

Section 1: Explaining the problem

Last summer, I read a book on how to use the pomodoro technique correctly. Pomodoro is a productivity technique that helps you divide your work time into focused work and breaks, while you eradicate distractions and maintain a laser-like focus on your work. You start with a 25 minutes of focused working with no distractions, and then follow it with 5 minutes of break and then do that again. I love using pomodoros whenever I am working. Any piece of work I have ever created since then was likely done while I had my small timer on, towards the right of my computer, counting down the minutes and seconds until my next break.

Another important aspect of using pomodoros, is keeping track of them. You typically have a list of tasks that you want to accomplish within a time frame, you estimate how many pomodoros (or 25 minute chunks of work) it will take you to do them and then you complete your task and see how many pomodoros it actually took you to finish it. This is how that looks like for me:

☑ Aa Task	👁 Likely to Procrastina...	📅 Deadline	👁 Goal	☰ Actual	☰ Fee
☑ SS111 LBA (2 Pomodoro)	Low	October 2, 2023	🍅🍅	🍅🍅🍅🍅🍅🍅	Good t working would l assign
☑ CS156 finish	Medium	October 2, 2023	🍅🍅	🍅🍅	Done v
☑ Figure out mentor thingy	Medium	October 2, 2023	🍅	🍅	Uno, d
☑ SMEDA review by this weekend	High 🔥	October 1, 2023	🍅🍅🍅	🍅🍅🍅🍅🍅🍅🍅🍅🍅🍅	Probab than th show, b Oct 5 f
☑ Work-study (2 Pomodoro)	High 🔥	October 1, 2023	🍅🍅	🍅🍅	One m great a produc
☑ Finance wrap-up + boxing route	High 🔥	October 1, 2023	🍅	🍅	halfwa
☑ Capstone (work for 4-5 Pomodoro)	Medium	October 1, 2023	🍅🍅🍅🍅🍅	🍅🍅	Could overall
☑ SS111 (2 sessions)	Medium	October 1, 2023	🍅🍅	🍅🍅🍅🍅	Only o throug

One problem I run into a lot is that I grossly underestimate how long it will take me to do a certain task. Many tasks take a lot more pomodoros than I anticipate, some less, so I wondered if I could design a machine learning algorithm that can analyze my estimates and actual pomodoros, and then be able to accurately predict how many pomodoros it will take me to actually do a certain task based on how many I estimate.

Section 2: Explaining the data

Since all my data was in the form of tomato emojis on a notion chart, it would've been quite complicated to export it and convert it to numerical data. Although there are python libraries that can help incorporate emojis in your code like pythonji. I only have 15 data points, which means that it is doable and much easier to export them into python manually, which I did in a list of lists. Then when plotting them, I broke the list of lists into two lists, one for the x values in the data and another for the y values. The reason why I decided on a list was because it was exceptionally important my data was ordered, since each x value is linked to a y value and if they were jumbled then my data would be meaningless. Lists also allow duplicate values which was important since my data does have duplicate values which is helpful for making stronger predictions. It is also convenient that a list is changeable, since it would help me add more data points in the future and make my predictions better.

Here is a plot of my last 15 pomodoros. On the x-axis we have the estimated times and on the y-axis we have the actual times.

```
In [62]: 1 #Lists of lists created manually from my notion.
2 Total = ([[3, 5],
3           [5, 4],
4           [1, 2],
5           [2, 3],
6           [1, 1],
7           [1, 2],
8           [2, 6],
9           [1, 3],
10          [2, 3],
11          [1, 1],
12          [1, 2],
13          [1, 1],
14          [3, 10],
15          [2, 3],
16          [1, 1]])
17
18
```

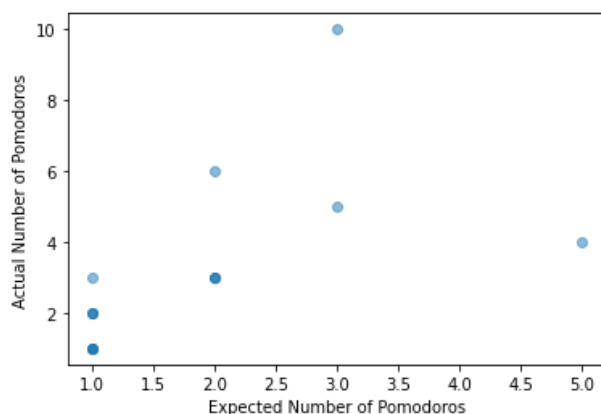
Section 3:

As mentioned above, I broke down my list of lists into two separate lists for x and y values, and then I plotted them to get an idea of the data visually. The visualization was quite anti-climatic since we have a lot of repeating values. I made the data points transparent, by setting alpha to 0.5, so that the repetitions are more obvious. The data is discrete, which can make it harder to process, and it is repeating quite a lot as well which is not a bad thing since that tells me that there is a regular pattern to how I estimate a pomodoro and how the actual number of pomodoros play out.

```
In [63]: 1 import matplotlib.pyplot as plt
2
3 #Breaking down the lists for plotting
4 DataX = [x[0] for x in Total]
5 DataY = [x[1] for x in Total]
6
7 print("List 1:", DataX)
8 print("List 2:", DataY)
9
10 plt.xlabel("Expected Number of Pomodoros")
11 plt.ylabel("Actual Number of Pomodoros")
12 plt.scatter(DataX, DataY, alpha=0.5)
13 plt.show
```

```
List 1: [3, 5, 1, 2, 1, 1, 2, 1, 2, 1, 1, 1, 3, 2, 1]
List 2: [5, 4, 2, 3, 1, 2, 6, 3, 3, 1, 2, 1, 10, 3, 1]
```

```
Out[63]: <function matplotlib.pyplot.show(close=None, block=None)>
```



Explaining my decisions:

My data and goal, makes this a regression problem. We have two variables and we want to model a relationship between them in a way that makes it easier to predict one based on the other.

Out of all the potential tools I could've used to solve this, I ended up deciding on a Regression Analysis and a Regression Tree. The reason I used these two methods, was because I wasn't sure if my data was linear or not and couldn't decide which one would give me more accurate results. Other methods like Random forests, Support vector

machines, and neural networks could also be used, and could potentially provide better results, but I don't have a lot of data points, which would undermine the use of these more complicated models since using these models would take a longer time, more processing and their effectiveness would be capped by the limited data. Thus the two chosen methods are adequate for the task at hand.

I analysed the results of my two models through mean absolute error and root mean squared error.

