

Weekly Progress Report

Agrotech Live :: Wiggle Labs

2025-04-16

Abstract

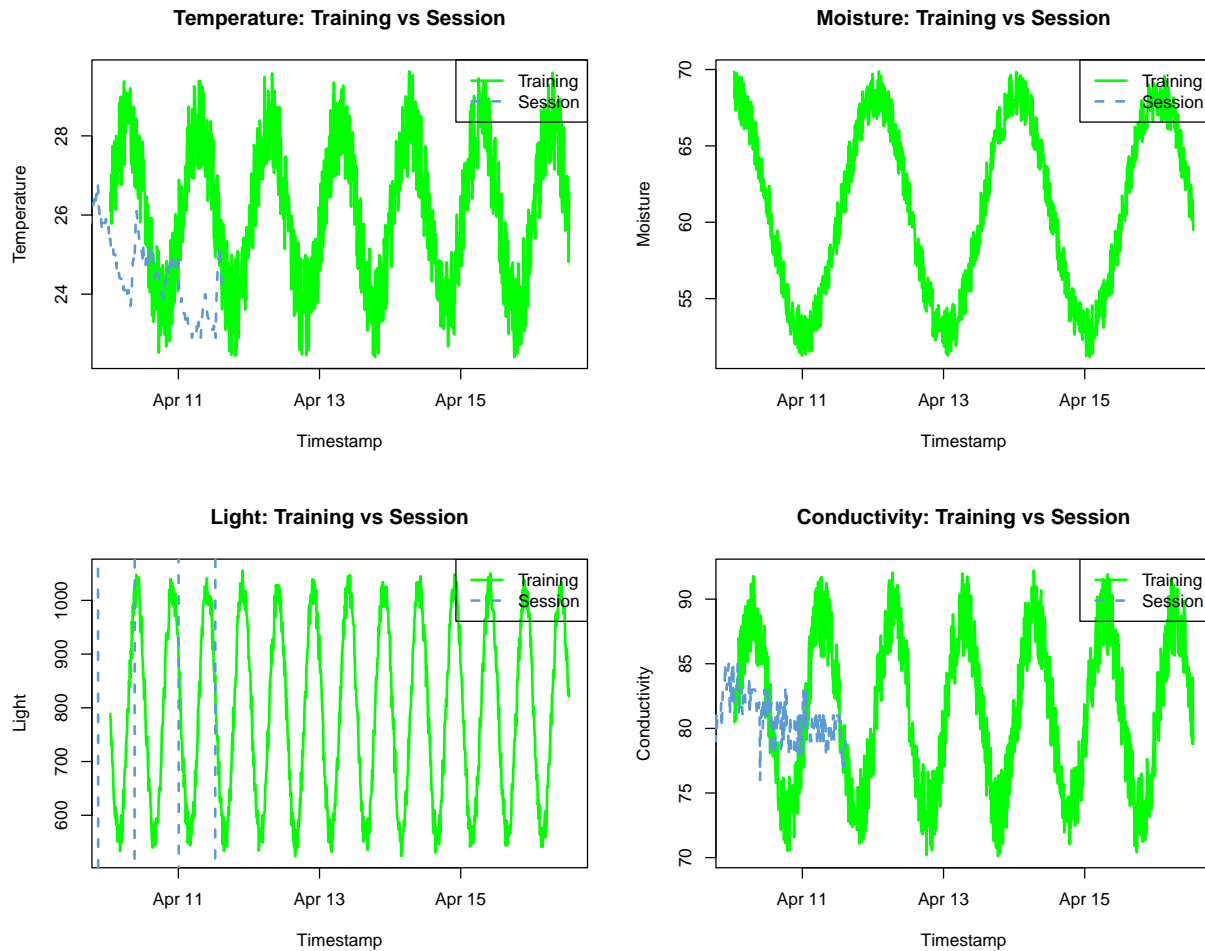
In an ongoing effort to achieve food sovereignty, Wiggle Labs has developed Agrotech Live to monitor the soil health of different plants and crops. This tool collects four data points; Temperature, Moisture, Light and Conductivity from sensors placed near a subject. The training data for this program are the ideal conditions for the subject, and the performance of the experiment is based on how close the collected sensor (testing) data is to the input care training data. Input data is identical in structure the testing data (collected during a session), except it represents only perfect conditions for the subject. A score is generated (along with other statistical results) periodically to communicate how well the experiment is performing.

This report examines the relationships between the session and training data features over the past 7 days, and forecasts the next 3 days. Clustering will be used to prepare the data for classification, enabling the communication of health quality thresholds (e.g., ideal, good, OK, fair, poor).

Stage 1: Time Series, Correlation, Covariance and Heatmap

Time Series Comparison

This section shows the session and training data features on top of one another to understand how they compare over time.



Correlation and Covariance Matrices

In this correlation matrix, score of 1.0 or -1.0 represents a perfect (positive or negative) self-correlation and values closer to 0 show less to no correlation. The matrix above reflects the same relationships as the covariance matrix above.

