

# Data Mining Final Project

By Camilla Sisemore



# Attempting to predict the Occupancy of a Room by using Sensor Data

One February, likely in 2015 (based on the publication date of their paper), Luis M. Candanedo and Véronique Feldheim, both of Université de Mons in Mons, Belgium performed an experiment to see if they could determine if a room was occupied based on readings from sensors that measured temperature, humidity, carbon dioxide, and light. The readings were taken three to four times each minute and averaged to give one result per minute per feature. A picture was taken once per minute so that occupancy could be determined by the researchers. The dataset also includes humidity ratio. I will be attempting to determine the occupancy based on combinations of temperature, humidity, carbon dioxide, light, and the humidity ratio. I suspect that CO<sub>2</sub> will definitely be one of the indicators.

<sup>1</sup> (Accurate occupancy detection of an office room from light, temperature, humidity and CO2 measurements using statistical learning models, Energy and Buildings Volume 112, 2016)



# The Data

The Train set initially consisted of 8,143 rows and 7 columns. The Train set needed to be cleaned. Every time I thought I had it clean, I was wrong. After cleaning, the train set had 8,140 rows. I used Python with the Pandas library in Jupyter Notebook.

Commands used (some more than once):

```
df.shape
df.isna().any()
df.dropna(axis = 0, how = 'any', subset=['Occupied'], inplace=True)
df['Occupied'].value_counts()
df['Occupied'] = df['Occupied'].str.replace(' ', '')
df['Occupied'] = df['Occupied'].replace('Y', 'yes', regex=True)
df['Occupied'] = df['Occupied'] = df['Occupied'].replace(['N'], 'no')
df['Occupied'] = df['Occupied'].replace(['yeses'], 'yes')
df['Occupied'] = df['Occupied'].replace(['No'], 'no')
```



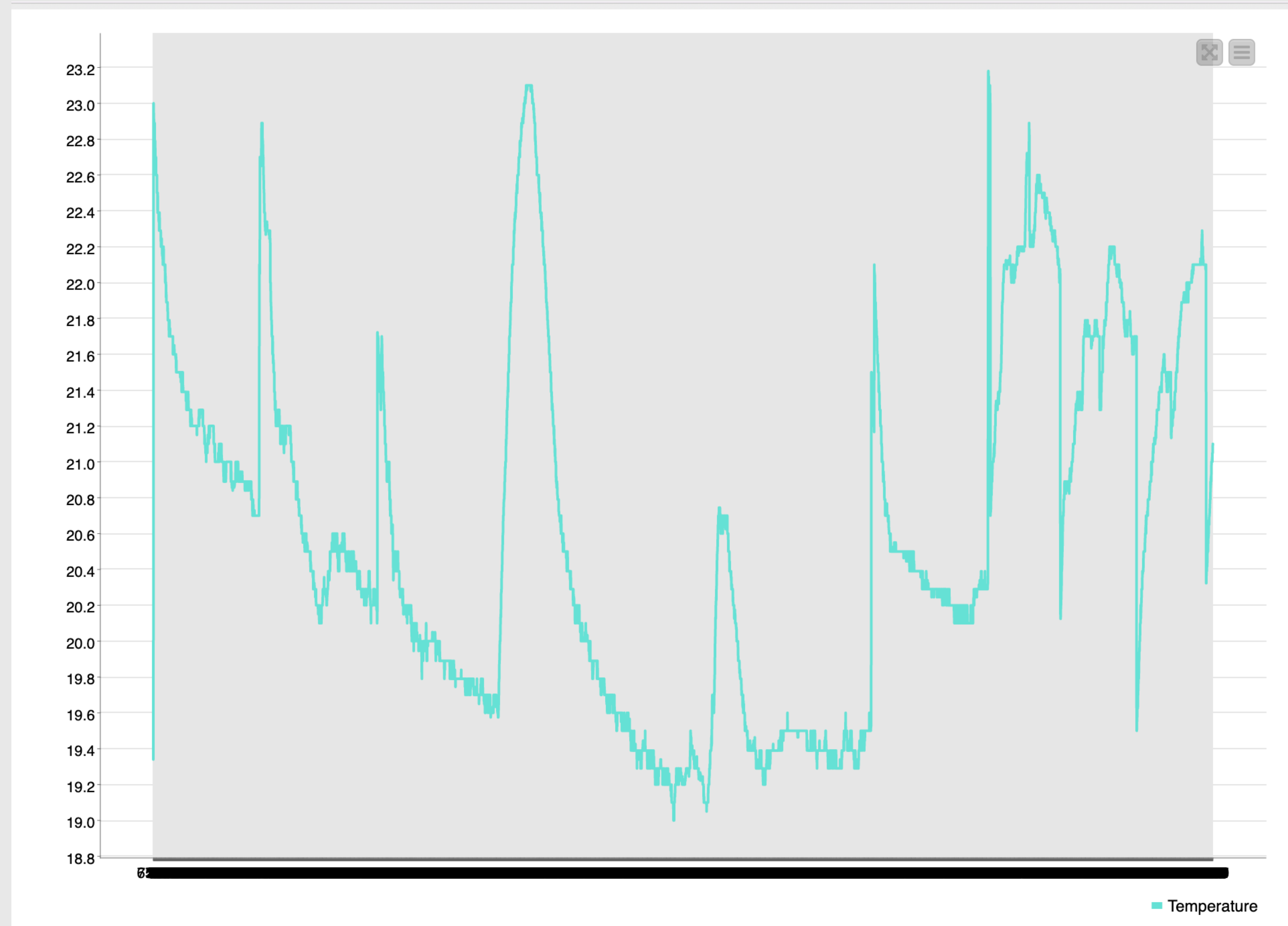
# Training Set Exploration

I needed to get the mean and median of the features. I used `df.mean()` and `df.median()` to do this.