I was also quite nervous about my limited data and felt like splitting it 80/20 for training and testing would make the dataset even smaller. So I utilized a cross-validation technique of LOOCV or Leave One Out Cross Validation to further check the robustness of my models. I avoided K-fold techniques cause they expect the data to be split into many 'sets' which I didn't want to do, given my limited data.

Once I was done, I wondered if my predictions and models could get even better if I separated my data set into the tasks that I added another variable with binary values. 1 for if I genuinely enjoyed doing a task and 0 for the tasks that were chores. Thus I divided my data and then did both, the regression analysis and the regression tree on this split data set. I also made use of LOOCV as a cross-validation technique for both of these models as well.

I start out in each part by fitting a good old descriptive regression to my data just to see what it looked like.

Explaining the math:

1. Linear Regression based on an 80/20 Split:

- Linear regression involves fitting a linear equation to the observed data. An 80/20 split indicates that 80% of the data is used for training the model, and the remaining 20% is used for testing its performance.
- Formula: The formula for simple linear regression is $(Y = \beta_0 + \beta_1 X + \epsilon)$, where (Y) is the dependent variable, (X) is the independent variable, (β_0) and (β_1) are the regression coefficients, and (ϵ) is the error term.

2. Regression Trees based on an 80/20 Split:

- Regression trees involve partitioning the predictor space into distinct regions based on the values of the predictors. An 80/20 split implies that 80% of the data is used for training the regression tree, and the remaining 20% is used for testing the tree's performance.
- Formula: The splitting criterion in this case involves measures of variance or mean squared error, to determine the optimal splits at each node.

3. LOOCV (Leave-One-Out Cross-Validation):

- LOOCV is a method for assessing the performance of a model. It involves splitting the dataset into two parts, training and testing, where one data point is used for testing and the remaining points are used for training. This process is repeated for each data point, and the average error is computed.
- Formula: The LOOCV error is computed as the average of the errors obtained from each run of the leave-one-out procedure.

4. MAE (Mean Absolute Error):

- It measures the average of the absolute differences between predicted and actual values. It gives an idea of the magnitude of the error, without considering their direction.
- Formula: $(MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|)$, where (n) is the number of data points, (y_i) is the actual value, and (\hat{y}_i) is the predicted value.

5. RMSE (Root Mean Squared Error):

- It measures the square root of the average of the squared differences between predicted and actual values. It gives a relatively higher weight to large errors.
- Formula: $(RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2})$, where (n) is the number of data points, (y_i) is the actual value, and (\hat{y}_i) is the predicted value.

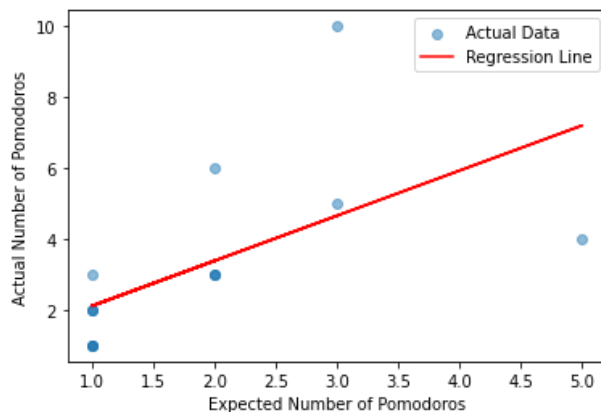
Sections 4, 5, 6, 7 and 8:

The task that would work best with the data and the goal would be a regression since I am trying to get a clear idea of how much I underestimate when assigning pomodoros to a task. Although the data is discrete, since Pomodoros are typically not divided, I decided to run the regression as if they were continuous and then round up to the nearest whole number when I get a good prediction model that predicts the number of pomodoros required for a task.

Visualizing a Linear Regression

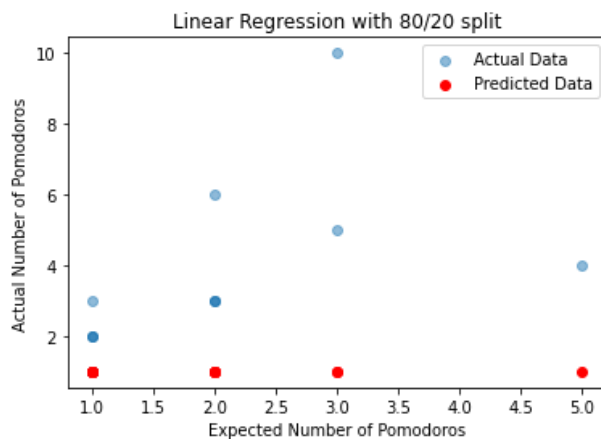
```
In [64]: 1 import numpy as np
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import mean_squared_error, mean_absolute_error
4 import matplotlib.pyplot as plt
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import PolynomialFeatures
7 from sklearn.model_selection import LeaveOneOut
8 from sklearn.tree import DecisionTreeRegressor
9
10 # Reshape the data to be 2D (required by scikit-learn)
11 DataX = np.array(DataX).reshape(-1, 1)
12 DataY = np.array(DataY)
13
14 # Perform an 80/20 train-test split
15 X_train, X_test, y_train, y_test = train_test_split(DataX, DataY, test_size=0.8)
16
17
```

```
In [65]: 1 # Create a linear regression model
2 model = LinearRegression()
3
4 # Fit the model to your data
5 model.fit(DataX, DataY)
6
7 # Make predictions
8 y_pred = model.predict(DataX)
9
10 # Visualize the data and regression line
11 plt.scatter(DataX, DataY, alpha=0.5, label="Actual Data")
12 plt.plot(DataX, y_pred, color='red', label="Regression Line")
13 plt.xlabel("Expected Number of Pomodoros")
14 plt.ylabel("Actual Number of Pomodoros")
15 plt.legend()
16 plt.show()
17
18
```



A basic linear regression is not very impressive, just looking at the line of best fit tells us that this prediction is quite inaccurate and the mean squared error is also quite high.

```
In [66]: 1 # Create polynomial features
2 poly_features = PolynomialFeatures(degree=2)
3 X_train_poly = poly_features.fit_transform(X_train)
4 X_test_poly = poly_features.transform(X_test)
5
6 # Create a polynomial regression model
7 model = LinearRegression()
8
9 # Fit the model to the training data
10 model.fit(X_train_poly, y_train)
11
12 # Make predictions on the testing data
13 y_pred = model.predict(X_test_poly)
14
15 mse = mean_squared_error(y_test, y_pred)
16 # Calculate root mean squared error
17 rmse = np.sqrt(mse)
18
19
20 # Calculate mean absolute error
21 mae = mean_absolute_error(y_test, y_pred)
22
23 # Visualize the data and regression line
24 plt.scatter(DataX, DataY, alpha=0.5, label="Actual Data")
25 plt.scatter(X_test, y_pred, color='red', label="Predicted Data")
26 plt.xlabel("Expected Number of Pomodoros")
27 plt.ylabel("Actual Number of Pomodoros")
28 plt.title('Linear Regression with 80/20 split')
29 plt.legend()
30 plt.show()
31
32 # Print the regression coefficients, RMSE, and MAE
33 print("Intercept:", model.intercept_)
34 print("Coefficients:", model.coef_)
35 print("Root Mean Squared Error:", rmse)
36 print("Mean Absolute Error:", mae)
37
```