CROSS-COVARIANCE MATRIX (Session vs Training)

##	Temperature	Moisture	Light	Conductivity
## Temperature	0.9615548	1.4337150	9.004801	3.377057
## Moisture	0.6396479	-3.0793977	-51.078258	1.875694
## Light	-206.5048534	2101.2970940	83063.171110	-435.946003
## Conductivity	0.8880845	-0.1147215	-59.493709	3.097508

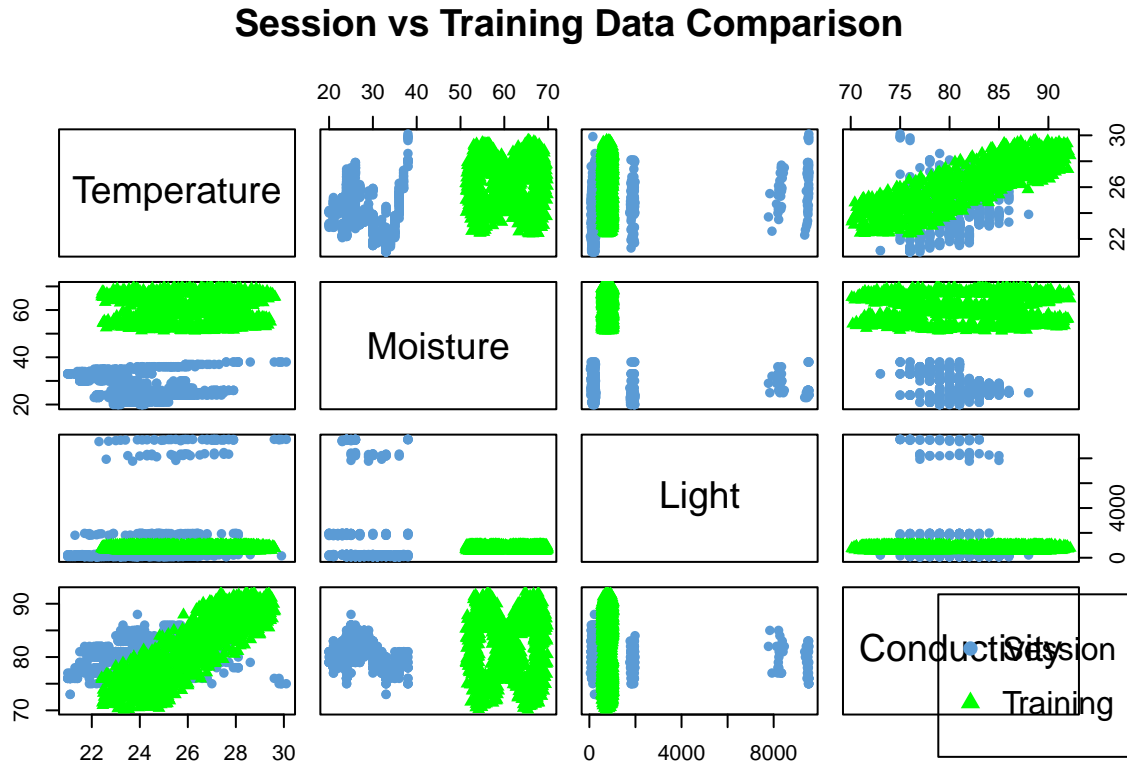
##

CROSS-CORRELATION MATRIX (Session vs Training)

##	Temperature	Moisture	Light	Conductivity
## Temperature	0.32783027	0.158592364	0.03300803	0.35027163
## Moisture	0.07647516	-0.119450951	-0.06565774	0.06822333

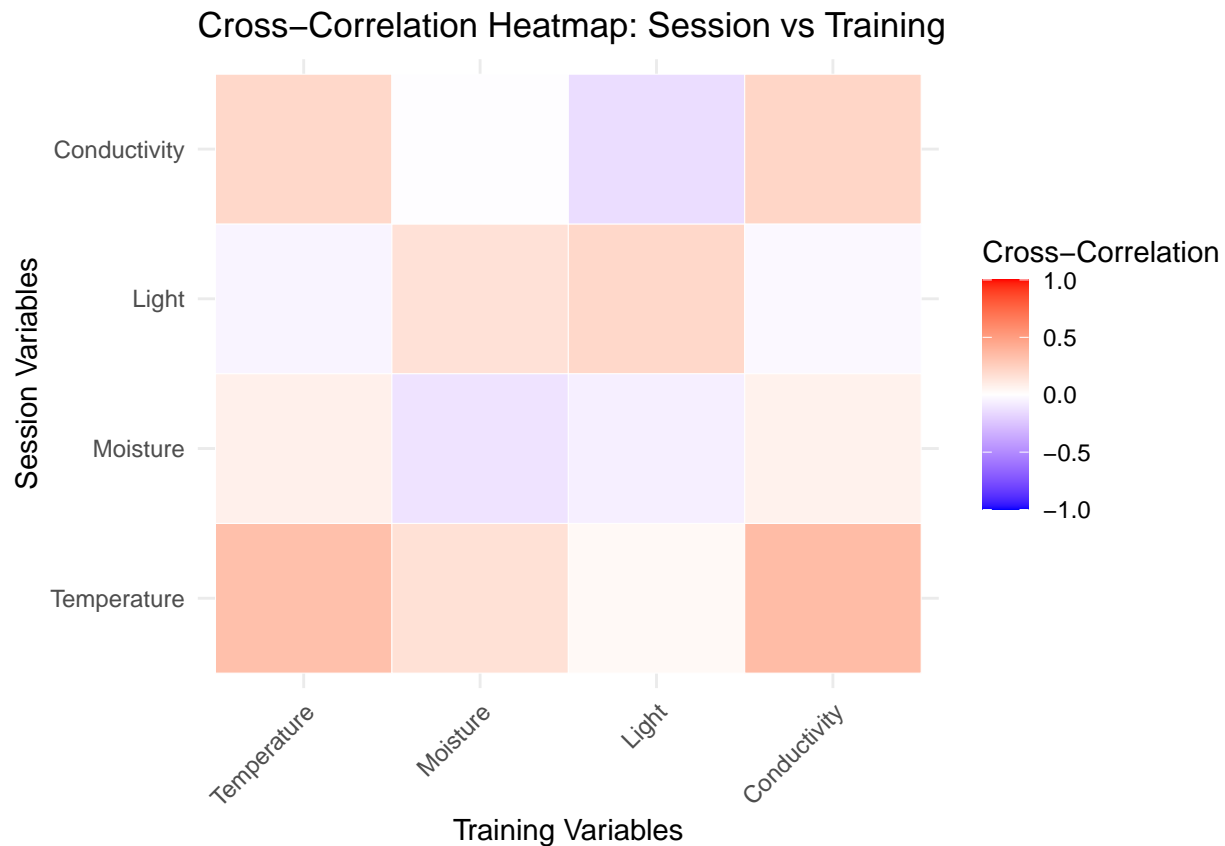
```
## Light      -0.04712967  0.155595144  0.20381826  -0.03026833
## Conductivity 0.20344965 -0.008526914 -0.14653593  0.21587714
```