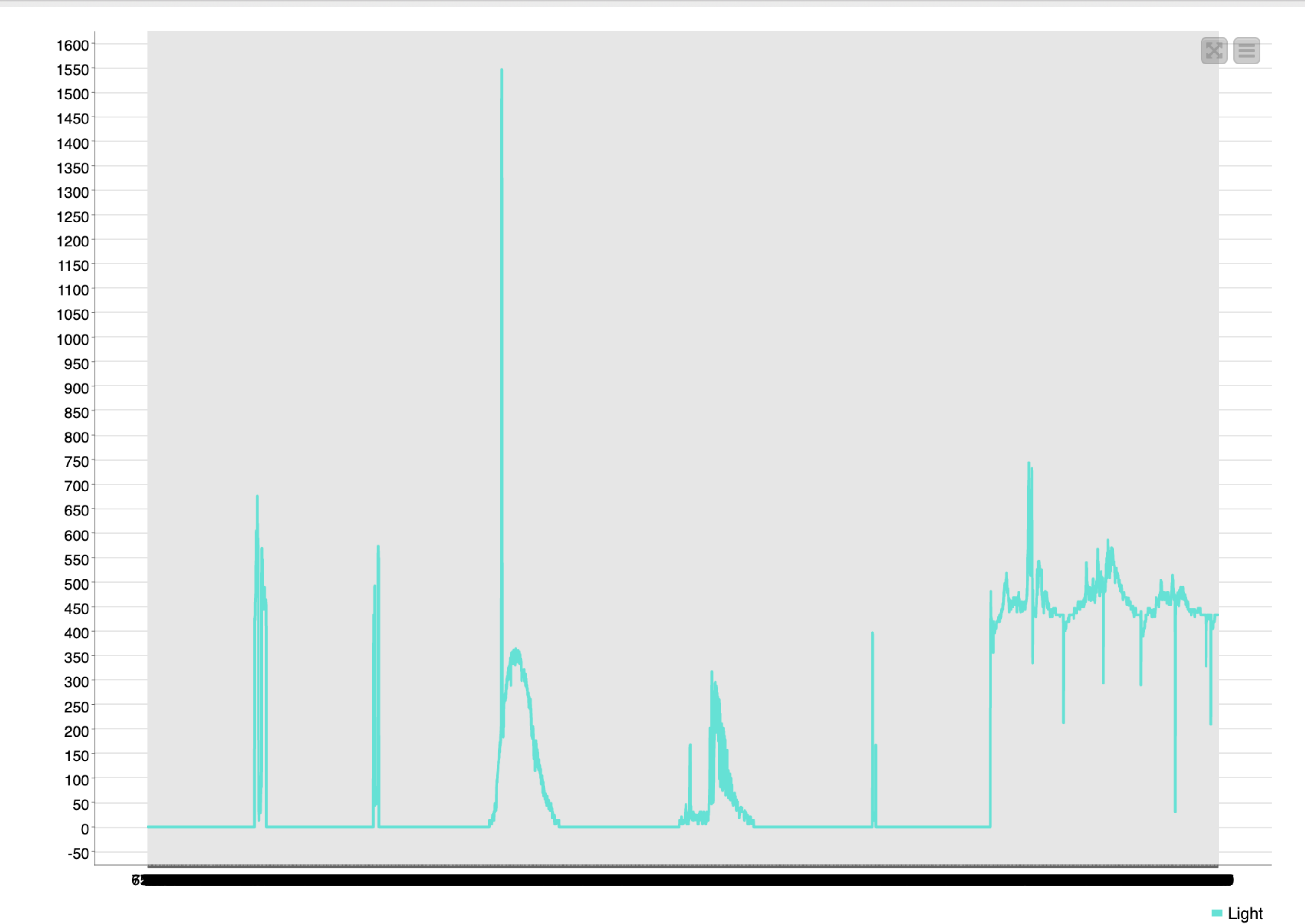
Feature	Mean	Median
Temperature	20.618656	20.390000
Humidity	25.731379	26.222500
Light	119.502981	0.000000
CO2	606.533668	453.500000
HumidityRatio	0.003862	0.003801