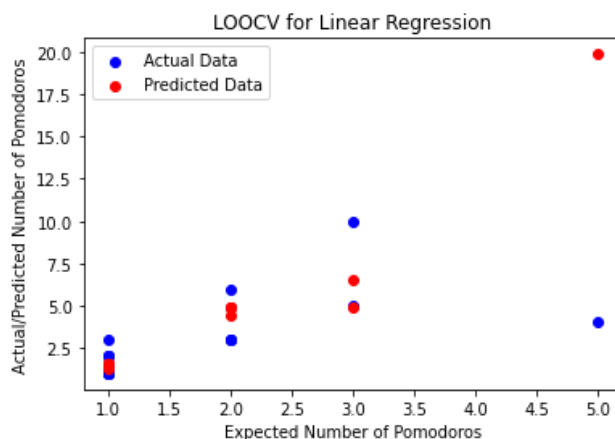
```
Intercept: 1.0
Coefficients: [0. 0. 0.]
Root Mean Squared Error: 3.5355339059327378
Mean Absolute Error: 2.6666666666666665
```

There is no "good" RMSE or MAE. Thus we will compare the models among themselves at the end and see which one does the best job.

```

In [67]: 1 # Initialize lists to store the predicted and actual values
2 y_pred_loocv = []
3 y_actual_loocv = []
4
5 loo = LeaveOneOut()
6 for train_index, test_index in loo.split(DataX):
7     X_train, X_test = DataX[train_index], DataX[test_index]
8     y_train, y_test = DataY[train_index], DataY[test_index]
9
10    poly_features = PolynomialFeatures(degree=2)
11    X_train_poly = poly_features.fit_transform(X_train)
12    X_test_poly = poly_features.transform(X_test)
13
14    model = LinearRegression()
15    model.fit(X_train_poly, y_train)
16
17    y_pred = model.predict(X_test_poly)
18
19    y_pred_loocv.append(y_pred[0])
20    y_actual_loocv.append(y_test[0])
21
22 # Convert lists to NumPy arrays
23 y_pred_loocv = np.array(y_pred_loocv)
24 y_actual_loocv = np.array(y_actual_loocv)
25
26 # Plotting the LOOCV predictions and actual values
27 plt.scatter(DataX, DataY, color='blue', label="Actual Data")
28 plt.scatter(DataX, y_pred_loocv, color='red', label="Predicted Data")
29 plt.xlabel("Expected Number of Pomodoros")
30 plt.ylabel("Actual/Predicted Number of Pomodoros")
31 plt.title('LOOCV for Linear Regression')
32 plt.legend()
33 plt.show()
34
35 # Calculate mean squared error for LOOCV
36 mse_loocv = mean_squared_error(y_actual_loocv, y_pred_loocv)
37 rmse_loocv = np.sqrt(mse_loocv)
38 mae_loocv = mean_absolute_error(y_actual_loocv, y_pred_loocv)
39
40 # Print mean squared error, root mean squared error, and mean absolute error
41 print(f"Root Mean Squared Error with LOOCV: {rmse_loocv}")
42 print(f"Mean Absolute Error with LOOCV: {mae_loocv}")
43

```



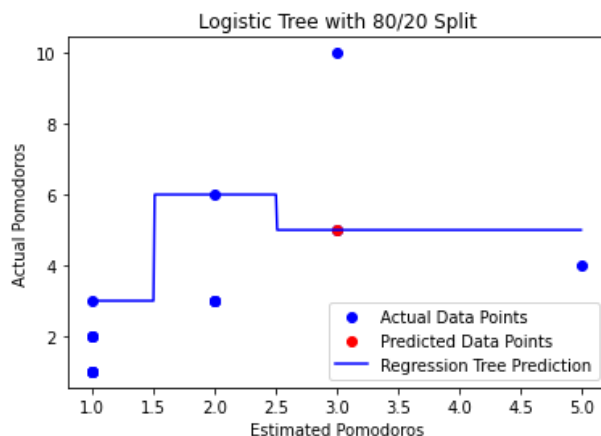
Root Mean Squared Error with LOOCV: 4.465143637755054
Mean Absolute Error with LOOCV: 2.3672462021880403

LOGISTIC REGRESSION

```
In [83]: 1 # Create and fit the regression tree on the training data
2 reg_tree_split = DecisionTreeRegressor(random_state=0)
3 reg_tree_split.fit(X_train, y_train)
4
5 # Make predictions on the test data
6 y_pred_split = reg_tree_split.predict(X_test)
7
8 # Compute evaluation metrics for the 80/20 split
9 mae_split = mean_absolute_error(y_test, y_pred_split)
10 mse_split = mean_squared_error(y_test, y_pred_split)
11 rmse_split = np.sqrt(mse_split)
12
13
14 # Print the evaluation metrics for the 80/20 split
15 print("Mean Absolute Error", mae_split)
16 print("Root Mean Squared Error:", rmse_split)
17
18
19 # Visualize the regression tree with the 80/20 split data
20 #X_grid = np.arange(min(DataX), max(DataX), 0.01).reshape(-1, 1)
21 plt.scatter(DataX, DataY, color='blue', label='Actual Data Points')
22 plt.scatter(X_test, y_pred_split, color='red', label='Predicted Data Points')
23 plt.plot(X_grid, reg_tree_split.predict(X_grid), color='blue', label='Regression Tree I
24 plt.title('Logistic Tree with 80/20 Split')
25 plt.xlabel('Estimated Pomodoros')
26 plt.ylabel('Actual Pomodoros')
27 plt.legend()
28 plt.show()
```