Plotting variables against each other.



Up until this point (excluding the matrecies) we have been treating the session and training data as independent data sets for comparison purposes. Instead of running similarity algorithms within each variable, both the independent (`session_data`) and the dependent (`training_data`) must be used together in order to produce valuable insights about their relationship.

Cross Correlation Heat Map



Stage 2: Similarities and RSME

In this stage the goal is to measure the similarities and differences between our session and training data. From there, we'll be able to classify and label the features for K means clustering in Stage 3.

Preprocessing

```
library(data.table)
```

```
##  
## Attaching package: 'data.table'  
  
## The following objects are masked from 'package:reshape2':  
##  
## dcast, melt  
  
## The following objects are masked from 'package:dplyr':  
##  
## between, first, last
```

```
library(dplyr)

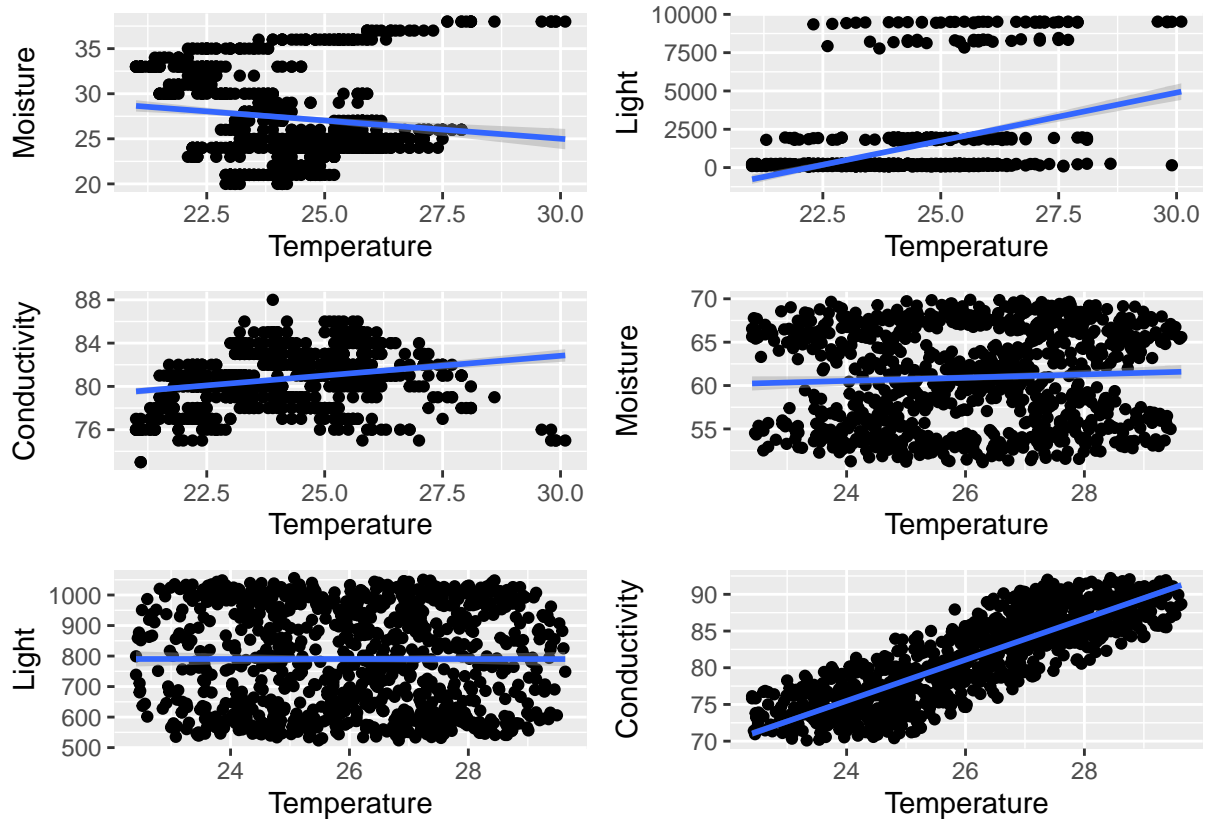
training_dt <- as.data.table(training_data)
session_dt <- as.data.table(session_data)
# head()
```

Cleaning the Data

To clean up the data, we'll do some indexing by setting up keys. Then to organize it, we put the fields we want to use for our analysis into their own data frames. Converting `sq_ft_lot` and `sale_price` to numeric will ensure no problems when plotting.