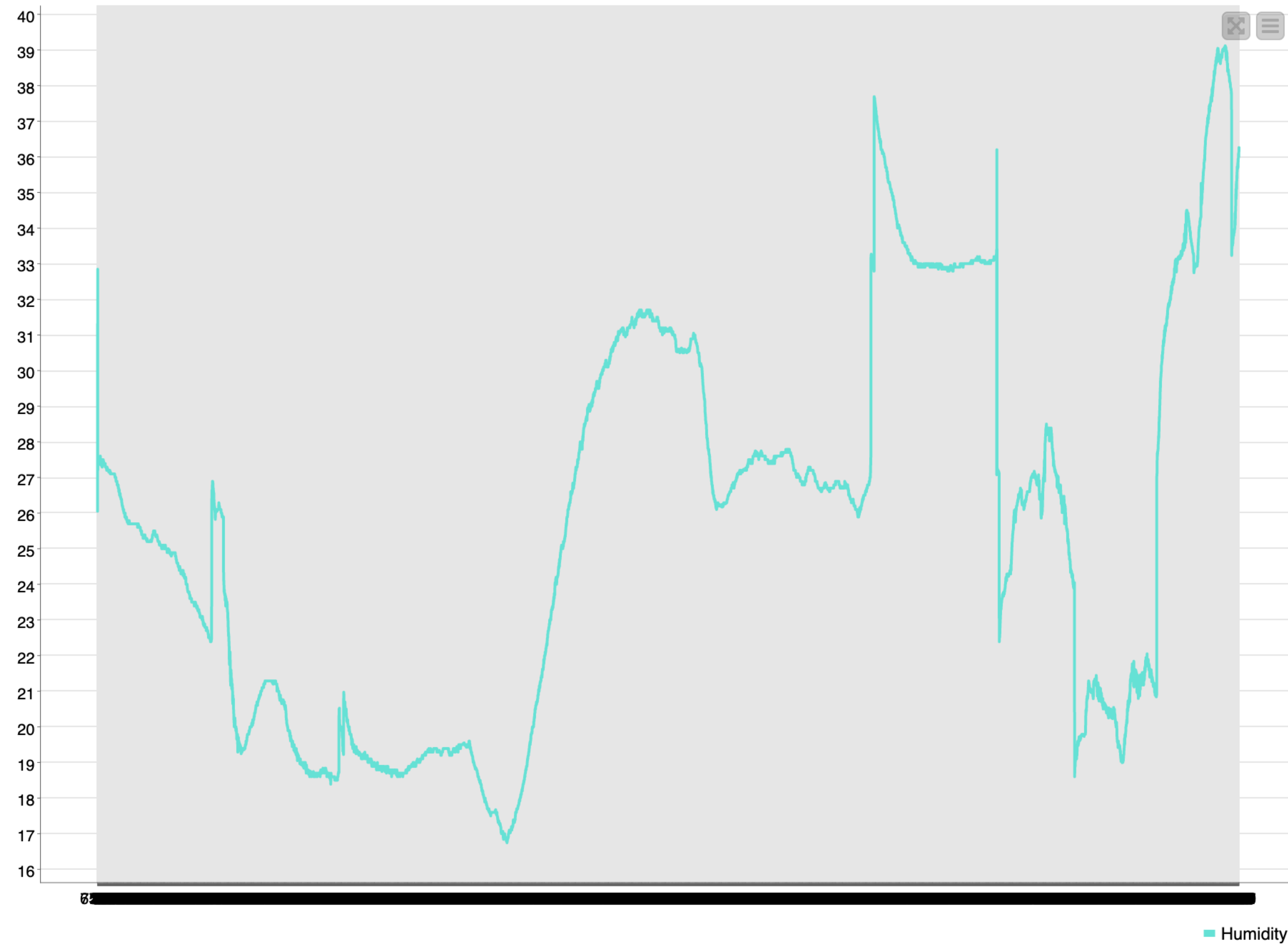
# Line Plot - Temperature



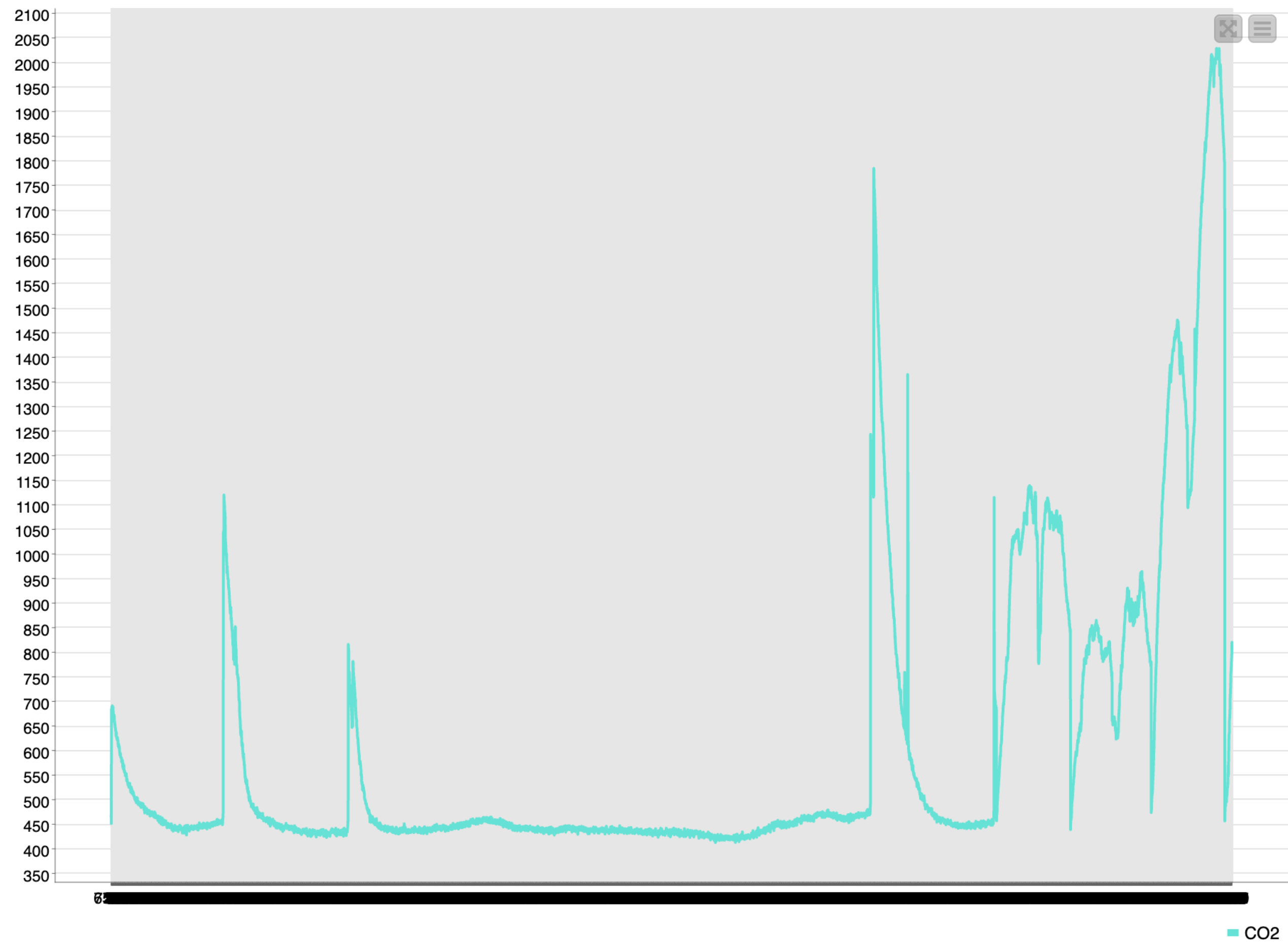
# Line Plot - Light



# Line Plot - Humidity

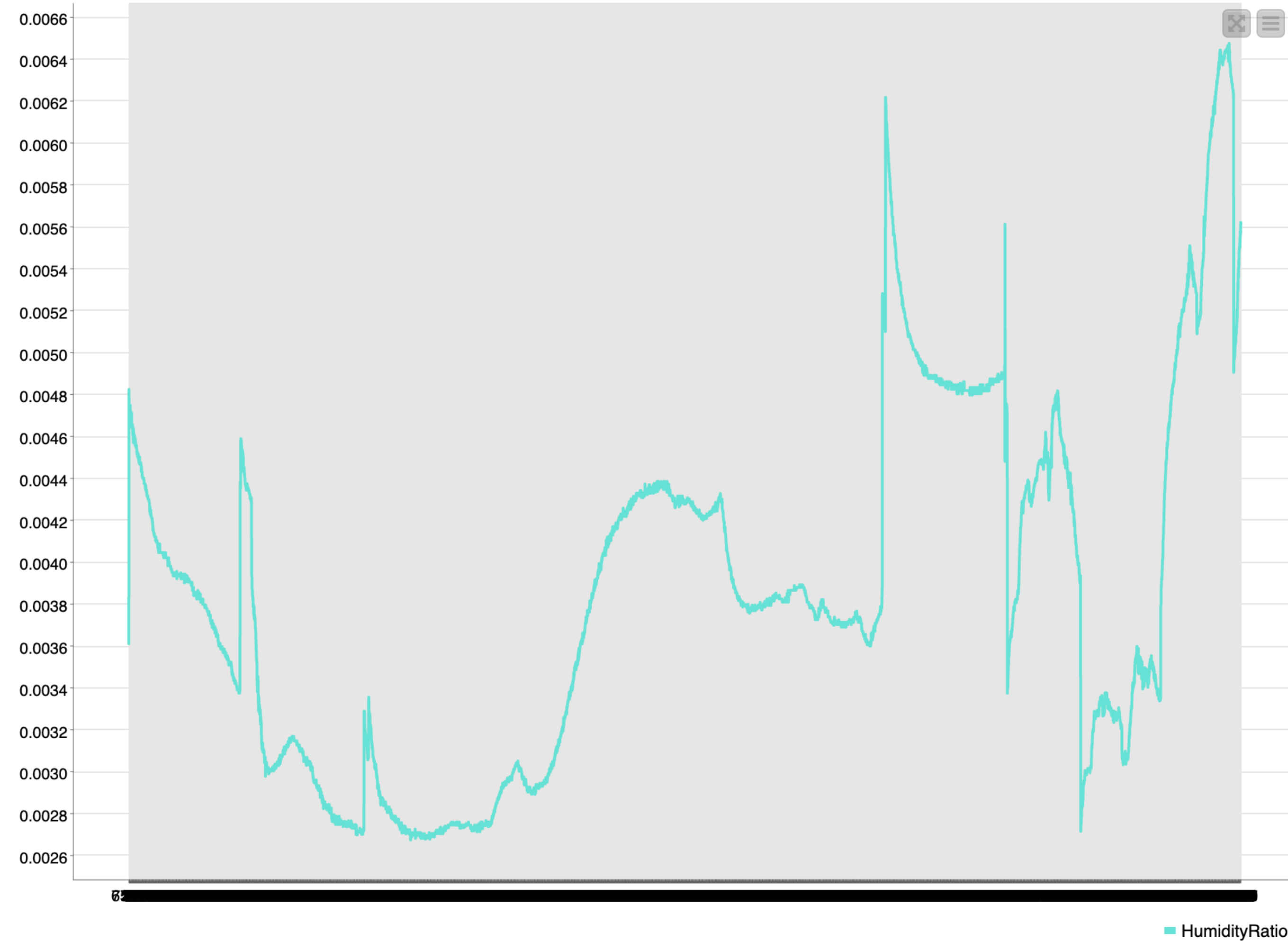


# Line Plot – Carbon Dioxide

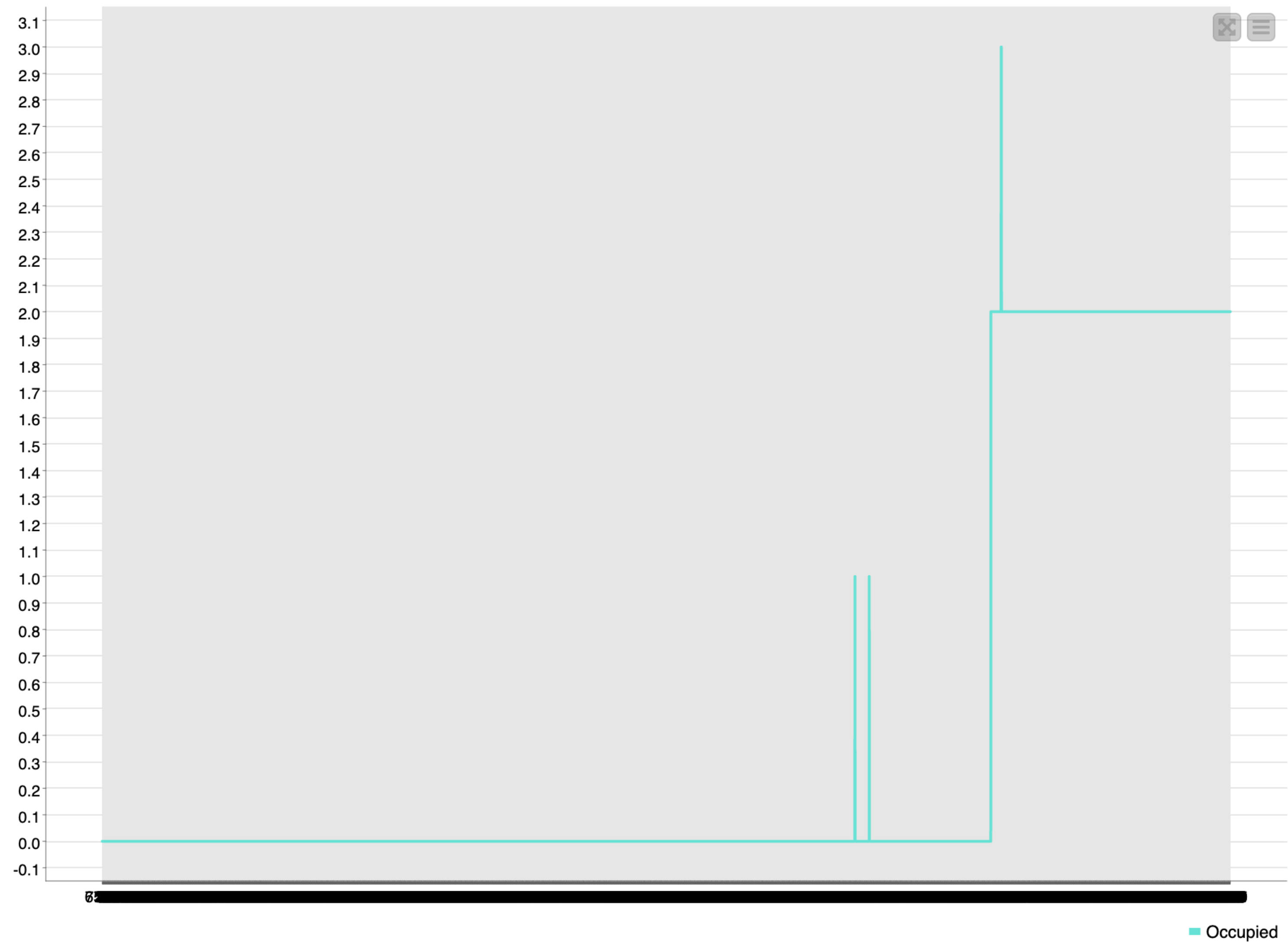




# Line Plot – Humidity Ratio



# Line Plot - Occupied



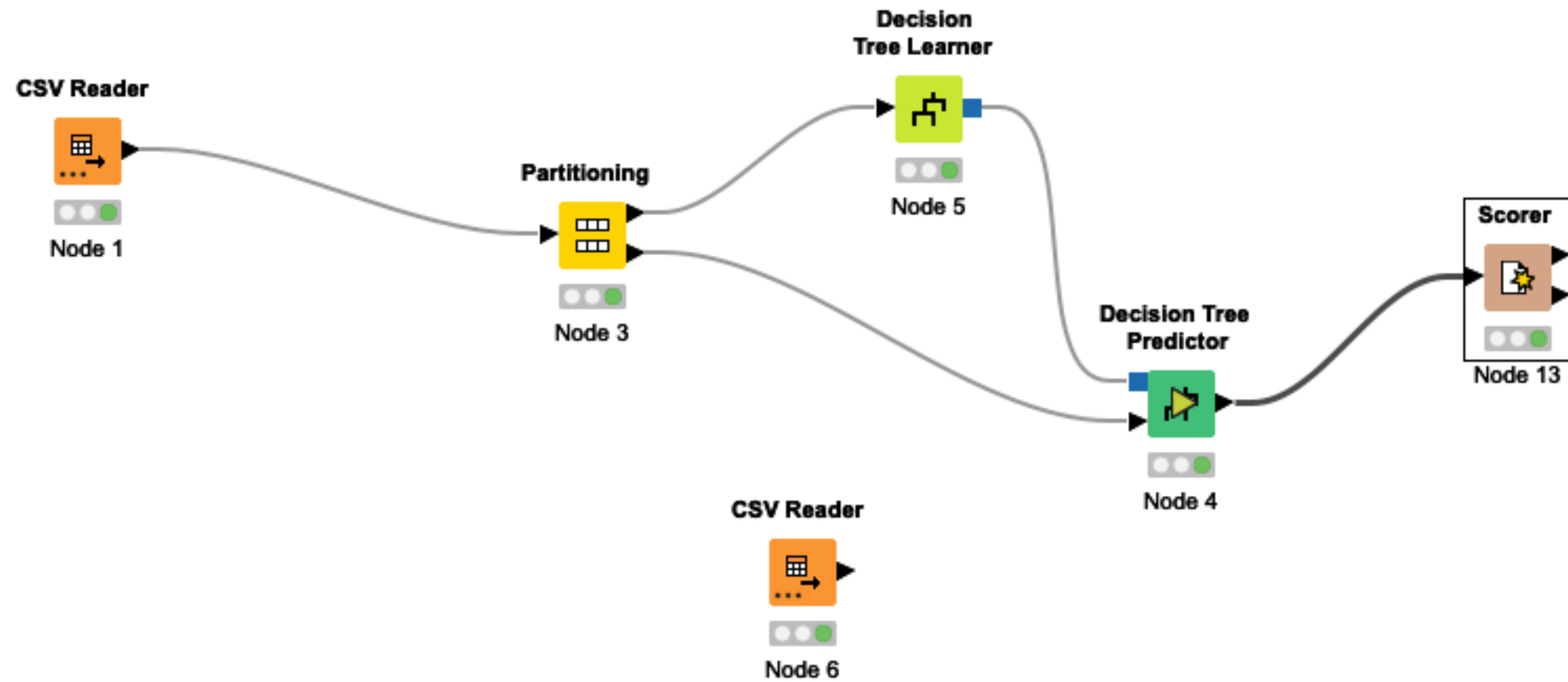