Mean Absolute Error 5.0

Root Mean Squared Error: 5.0

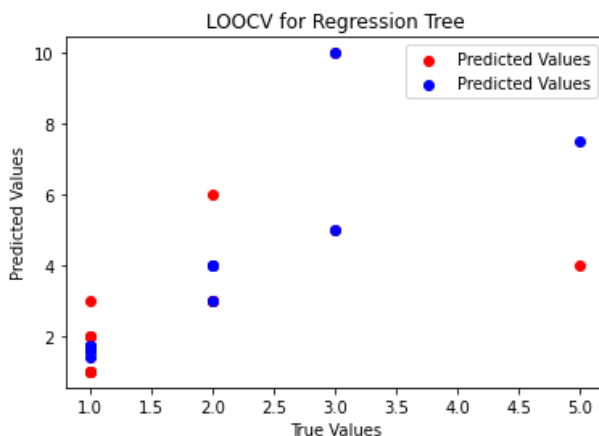


In [84]:

```
1 # Perform Leave-One-Out Cross-Validation
2 loo = LeaveOneOut()
3 y_true_cv, y_pred_cv = [], []
4
5 for train_index, test_index in loo.split(DataX):
6     X_train_cv, X_test_cv = DataX[train_index], DataX[test_index]
7     y_train_cv, y_test_cv = DataY[train_index], DataY[test_index]
8
9     # Fit the regression tree on the training data
10    reg_tree_cv = DecisionTreeRegressor(random_state=0)
11    reg_tree_cv.fit(X_train_cv, y_train_cv)
12
13    # Make predictions on the test data
14    y_pred_cv.extend(reg_tree_cv.predict(X_test_cv))
15    y_true_cv.extend(y_test_cv)
16
17 # Compute evaluation metrics for LOOCV
18 mae_cv = mean_absolute_error(y_true_cv, y_pred_cv)
19 mse_cv = mean_squared_error(y_true_cv, y_pred_cv)
20 rmse_cv = np.sqrt(mse_cv)
21
22 # Print the evaluation metrics for LOOCV
23 print(f"LOOCV - Mean Absolute Error: {mae_cv:.2f}")
24 print(f"LOOCV - Root Mean Squared Error: {rmse_cv:.2f}")
25
26 import matplotlib.pyplot as plt
27
28 # Create a scatter plot of true values vs predicted values for LOOCV
29 plt.scatter(DataX, DataY, color='red', label='Predicted Values')
30 plt.scatter(DataX, y_pred_cv, color='blue', label='Predicted Values')
31 plt.title('LOOCV for Regression Tree')
32 plt.xlabel('True Values')
33 plt.ylabel('Predicted Values')
34 plt.legend()
35 plt.show()
36
37
```

LOOCV - Mean Absolute Error: 1.68

LOOCV - Root Mean Squared Error: 2.30



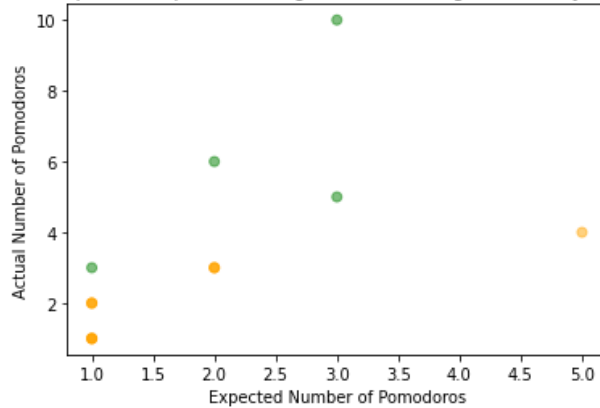
Adding more nuance for better predictions

This is better than the last model, but I am still not very happy with it. I won't trust this model to rightly predict the number of pomodoros for me and that makes me sad.

To make this model better, I have decided that I am going to add another variable, that I expect to be helpful. The variable would indicate what kind of task am I assigning pomodoros for, is it a task I am genuinely excited to work on or is it a task that feels more like a chore? The variable introduced as data z is binary, 1 indicating I was excited to work on the task and 0 if it was a chore.


```
In [70]: 1 DataZ = [1,0,0,0,0,0,1,1,0,0,0,0,1,0,0]
2
3 # Create a colormap
4 colors = ['green' if z == 1 else 'orange' for z in DataZ]
5
6 # Plotting the data
7 plt.scatter(DataX, DataY, c=colors, alpha=0.5)
8 plt.xlabel("Expected Number of Pomodoros")
9 plt.ylabel("Actual Number of Pomodoros")
10 plt.title('Expected and predicted pomodoros, green for exciting tasks and yellow for chores')
11
12 plt.show()
```

Expected and predicted pomodoros, green for exciting tasks and yellow for chores

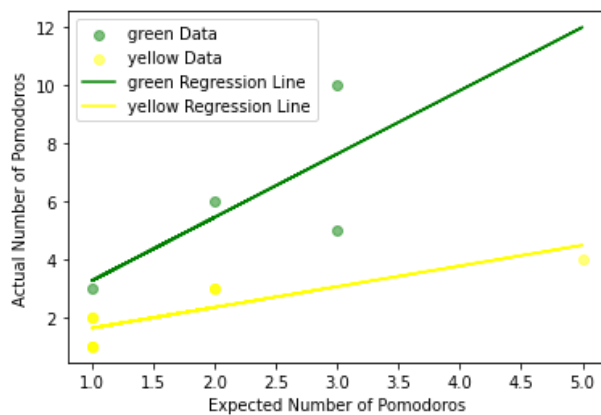


I am a little surprised to see that most of my tasks feel like a chore to me. I was expecting to find more of my tasks genuinely exciting. I should (at some other point in time) figure out why that is and see if I can figure out a way to make more of my tasks feel exciting.

```

In [71]: 1 # Filter data based on DataZ values
2 green_indices = [i for i, z in enumerate(DataZ) if z == 1]
3 yellow_indices = [i for i, z in enumerate(DataZ) if z == 0]
4
5 # Fit separate linear regression models for red and blue points
6 model_green = LinearRegression()
7 model_yellow = LinearRegression()
8
9 model_green.fit(DataX[green_indices], DataY[green_indices])
10 model_yellow.fit(DataX[yellow_indices], DataY[yellow_indices])
11
12 # Make predictions
13 y_pred_green = model_green.predict(DataX)
14 y_pred_yellow = model_yellow.predict(DataX)
15
16 # Calculate mean squared error
17 mse_green = mean_squared_error(DataY[green_indices], y_pred_green[green_indices])
18 mse_yellow = mean_squared_error(DataY[yellow_indices], y_pred_yellow[yellow_indices])
19
20 # Visualize the data and regression lines
21 plt.scatter(DataX[green_indices], DataY[green_indices], color='green', alpha=0.5, label='green Data')
22 plt.scatter(DataX[yellow_indices], DataY[yellow_indices], color='yellow', alpha=0.5, label='yellow Data')
23 plt.plot(DataX, y_pred_green, color='green', label="green Regression Line")
24 plt.plot(DataX, y_pred_yellow, color='yellow', label="yellow Regression Line")
25 plt.xlabel("Expected Number of Pomodoros")
26 plt.ylabel("Actual Number of Pomodoros")
27 plt.legend()
28 plt.show()
29

```



Oh wow, look at that. Our mean squared error actually reduces to be much lower when I differentiate between tasks that I genuinely enjoy verses tasks that feel like a chore.