Linear Regression Model & Errors

```
## 'geom_smooth()' using formula = 'y ~ x'
## 'geom_smooth()' using formula = 'y ~ x'
## 'geom_smooth()' using formula = 'y ~ x'
## 'geom_smooth()' using formula = 'y ~ x'
## 'geom_smooth()' using formula = 'y ~ x'
## 'geom_smooth()' using formula = 'y ~ x'
```



This plot is nice, but can only tell us so much about the data. Running a `summary()` will give us more information.

Session Summaries

AGT session data

```
##
## Please cite as:

## Hlavac, Marek (2022). stargazer: Well-Formatted Regression and Summary Statistics Tables.

## R package version 5.2.3. https://CRAN.R-project.org/package=stargazer
```

```
##
## Temperature Model: Simple vs Full
## =====
##                               Dependent variable:
##                               -----
##                               Temperature
##                               (1)                (2)
## -----
## Temperature                0.359***          0.423***
##                           (0.034)            (0.039)
##
## Moisture                   0.068***
##                           (0.013)
##
## Light                      -0.0001***
##                           (0.00003)
##
## Conductivity              0.107***
##                           (0.026)
##
## Constant                  17.496***          5.619***
##                           (0.816)            (2.176)
## -----
## Observations                937              937
## R2                        0.107              0.178
## Adjusted R2               0.107              0.174
## Residual Std. Error    1.693 (df = 935)      1.628 (df = 932)
## F Statistic           112.587*** (df = 1; 935) 50.323*** (df = 4; 932)
## =====
## Note:                      *p<0.1; **p<0.05; ***p<0.01
```

```
##
## Moisture Model: Simple vs Full
## =====
##                               Dependent variable:
##                               -----
##                               Moisture          Temperature
##                               (1)                (2)
## -----
## Temperature                0.423***
##                           (0.039)
```

```

##
## Moisture          -0.141***          0.068***
##                  (0.038)            (0.013)
##
## Light              -0.0001***
##                  (0.00003)
##
## Conductivity       0.107***
##                  (0.026)
##
## Constant          64.805***          5.619***
##                  (1.067)            (2.176)
##
## -----
## Observations      937                937
## R2                 0.014              0.178
## Adjusted R2        0.013              0.174
## Residual Std. Error 5.484 (df = 935)    1.628 (df = 932)
## F Statistic       13.534*** (df = 1; 935) 50.323*** (df = 4; 932)
## =====
## Note:                *p<0.1; **p<0.05; ***p<0.01

```

```

##
## Light Model: Simple vs Full

```

```

## =====
##                               Dependent variable:
##                               -----
##                               Light          Temperature
##                               (1)           (2)
##                               -----
## Temperature                0.423***
##                               (0.039)
##
## Moisture                    0.068***
##                               (0.013)
##
## Light                       0.014***          -0.0001***
##                               (0.002)            (0.00003)
##
## Conductivity                0.107***
##                               (0.026)
##
## Constant                    773.949***          5.619***
##                               (5.916)            (2.176)
##
## -----
## Observations      937                937
## R2                 0.042              0.178
## Adjusted R2        0.041              0.174
## Residual Std. Error 163.190 (df = 935)    1.628 (df = 932)
## F Statistic       40.525*** (df = 1; 935) 50.323*** (df = 4; 932)
## =====
## Note:                *p<0.1; **p<0.05; ***p<0.01

```

```
##
## Conductivity Model: Simple vs Full
## =====
##                               Dependent variable:
##                               -----
##                               Conductivity      Temperature
##                               (1)              (2)
##                               -----
## Temperature                                0.423***
##                                           (0.039)
##
## Moisture                                0.068***
##                                           (0.013)
##
## Light                                -0.0001***
##                                           (0.00003)
##
## Conductivity          0.522***          0.107***
##                      (0.077)          (0.026)
##
## Constant              39.385***          5.619***
##                      (6.227)          (2.176)
##
## -----
## Observations              937              937
## R2                        0.047              0.178
## Adjusted R2              0.046              0.174
## Residual Std. Error    5.752 (df = 935)      1.628 (df = 932)
## F Statistic            45.704*** (df = 1; 935) 50.323*** (df = 4; 932)
## =====
## Note:                               *p<0.1; **p<0.05; ***p<0.01
```

Residual Standard Error

Residuals represent the differences between the actual and predicted values of our dependent (response) variable, `sq_ft_lot`. A high RSE indicates a weak model for this prediction.

Multiple R

Multiple R, also called the correlation coefficient, measures the strength and direction of a relationship between variables. On a scale of -1 to +1, values closer to -1 or +1 represent perfect negative or positive correlation. 0 means no correlation at all.

Multiple R2 Error

R squared tells us the proportion for variance in `sq_ft_lot` explained by `sale_price`, the predictor variable. It can be on a scale of 0 to 1. A value closer to 1 represents greater variance, while closer to 0 represents less variance. Returning a whole 0 or 1 means none or perfect variance respectively.

Adjusted R2 Error

This error is a modified version of the above R-squared error that is able to accommodate for multiple predictors in a regression model.

Stage 3: Modeling, Classification and Metrics

In this next section, we will use the KNN algorithm to find the best K for a label. As this is a placeholder for future use, a good question to use this for is:

How can we classify a subjects health by having thresholds of poor, unsatisfactory, neutral, satisfactory or excellent?

