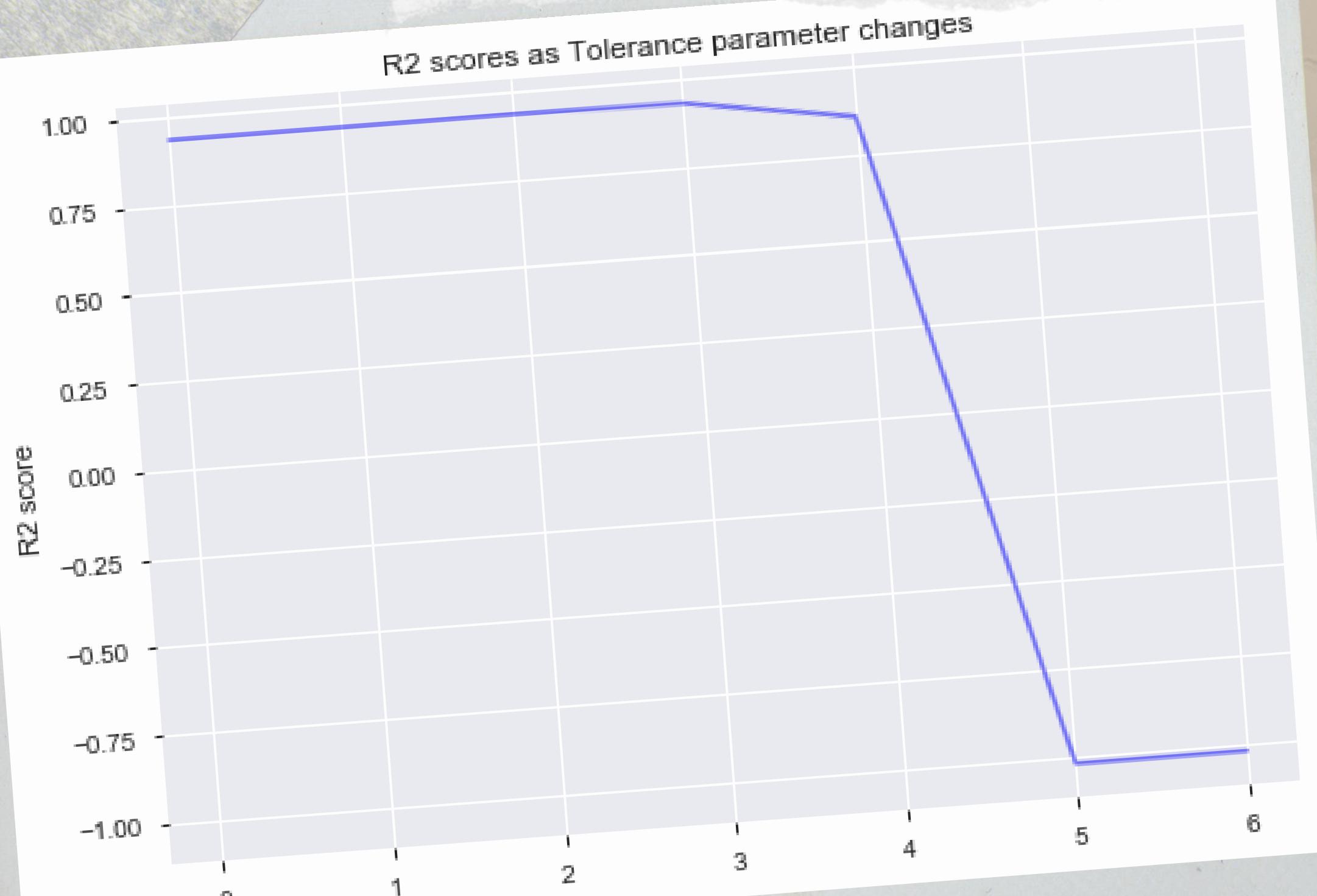
EPS = [0.0001,
0.001, 0.01, 0.1, 1,
10, 100]

[0.03773120379959882,
0.037550596552221864,
0.03618646257162872,
0.03295343254435603,
0.3298240318210512,
1.0228866433784263,
1.0228866433784325]