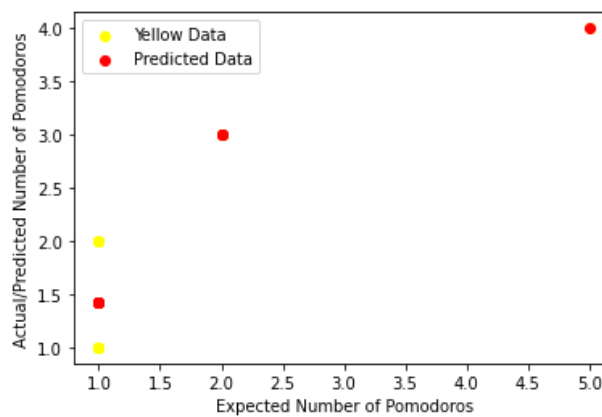
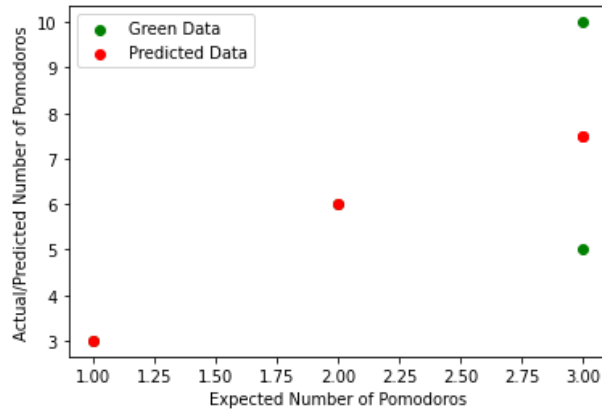
In [72]:

```
1 import numpy as np
2 from sklearn.preprocessing import PolynomialFeatures
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_squared_error
5 import matplotlib.pyplot as plt
6
7 # THE COLORS ARE CHANGED WITHIN THE VARIABLES SO I COULD TELL THE DIFFERENCE BETWEEN TI
8
9 # Your data
10 DataX = np.array([x[0] for x in Total]).reshape(-1, 1)
11 DataY = np.array([x[1] for x in Total])
12 DataZ = np.array([1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0])
13
14 # Red Data
15 red_X = DataX[DataZ == 1]
16 red_Y = DataY[DataZ == 1]
17
18 # Blue Data
19 blue_X = DataX[DataZ == 0]
20 blue_Y = DataY[DataZ == 0]
21
22 # Polynomial degree
23 degree = 2
24
25 # Create polynomial features for red data
26 poly_features_red = PolynomialFeatures(degree=degree)
27 red_X_poly = poly_features_red.fit_transform(red_X)
28 red_model = LinearRegression()
29 red_model.fit(red_X_poly, red_Y)
30 red_Y_pred = red_model.predict(red_X_poly)
31
32 # Create polynomial features for blue data
33 poly_features_blue = PolynomialFeatures(degree=degree)
34 blue_X_poly = poly_features_blue.fit_transform(blue_X)
35 blue_model = LinearRegression()
36 blue_model.fit(blue_X_poly, blue_Y)
37 blue_Y_pred = blue_model.predict(blue_X_poly)
38
39 # Visualize the data and predicted values for red data
40 plt.scatter(red_X, red_Y, color='green', label="Green Data")
41 plt.scatter(red_X, red_Y_pred, color='red', label="Predicted Data")
42 plt.xlabel("Expected Number of Pomodoros")
43 plt.ylabel("Actual/Predicted Number of Pomodoros")
44 plt.legend()
45 plt.show()
46
47 # Visualize the data and predicted values for blue data
48 plt.scatter(blue_X, blue_Y, color='yellow', label="Yellow Data")
49 plt.scatter(blue_X, blue_Y_pred, color='red', label="Predicted Data")
50 plt.xlabel("Expected Number of Pomodoros")
51 plt.ylabel("Actual/Predicted Number of Pomodoros")
52 plt.legend()
53 plt.show()
54
55
56 # Calculate RMSE and MAE for red data
57 rmse_red = np.sqrt(mean_squared_error(red_Y, red_Y_pred))
58 mae_red = np.mean(np.abs(red_Y - red_Y_pred))
59
60 # Calculate RMSE and MAE for blue data
61 rmse_blue = np.sqrt(mean_squared_error(blue_Y, blue_Y_pred))
62 mae_blue = np.mean(np.abs(blue_Y - blue_Y_pred))
63
64 # Print the regression coefficients and mean squared errors for red data
65 print("Green Data:")
66 print("Intercept:", red_model.intercept_)
67 print("Coefficient:", red_model.coef_)
68 print("Mean Squared Error:", mean_squared_error(red_Y, red_Y_pred))
69 print("Root Mean Squared Error:", rmse_red)
70 print("Mean Absolute Error:", mae_red)
71
72 # Print the regression coefficients and mean squared errors for blue data
73 print("\nYellow Data:")
74 print("Intercept:", blue_model.intercept_)
75 print("Coefficient:", blue_model.coef_)
76 print("Mean Squared Error:", mean_squared_error(blue_Y, blue_Y_pred))
```

```

77 print("Root Mean Squared Error:", rmse_blue)
78 print("Mean Absolute Error:", mae_blue)
79
80
81
82 # Calculate average RMSE and MAE
83 avg_rmse = (rmse_red + rmse_blue) / 2
84 avg_mae = (mae_red + mae_blue) / 2
85
86 # Print the average errors
87 print("\nAverage RMSE:", avg_rmse)
88 print("Average MAE:", avg_mae)
89

```



Green Data:

Intercept: -1.4999999999999636

Coefficient: [0. 5.25 -0.75]

Mean Squared Error: 3.125

Root Mean Squared Error: 1.7677669529663689

Mean Absolute Error: 1.2500000000000029

Yellow Data:

Intercept: -0.7619047619047592

Coefficient: [0. 2.5 -0.30952381]