To measure this hypothetical threshold, we would use how closely or different the training data is from testing. This was explored in the previous stage, and is required to be able to have a label to predict for in KNN.

```
## Session Cluster Distribution (%):

##
##      1      2
## 8.32444 91.67556

##
## Training Cluster Distribution (%):

##
##      1      2
## 52.40128 47.59872

##
## Session Cluster Centers (scaled):

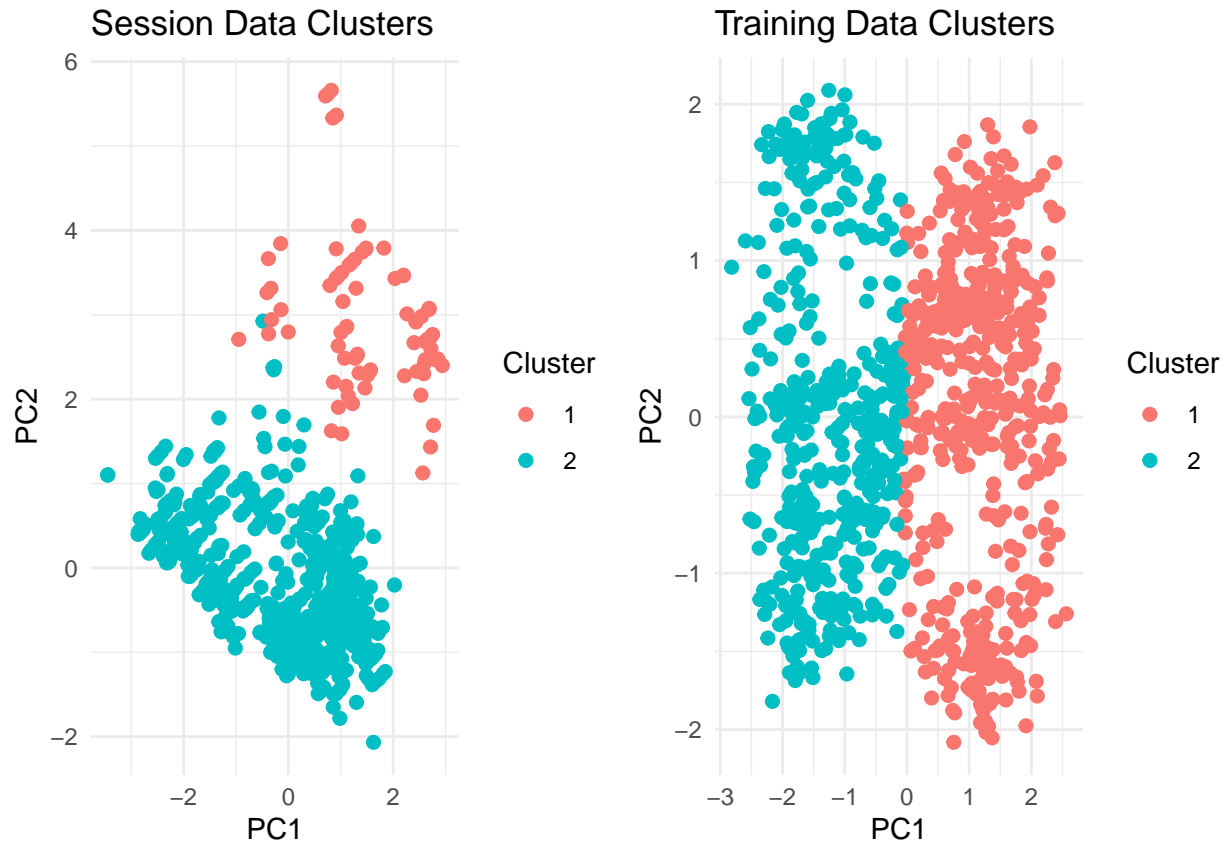
##   Temperature    Moisture      Light Conductivity
## 1    1.223590 -0.009060060  3.2021365  -0.50693761
## 2   -0.111106  0.000822683 -0.2907644   0.04603159

##
## Training Cluster Centers (scaled):

##   Temperature    Moisture      Light Conductivity
## 1    0.7941567  0.1221828  0.009808624   0.8176824
## 2   -0.8742846 -0.1345106 -0.010798283  -0.9001840

##
## Attaching package: 'gridExtra'

## The following object is masked from 'package:dplyr':
##
##      combine
```



In order to begin preprocessing and find the right k , we will only keep the relevant features and label the cluster. This analysis uses the session data as testing data and the generated ideal conditions as the training data.

We will use two approaches. First with the built in K nearest neighbor functions with R and then manually calculating the accuracy of each K values.

Method 1: Built in knn() Function

```
# Select numeric features and the cluster label
features_train <- training_data[, c("Temperature", "Moisture", "Light", "Conductivity")]
labels_train <- as.factor(training_data$Cluster)

test_features <- session_data[, c("Temperature", "Moisture", "Light", "Conductivity")]
labels_test <- as.factor(session_data$Cluster)
```

Normalizing the data will help scale larger values down to a comparable size.

```
# Normalize (scale) the features
features_train_scaled <- as.data.frame(scale(features_train))
test_features_scaled <- as.data.frame(scale(test_features))
```

For our train/test split,

```
train_features <- features_train_scaled
train_labels <- labels_train
```

Now we'll proceed with modeling KNN and evaluating its performance. Before we start though, this cross correlation data needs to be fitted a bit more in order to pass through the KNN model.

```
# Make sure train and test sets have the same number of rows as their respective labels
min_train_rows <- min(nrow(train_features), length(train_labels))
min_test_rows <- min(nrow(test_features), length(labels_test))

# Trim all data to match
train_features <- train_features[1:min_train_rows, ]
train_labels <- train_labels[1:min_train_rows]
test_features <- test_features[1:min_test_rows, ]
labels_test <- labels_test[1:min_test_rows]

# Convert labels to factors after trimming
train_labels <- as.factor(train_labels)
labels_test <- as.factor(labels_test)
```

```
library(class)

# Training the KNN model
k_value <- 3 # You can experiment with different values of k
knn_model <- knn(train_features, test_features, train_labels, k = k_value)
# Evaluate the model
confusion_matrix <- table(Predicted = knn_model, Actual = labels_test)
accuracy <- sum(diag(confusion_matrix)) / sum(confusion_matrix)

cat("Accuracy on Session Data (Test Data):", accuracy)
```

```
## Accuracy on Session Data (Test Data): 0.0832444
```

Method 2: Manual K Value Accuracy Computation

K Nearest Neighbor (KNN) works by computing the Euclidean distance between the test and training points. Then after selecting the proper k that is the shortest distance, it assigns those closest neighbors most common label to the test point.