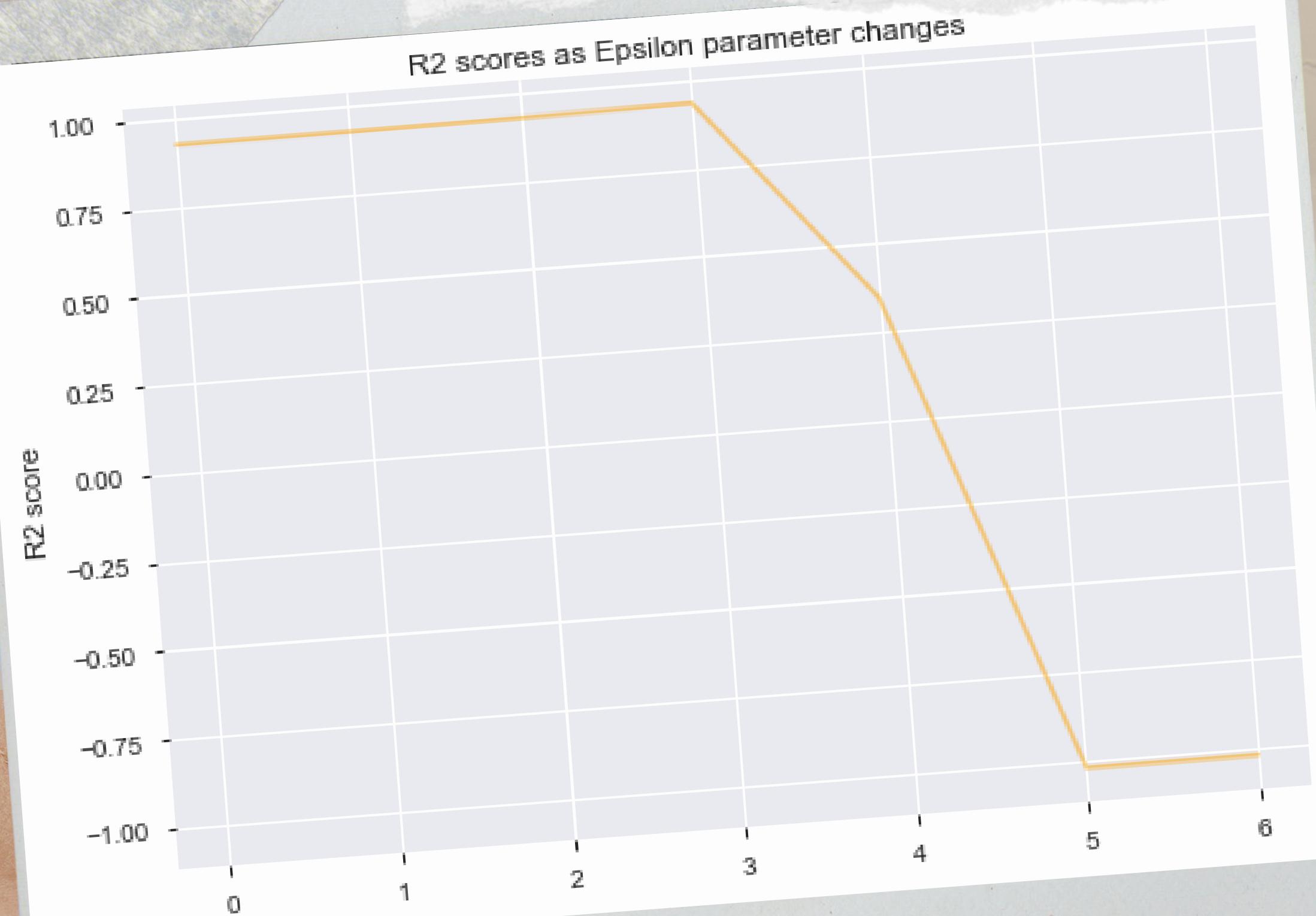
Plotting R2 Scores with C changes



Plotting R2 Scores with
Tolerance changes



Plotting R2 Scores with Epsilon
changes



WHAT I'VE LEARNED



R2 SCORE AND MSE

R2 is the accuracy in regression. MSE is how far your prediction is with actual values.

C, TOLERANCE, EPSILON

Own ways of affecting R2 and MSE

SVR IS MORE FLEXIBLE

Has more hyper parameters that you can change compared to LinearRegression

ROBUST SCALER

When dealing with huge values, it's easier to interpret MSE when you've scaled your features



THANKS

QUESTIONS?