Mean Squared Error: 0.15584415584415587

Root Mean Squared Error: 0.3947710169758614

Mean Absolute Error: 0.31168831168831224

Average RMSE: 1.0812689849711152

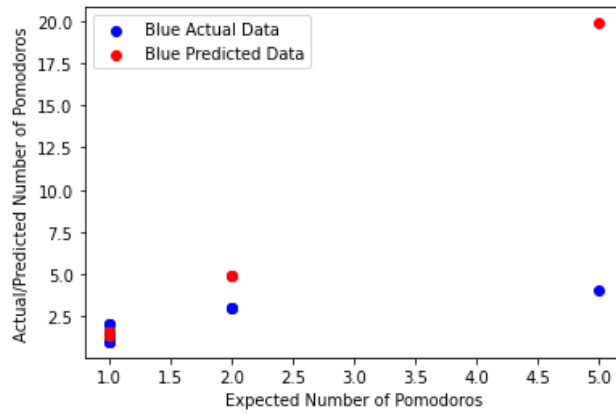
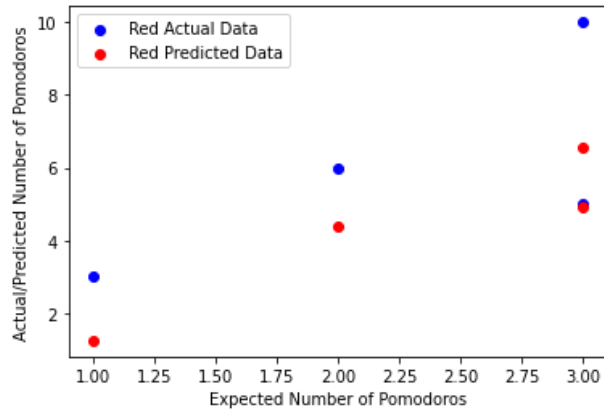
Average MAE: 0.7808441558441576

This is giving us even better error values, and I am willing to trusting our model more and more.

In [73]:

```
1 # Your data
2 DataX = np.array([x[0] for x in Total]).reshape(-1, 1)
3 DataY = np.array([x[1] for x in Total])
4 DataZ = np.array([1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0])
5
6 # Initialize lists to store the predicted and actual values
7 red_pred_loocv = []
8 red_actual_loocv = []
9 blue_pred_loocv = []
10 blue_actual_loocv = []
11
12 loo = LeaveOneOut()
13 for train_index, test_index in loo.split(DataX):
14     X_train, X_test = DataX[train_index], DataX[test_index]
15     y_train, y_test = DataY[train_index], DataY[test_index]
16     z_test = DataZ[test_index]
17
18     poly_features = PolynomialFeatures(degree=2)
19     X_train_poly = poly_features.fit_transform(X_train)
20     X_test_poly = poly_features.transform(X_test)
21
22     model = LinearRegression()
23     model.fit(X_train_poly, y_train)
24
25     y_pred = model.predict(X_test_poly)
26
27     if z_test == 1:
28         red_pred_loocv.append(y_pred[0])
29         red_actual_loocv.append(y_test[0])
30     else:
31         blue_pred_loocv.append(y_pred[0])
32         blue_actual_loocv.append(y_test[0])
33
34 # Convert lists to NumPy arrays
35 red_pred_loocv = np.array(red_pred_loocv).reshape(-1, 1)
36 red_actual_loocv = np.array(red_actual_loocv).reshape(-1, 1)
37 blue_pred_loocv = np.array(blue_pred_loocv).reshape(-1, 1)
38 blue_actual_loocv = np.array(blue_actual_loocv).reshape(-1, 1)
39
40 # Plotting the red LOOCV predictions and actual values
41 plt.scatter(DataX[DataZ == 1], red_actual_loocv, color='blue', label="Red Actual Data")
42 plt.scatter(DataX[DataZ == 1], red_pred_loocv, color='red', label="Red Predicted Data")
43 plt.xlabel("Expected Number of Pomodoros")
44 plt.ylabel("Actual/Predicted Number of Pomodoros")
45 plt.legend()
46 plt.show()
47
48 # Plotting the blue LOOCV predictions and actual values
49 plt.scatter(DataX[DataZ == 0], blue_actual_loocv, color='blue', label="Blue Actual Data")
50 plt.scatter(DataX[DataZ == 0], blue_pred_loocv, color='red', label="Blue Predicted Data")
51 plt.xlabel("Expected Number of Pomodoros")
52 plt.ylabel("Actual/Predicted Number of Pomodoros")
53 plt.legend()
54 plt.show()
55
56 # Calculate mean squared error for LOOCV for red and blue data
57 red_mse_loocv = mean_squared_error(red_actual_loocv, red_pred_loocv)
58 blue_mse_loocv = mean_squared_error(blue_actual_loocv, blue_pred_loocv)
59
60 # Calculate RMSE and MAE for red data
61 red_rmse_loocv = np.sqrt(red_mse_loocv)
62 red_mae_loocv = np.mean(np.abs(red_actual_loocv - red_pred_loocv))
63
64 # Calculate RMSE and MAE for blue data
65 blue_rmse_loocv = np.sqrt(blue_mse_loocv)
66 blue_mae_loocv = np.mean(np.abs(blue_actual_loocv - blue_pred_loocv))
67
68
69 # Calculate average RMSE and MAE
70 avg_rmse = (red_rmse_loocv + blue_rmse_loocv) / 2
71 avg_mae = (red_mae_loocv + blue_mae_loocv) / 2
72
73 # Print the average errors
74 print("\nAverage RMSE:", avg_rmse)
```

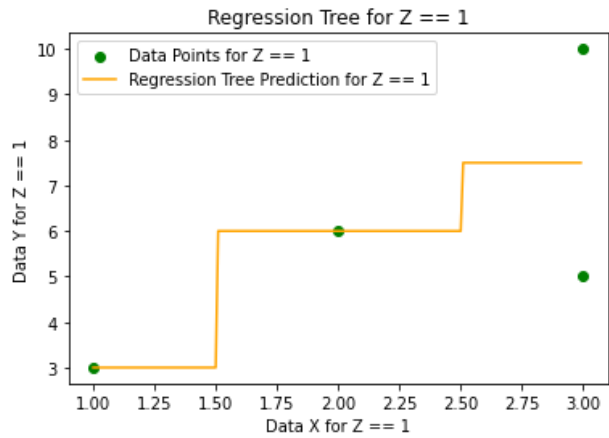
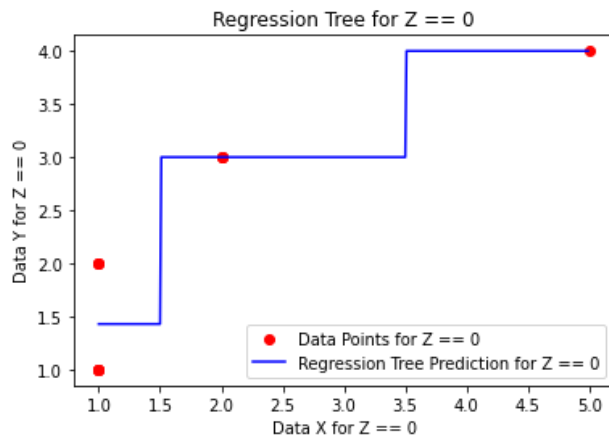
```
75 print("Average MAE:", avg_mae)
```