Using `knn()` the `class` package automatically computes the Euclidean Distance between two points. We will adjust in the input parameters to ingest the train and test data we have already prepared.

```
# Inputs for manual model
train_data <- features_train_scaled
test_data <- test_features_scaled
train_labels <- as.numeric(labels_train)
test_labels <- as.numeric(labels_test)
```

```
euc_dis <- function(p1, p2) {
  sqrt(sum((p1 - p2)^2))
}
```

In this next section we're implementing the KNN Classifier manually to train the binary classifier data. At this stage, the classifier logic is being defined below.

Here is where we'll compute accuracy for the session data using manual KNN.

```
k_values <- seq(1, 15, 2) # or however many k's you want

cat("Train samples:", nrow(train_data), "Label count:", length(train_labels), "\n")
```

```
## Train samples: 937 Label count: 937
```

```
accuracy_results <- c()

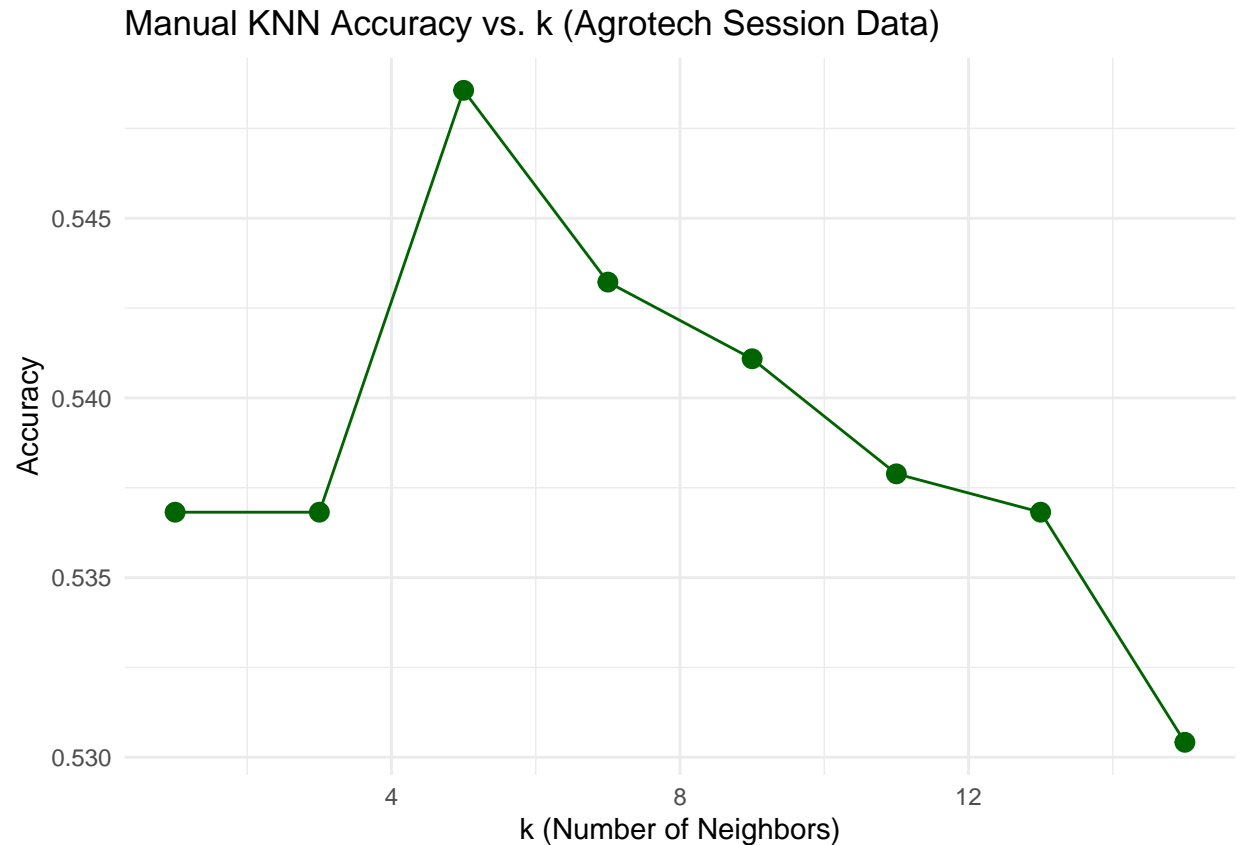
for (k in k_values) {
  predictions <- knn(train = train_data, test = test_data, cl = train_labels, k = k)
  acc <- mean(predictions == test_labels)
  accuracy_results <- c(accuracy_results, acc)
  cat("k =", k, ", Accuracy =", round(acc * 100, 2), "%\n")
}
```

```
## k = 1 , Accuracy = 53.68 %
## k = 3 , Accuracy = 53.68 %
## k = 5 , Accuracy = 54.86 %
## k = 7 , Accuracy = 54.32 %
## k = 9 , Accuracy = 54.11 %
## k = 11 , Accuracy = 53.79 %
## k = 13 , Accuracy = 53.68 %
## k = 15 , Accuracy = 53.04 %
```