# Predicting Occupancy

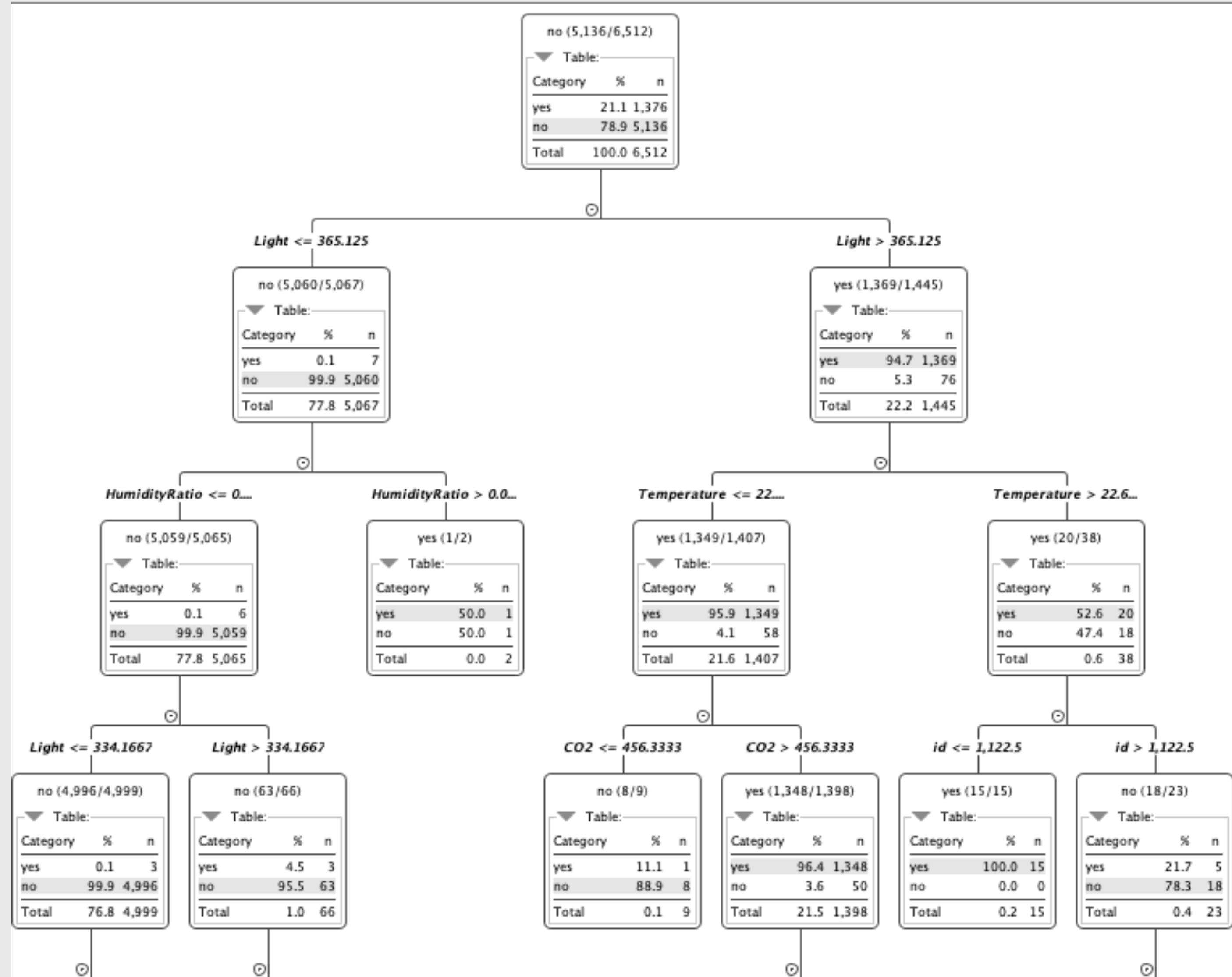
I will be building models to predict the occupancy of the room. I will be using a train set, which is the dataset that I had to clean, and a test set. The train set is the data that was used to train/teach the model – if only one dataset is provided to work with, the train set typically consumes 70-80% of the dataset.

The test set is what is used to test the accuracy of the model created by the train set. If the model is highly accurate, then it can be used again with new data. Due to my inability to figure out the nuances of KNIME and narrowing down the dataset down to specific columns, I manually had to drop columns from the dataset using `df.drop(['Light', 'CO2', 'HumidityRatio'], axis = 1)` as an example.