Average RMSE: 3.9134020704410872
Average MAE: 2.4110211329568534

In [74]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.tree import DecisionTreeRegressor
4
5 z = np.array([1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0])
6 DataX = np.array([x[0] for x in Total]).reshape(-1, 1)
7 DataY = np.array([x[1] for x in Total])
8
9 # Create and fit the first regression tree for z == 0
10 reg_tree_0 = DecisionTreeRegressor(random_state=0)
11 reg_tree_0.fit(DataX[z == 0], DataY[z == 0])
12
13 # Create and fit the second regression tree for z == 1
14 reg_tree_1 = DecisionTreeRegressor(random_state=0)
15 reg_tree_1.fit(DataX[z == 1], DataY[z == 1])
16
17 # Visualize the first regression tree for z == 0
18 X_grid_0 = np.arange(min(DataX[z == 0]), max(DataX[z == 0]), 0.01).reshape(-1, 1)
19 plt.scatter(DataX[z == 0], DataY[z == 0], color='red', label='Data Points for Z == 0')
20 plt.plot(X_grid_0, reg_tree_0.predict(X_grid_0), color='blue', label='Regression Tree 0')
21 plt.title('Regression Tree for Z == 0')
22 plt.xlabel('Data X for Z == 0')
23 plt.ylabel('Data Y for Z == 0')
24 plt.legend()
25 plt.show()
26
27 # Visualize the second regression tree for z == 1
28 X_grid_1 = np.arange(min(DataX[z == 1]), max(DataX[z == 1]), 0.01).reshape(-1, 1)
29 plt.scatter(DataX[z == 1], DataY[z == 1], color='green', label='Data Points for Z == 1')
30 plt.plot(X_grid_1, reg_tree_1.predict(X_grid_1), color='orange', label='Regression Tree 1')
31 plt.title('Regression Tree for Z == 1')
32 plt.xlabel('Data X for Z == 1')
33 plt.ylabel('Data Y for Z == 1')
34 plt.legend()
35 plt.show()
36
37 from sklearn.metrics import mean_squared_error, mean_absolute_error
38
39 # Make predictions for z == 0
40 y_pred_0 = reg_tree_0.predict(DataX[z == 0])
41 rmse_0 = np.sqrt(mean_squared_error(DataY[z == 0], y_pred_0))
42 mae_0 = mean_absolute_error(DataY[z == 0], y_pred_0)
43 print(f"For Z == 0 - RMSE: {rmse_0:.2f}")
44 print(f"For Z == 0 - MAE: {mae_0:.2f}")
45
46 # Make predictions for z == 1
47 y_pred_1 = reg_tree_1.predict(DataX[z == 1])
48 rmse_1 = np.sqrt(mean_squared_error(DataY[z == 1], y_pred_1))
49 mae_1 = mean_absolute_error(DataY[z == 1], y_pred_1)
50 print(f"For Z == 1 - RMSE: {rmse_1:.2f}")
51 print(f"For Z == 1 - MAE: {mae_1:.2f}")
52
53
54 # Calculate average RMSE and MAE
55 avg_rmse = (rmse_0 + rmse_1) / 2
56 avg_mae = (mae_0 + mae_1) / 2
57
58 # Print the average errors
59 print("\nAverage RMSE:", avg_rmse)
60 print("Average MAE:", avg_mae)
```



For Z == 0 - RMSE: 0.39

For Z == 0 - MAE: 0.31

For Z == 1 - RMSE: 1.77

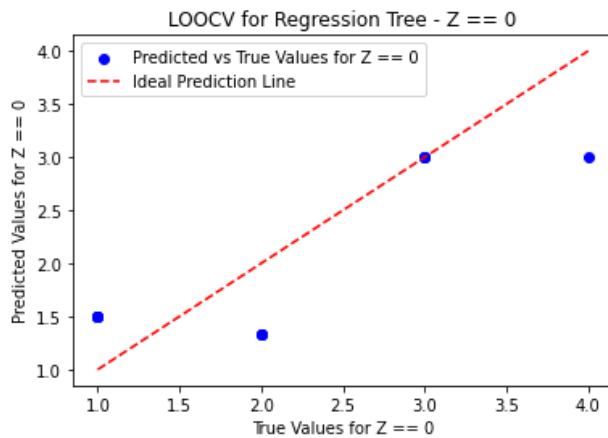
For Z == 1 - MAE: 1.25

Average RMSE: 1.0812689849711152

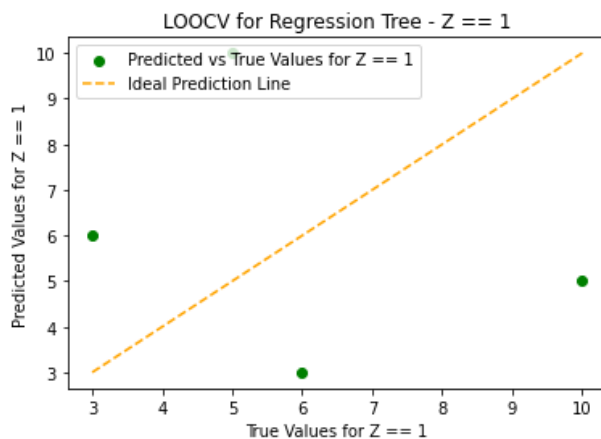
Average MAE: 0.7808441558441559

In [75]:

```
1 import matplotlib.pyplot as plt
2
3 # LOOCV for z == 0
4 loo = LeaveOneOut()
5 y_true, y_pred = [], []
6
7 for train_index, test_index in loo.split(DataX[z == 0]):
8     X_train, X_test = DataX[z == 0][train_index], DataX[z == 0][test_index]
9     y_train, y_test = DataY[z == 0][train_index], DataY[z == 0][test_index]
10
11     reg_tree_0_cv = DecisionTreeRegressor(random_state=0)
12     reg_tree_0_cv.fit(X_train, y_train)
13
14     y_pred.extend(reg_tree_0_cv.predict(X_test))
15     y_true.extend(y_test)
16
17 # Visualize LOOCV for z == 0
18 plt.scatter(y_true, y_pred, color='blue', label='Predicted vs True Values for Z == 0')
19 plt.plot([min(y_true), max(y_true)], [min(y_true), max(y_true)], color='red', linestyle='solid')
20 plt.title('LOOCV for Regression Tree - Z == 0')
21 plt.xlabel('True Values for Z == 0')
22 plt.ylabel('Predicted Values for Z == 0')
23 plt.legend()
24 plt.show()
25
26 # Compute evaluation metrics for LOOCV - z == 0
27 rmse_0_cv = np.sqrt(mean_squared_error(y_true, y_pred))
28 mae_0_cv = mean_absolute_error(y_true, y_pred)
29 print(f"LOOCV - Z == 0 - RMSE: {rmse_0_cv:.2f}")
30 print(f"LOOCV - Z == 0 - MAE: {mae_0_cv:.2f}")
31
32 # LOOCV for z == 1
33 y_true, y_pred = [], []
34
35 for train_index, test_index in loo.split(DataX[z == 1]):
36     X_train, X_test = DataX[z == 1][train_index], DataX[z == 1][test_index]
37     y_train, y_test = DataY[z == 1][train_index], DataY[z == 1][test_index]
38
39     reg_tree_1_cv = DecisionTreeRegressor(random_state=0)
40     reg_tree_1_cv.fit(X_train, y_train)
41
42     y_pred.extend(reg_tree_1_cv.predict(X_test))
43     y_true.extend(y_test)
44
45 # Visualize LOOCV for z == 1
46 plt.scatter(y_true, y_pred, color='green', label='Predicted vs True Values for Z == 1')
47 plt.plot([min(y_true), max(y_true)], [min(y_true), max(y_true)], color='orange', linestyle='solid')
48 plt.title('LOOCV for Regression Tree - Z == 1')
49 plt.xlabel('True Values for Z == 1')
50 plt.ylabel('Predicted Values for Z == 1')
51 plt.legend()
52 plt.show()
53
54 # Compute evaluation metrics for LOOCV - z == 1
55 rmse_1_cv = np.sqrt(mean_squared_error(y_true, y_pred))
56 mae_1_cv = mean_absolute_error(y_true, y_pred)
57 print(f"LOOCV - Z == 1 - RMSE: {rmse_1_cv:.2f}")
58 print(f"LOOCV - Z == 1 - MAE: {mae_1_cv:.2f}")
59
60 # Calculate average RMSE and MAE
61 avg_rmse = (rmse_0_cv + rmse_1_cv) / 2
62 avg_mae = (mae_0_cv + mae_1_cv) / 2
63
64 # Print the average errors
65 print("\nAverage RMSE:", avg_rmse)
66 print("Average MAE:", avg_mae)
```



LOOCV - Z == 0 - RMSE: 0.55
 LOOCV - Z == 0 - MAE: 0.45



LOOCV - Z == 1 - RMSE: 4.12
 LOOCV - Z == 1 - MAE: 4.00

Average RMSE: 2.3367937540904204
 Average MAE: 2.227272727272727

I was initially expecting LOOCV to do much better than testing and training since we have a small data-set and LOOCV is usually recommended for smaller data-sets, but it not performing as well which can be justified by the high variability in our data and LOOCV being more prone to overfitting and bad with highly variable data.

Section 9:

I defined the problem and my idea of what the solution would be I started out by plotting the data I manually input. Then the next steps followed this order

Regressions:

Linear Regression (80/20)

LOOCV

Regression Tree (80/20)

LOOCV

Not happy, split the data

Linear Regression (80/20)

LOOCV

Regression Tree (80/20)

LOOCV

Each step was followed by a visualization and a RMSE and MSE. Here is a chart detailing them all.

	Method	RMSE	MAE
Before splitting: Linear Regression (80/20 split)		1.8614445587797204	3.0288461538461533
Before splitting: LOOCV on Linear Regression		4.465143637755054	2.3672462021880403
Before splitting: Regression Tree (80/20 split)		0.7142857142857142	0.7142857142857142
Before splitting: LOOCV on Regression Tree		1.68	2.30
After splitting: Linear Regression (80/20 split)		1.0812689849711152	0.7808441558441576
After splitting: LOOCV on Linear Regression		3.9134020704410872	2.4110211329568534
After splitting: Regression Tree (80/20 split)		1.0812689849711152	0.7808441558441559
After splitting: LOOCV on Regression Tree		2.3367937540904204	2.227272727272727

Surprisingly the results our RMSE and MAE are the lowest for the regression tree before we split our data into tasks that I am excited to work on and tasks that I am not excited to work on. I was hoping for the split to be more helpful but it was not. It is also important to note that I averaged the results after the split so in some instances the RMSE and MAE values might be lower for one of the splits and higher for the other.

Proper predictor model:

```
In [76]: 1 # Assuming you have X_train, X_test, y_train, y_test already defined
2
3 # Create and fit the regression tree on the training data
4 reg_tree_split = DecisionTreeRegressor(random_state=0)
5 reg_tree_split.fit(X_train, y_train)
6
7 def predict_pomodoros(input_value):
8     predicted_pomodoros = reg_tree_split.predict(np.array(input_value).reshape(1, -1))
9     return predicted_pomodoros[0]
10
11 # Sample usage of the function
12 user_input = float(input("Enter the estimated number of pomodoros: "))
13 predicted_value = predict_pomodoros(user_input)
14 print(f"The predicted number of pomodoros is: {predicted_value}")
15
```

```
Enter the estimated number of pomodoros: 1
The predicted number of pomodoros is: 3.0
```

A problem with our predictor model is that it gives us, bad results when we input 5 or more pomodoros since the training data on it, is just one data point. This is why it should be used with caution and shouldn't be used to estimate more than 5 pomodoros. Ideally, if you follow the rules of the pomodoro, if you have a task that is estimated to take that long, then you are supposed break it up into two separate tasks to keep things manageable.

I like to think that this predictor is almost like Thor's hammer (which only Thor could pick up) in the sense that it only works for people who respect the rule book of the pomodoro and don't bite off more than they can chew, and thus they are deserving of getting better time estimates for their valuable time.

Section 10:

References:

https://www.w3schools.com/python/python_lists.asp

<https://www.datarobot.com/blog/how-much-data-is-needed-to-train-a-good-model/#:~:text=For%20example%2C%20if%20you%20have,observations%20to%20train%20a%20model>

<https://betterprogramming.pub/emojis-as-python-variables-sure-why-not-96ce955dada1>

AI statement:

I did use Chat GPT to help with some of the coding syntax, but I do have conceptual understanding of the code, I commented the code myself to make it more clear and I didn't use AI to write any of the explanations.