The accuracy of this model varies based in the k value used. Plotting them will provide a better idea of which to use to get the best accuracy.

```
accuracy_data <- data.frame(
  k_values = k_values,
  accuracy = accuracy_results,
  dataset = "Session/Training"
)

ggplot(accuracy_data, aes(x = k_values, y = accuracy)) +
  geom_line(color = "darkgreen") +
  geom_point(size = 3, color = "darkgreen") +
  labs(title = "Manual KNN Accuracy vs. k (Agrotech Session Data)",
       x = "k (Number of Neighbors)", y = "Accuracy") +
  theme_minimal()
```



Stage 4: Forecasting

Time Series Analysis and Forecasting

Number of unique days in dataset: 8

Date range: 20182 to 20189

Average observations per day: 117.12

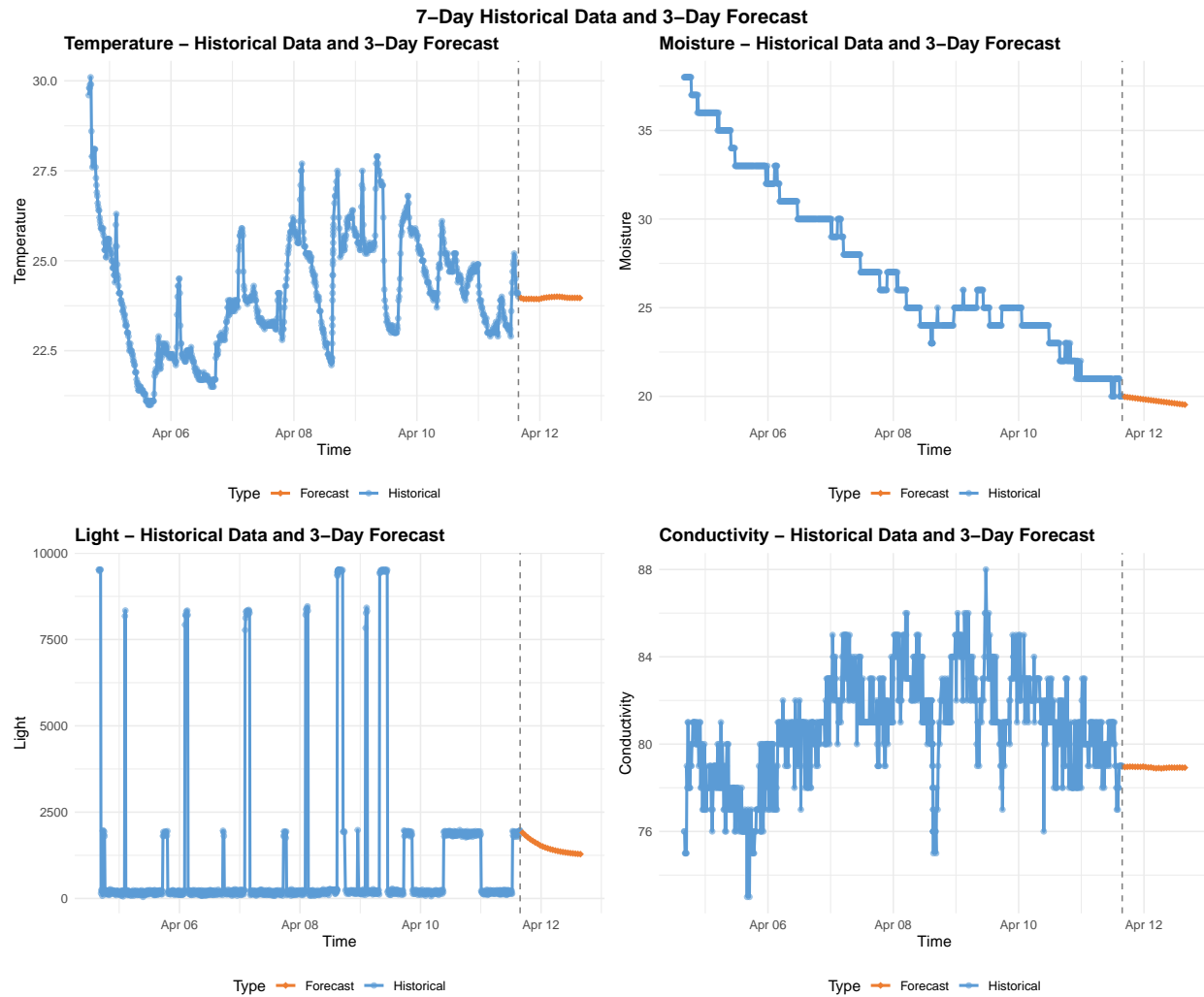
Creating Forecast Models

For each sensor variable, we'll use an appropriate time series forecasting method. We'll evaluate ARIMA, ETS (Exponential Smoothing), and Prophet models to find the best approach for our data.