# Decision Tree One: Temperature and Humidity

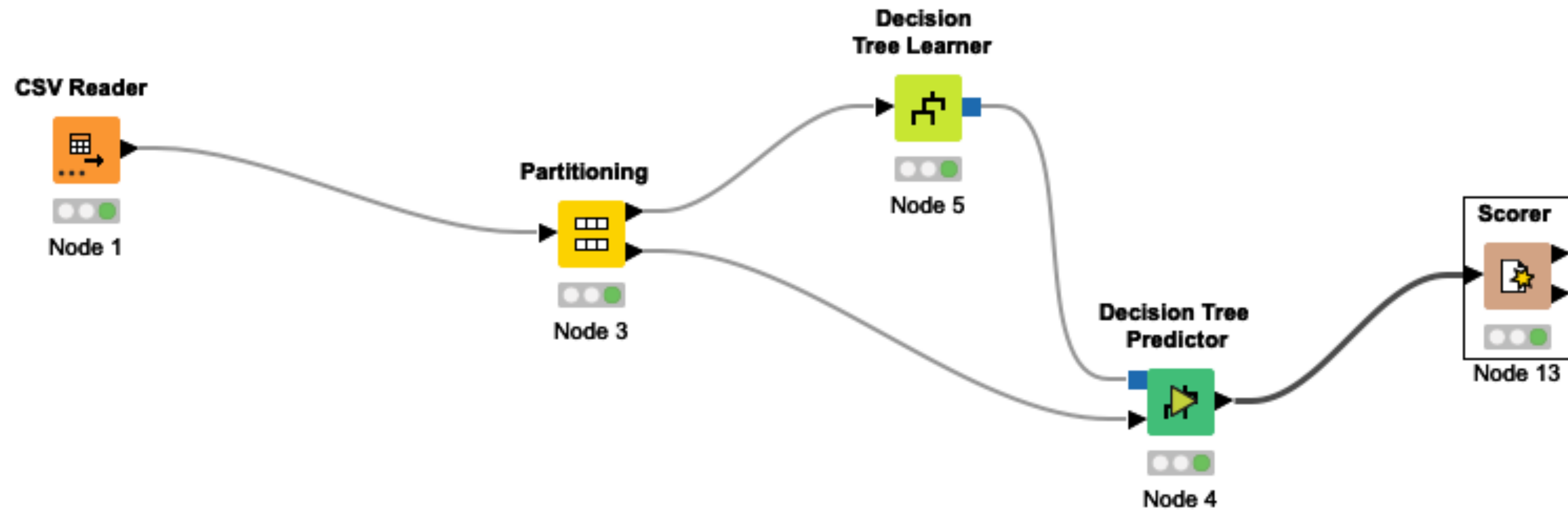


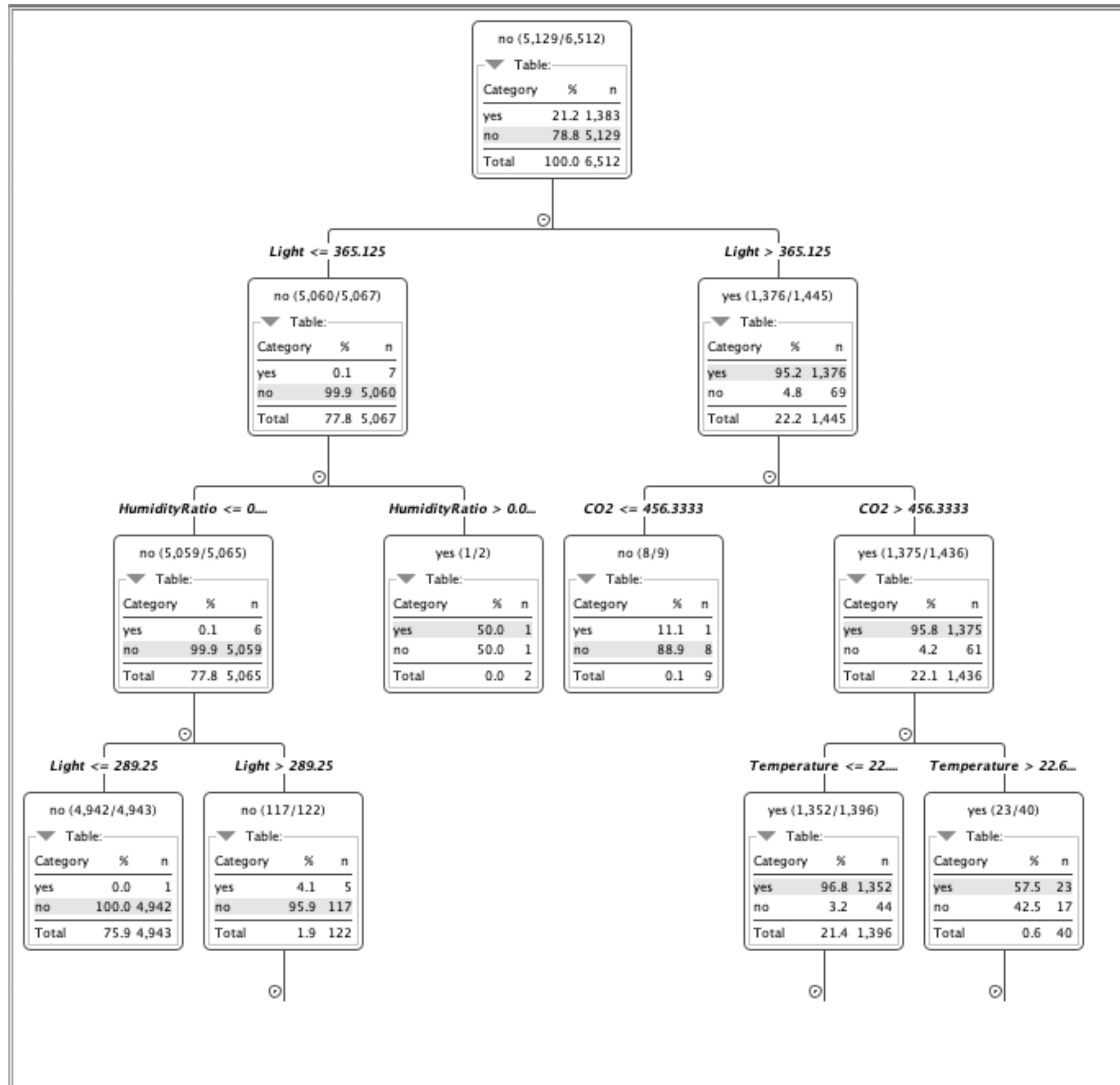


# Accuracy

Accuracy statistics - 4:13 - Scorer													
File Edit Hilite Navigation View													
Table "default" - Rows: 3 Spec - Columns: 11 Properties Flow Variables													
Row ID	I TrueP...	I FalseP...	I TrueN...	I False...	D Recall	D Precisi...	D Sensiti...	D Specifi...	D F-me...	D Accur...	D Cohen...		
yes	350	2	1274	2	0.994	0.994	0.994	0.998	0.994	?	?		
no	1274	2	350	2	0.998	0.998	0.998	0.994	0.998	?	?		
Overall	?	?	?	?	?	?	?	?	?	0.998	0.993		

## Decision Tree Two: Carbon Dioxide and Humidity Ratio

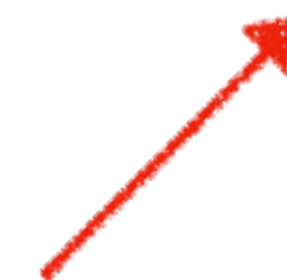




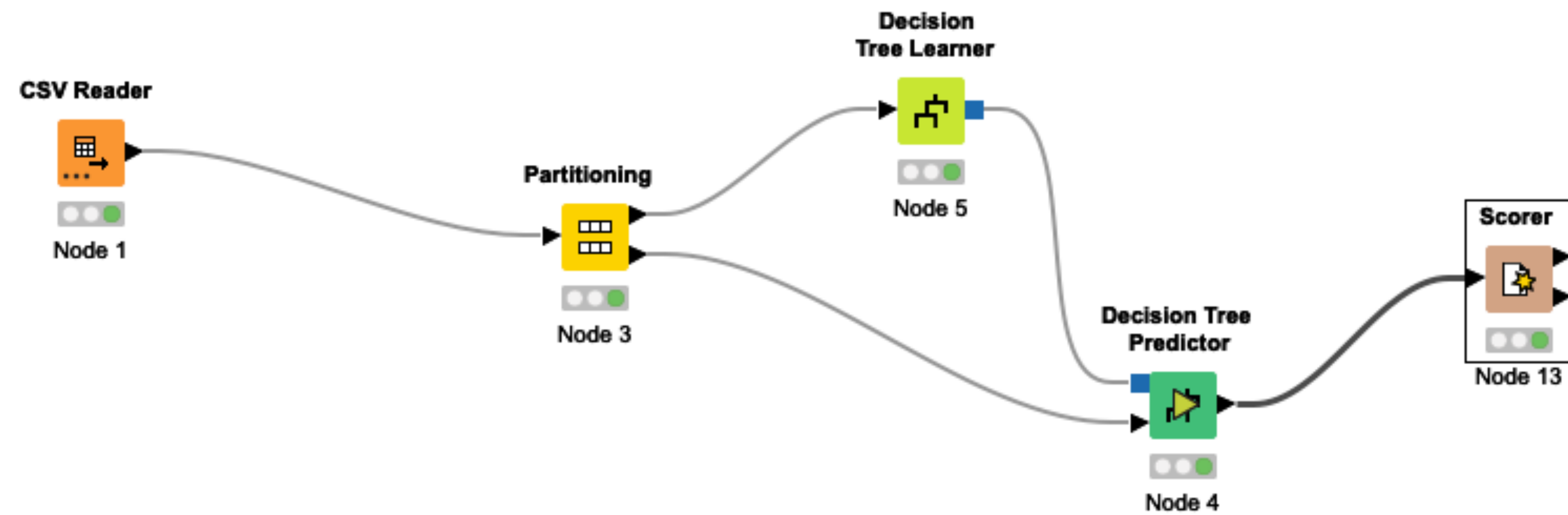


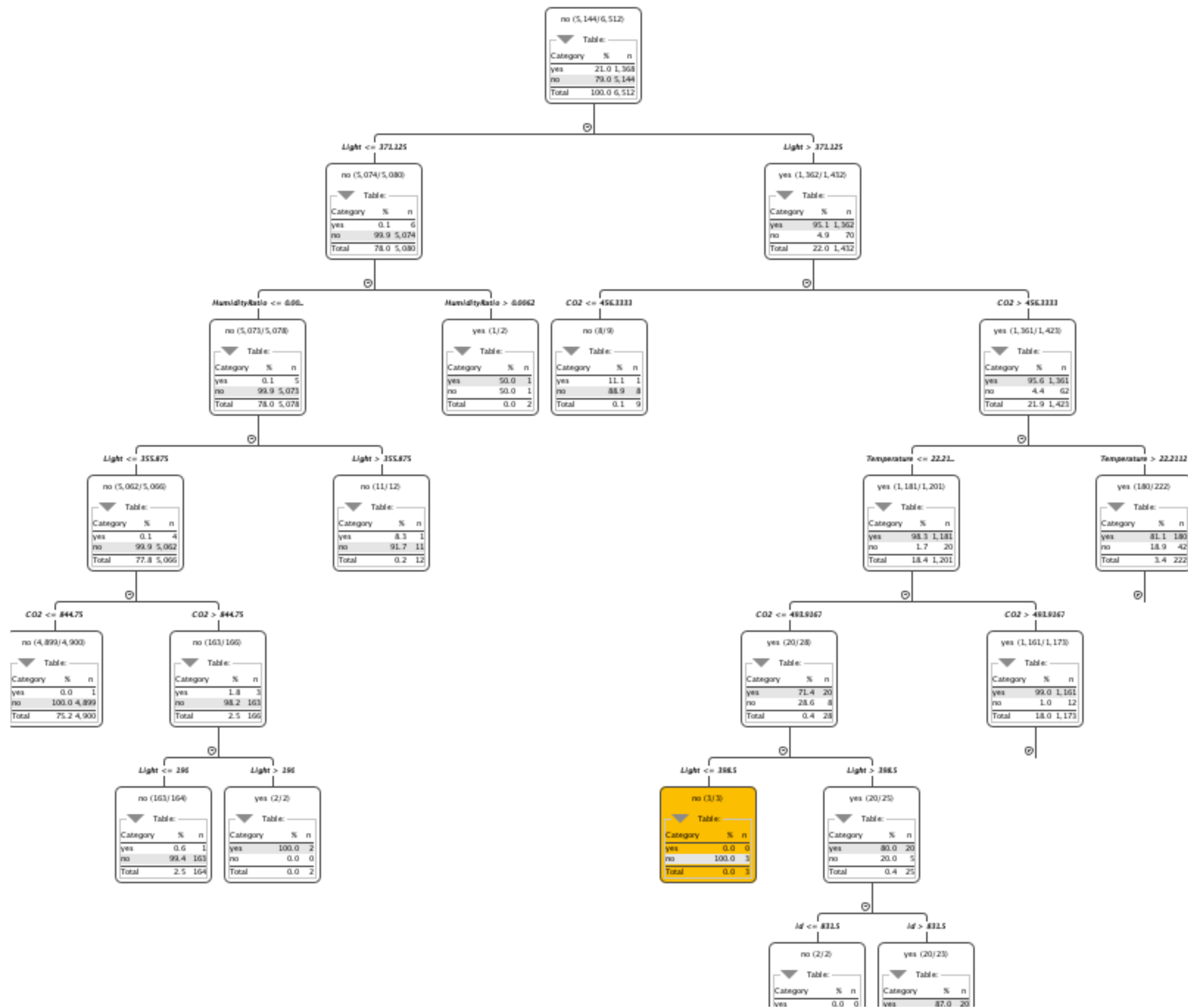
# Accuracy

Accuracy statistics - 0:13 - Scorer													
File Edit Hilite Navigation View													
Table "default" - Rows: 3 Spec - Columns: 11 Properties Flow Variables													
Row ID	I TrueP...	I FalseP...	I TrueN...	I False...	D Recall	D Precisi...	D Sensiti...	D Specifi...	D F-me...	D Accur...	D Cohen..		
yes	342	6	1277	3	0.991	0.983	0.991	0.995	0.987	?	?		
no	1277	3	342	6	0.995	0.998	0.995	0.991	0.996	?	?		
Overall	?	?	?	?	?	?	?	?	?	0.994	0.984		



## Decision Tree Three: Humidity, Light, and Carbon Dioxide



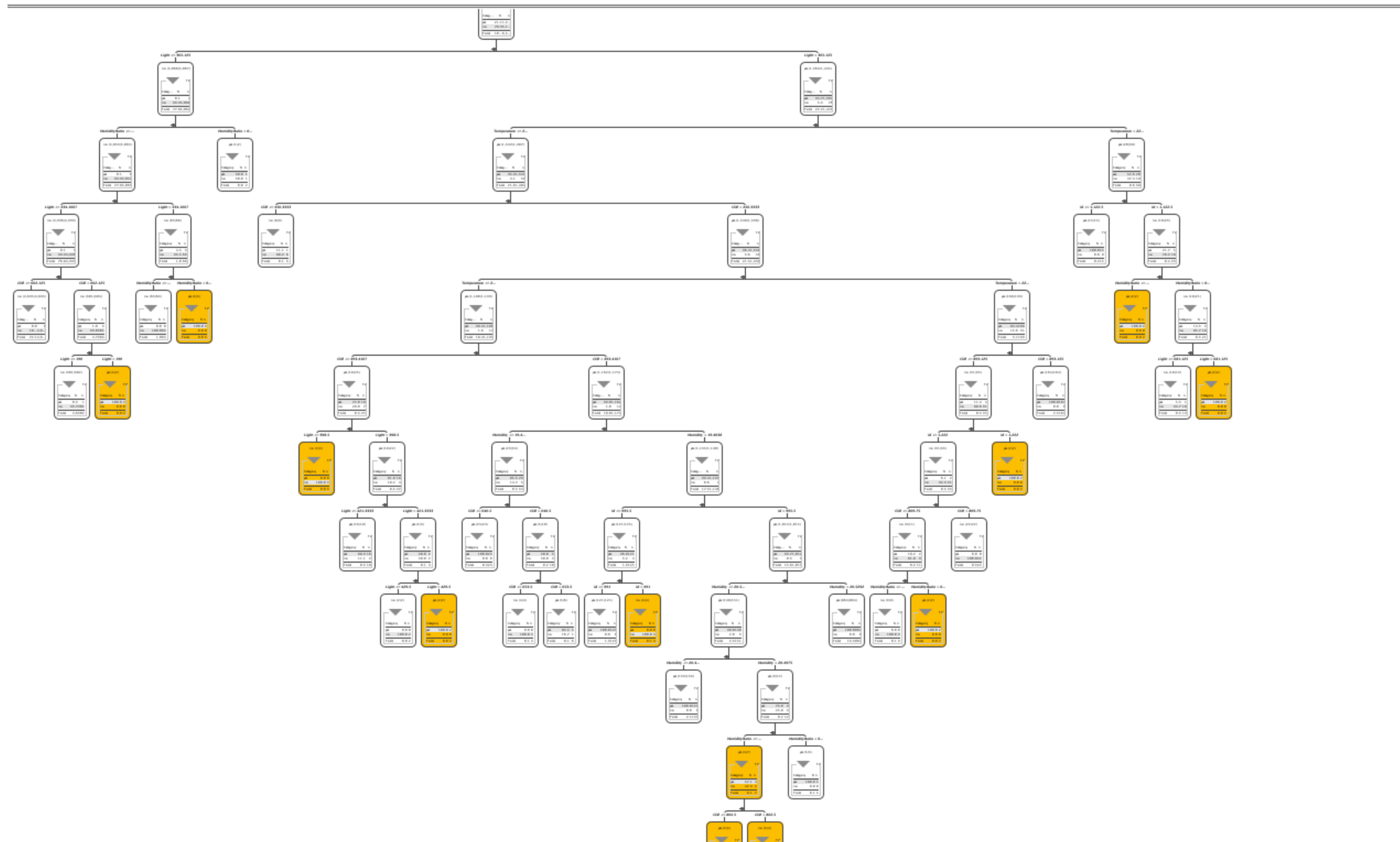


# Accuracy

The screenshot shows the Scorer application window titled "Accuracy statistics - 0:13 - Scorer". The interface includes a menu bar (File, Edit, Hilite, Navigation, View) and a toolbar with tabs for "Table 'default' - Rows: 3", "Spec - Columns: 11", "Properties", and "Flow Variables".

Row ID	I TrueP...	I FalseP...	I TrueN...	I False...	D Recall	D Precisi...	D Sensiti...	D Specifi...	D F-me...	D Accur...	D Cohen...
yes	354	5	1263	6	0.983	0.986	0.983	0.996	0.985	?	?
no	1263	6	354	5	0.996	0.995	0.996	0.983	0.996	?	?
Overall	?	?	?	?	?	?	?	?	?	0.993	0.98

A red arrow points from the bottom right towards the "Overall" row of the table.



Of these three decision trees, I would have to go with temperature and humidity, as it had the highest accuracy.

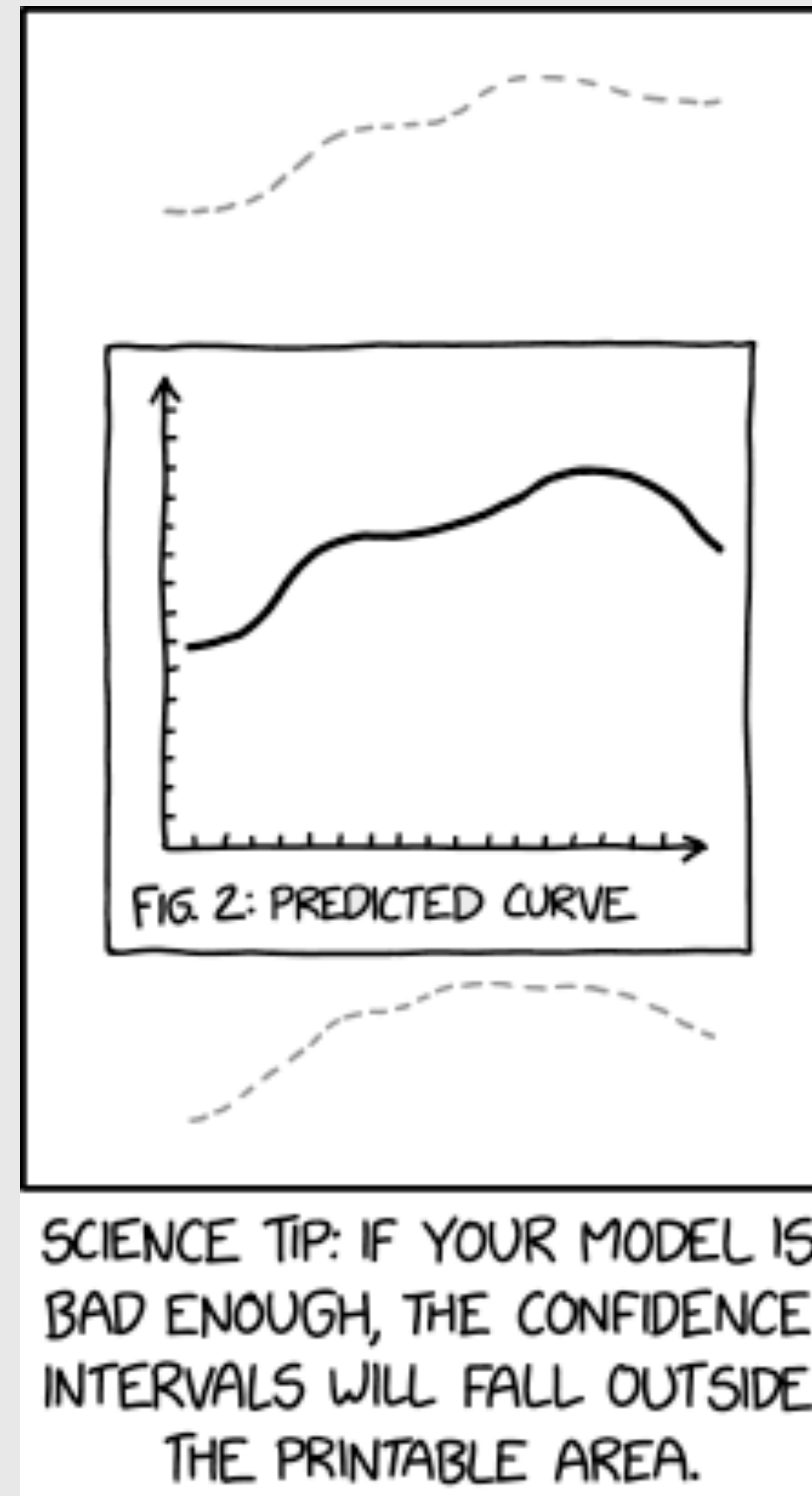


# Conclusion

I was extremely surprised that the model using temperature and humidity had the highest accuracy, as I had hypothesized to myself that carbon dioxide and humidity ratio would have been higher.

In this project, I cleaned, recleaned, and took the dataset to the dry cleaners. I learned how little I know about KNIME, and that I still have a lot to learn in the Classification area.





<https://xkcd.com/2311/>