Visualization of Historical Data and Forecasts

```
## Warning: 'aes_string()' was deprecated in ggplot2 3.0.0.  
## i Please use tidy evaluation idioms with 'aes()'.  
## i See also 'vignette("ggplot2-in-packages")' for more information.  
## This warning is displayed once every 8 hours.  
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was  
## generated.
```

```
## Warning: Using 'size' aesthetic for lines was deprecated in ggplot2 3.4.0.
## i Please use 'linewidth' instead.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.
```



Forecast Accuracy and Confidence Intervals

```
##
## Temperature Forecast Accuracy Metrics:
##           ME      RMSE      MAE      MPE      MAPE      MASE
## Training set -0.003616977 0.1969759 0.1301407 -0.01469972 0.5355604 0.1164202
##           ACF1
## Training set 0.006182006

##
## Moisture Forecast Accuracy Metrics:
##           ME      RMSE      MAE      MPE      MAPE      MASE
## Training set 4.556916e-05 0.2540073 0.1063922 -0.001173416 0.4101601 0.1646374
```

```

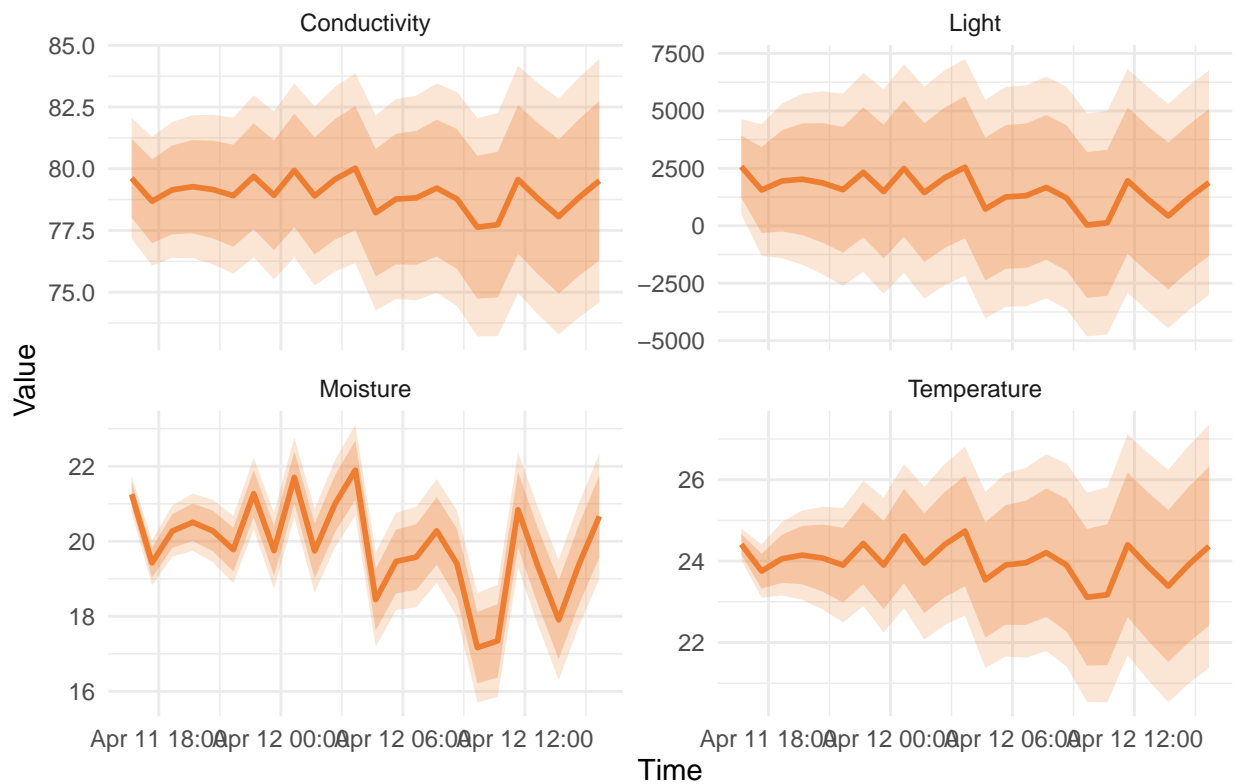
##                               ACF1
## Training set -0.0117104

##
## Light Forecast Accuracy Metrics:
##                               ME    RMSE    MAE    MPE    MAPE    MASE    ACF1
## Training set -10.1262 1052.6 313.7139 -95.74347 106.4918 0.181823 -0.008651945

##
## Conductivity Forecast Accuracy Metrics:
##                               ME    RMSE    MAE    MPE    MAPE    MASE
## Training set 0.008766796 1.244855 0.9799116 -0.008150268 1.215079 0.5272005
##                               ACF1
## Training set -0.006767934

```

3-Day Forecasts with 80% and 95% Confidence Intervals



Forecast Table Summary

Table 1: Daily Forecast Summary for the Next 3 Days

Day	Variable	Min	Mean	Max
1	Temperature	23.93	23.97	24.00
2	Temperature	NA	NA	NA
3	Temperature	NA	NA	NA

Day	Variable	Min	Mean	Max
1	Moisture	19.53	19.75	19.97
2	Moisture	NA	NA	NA
3	Moisture	NA	NA	NA
1	Light	1284.85	1480.36	1902.90
2	Light	NA	NA	NA
3	Light	NA	NA	NA
1	Conductivity	78.89	78.93	78.97
2	Conductivity	NA	NA	NA
3	Conductivity	NA	NA	NA

Comparison to Training (Ideal) Data

```
## Temperature      Moisture      Light Conductivity
##      26.13298      60.93523      790.28120      81.46064
```

```
##      Variable Ideal_Mean Ideal_Min Ideal_Max
## 1 Temperature      26.13298      NA      NA
## 2 Moisture      60.93523      NA      NA
## 3 Light      790.28120      NA      NA
## 4 Conductivity      81.46064      NA      NA
```

```
## Ideal Conditions (from Training Data):
```

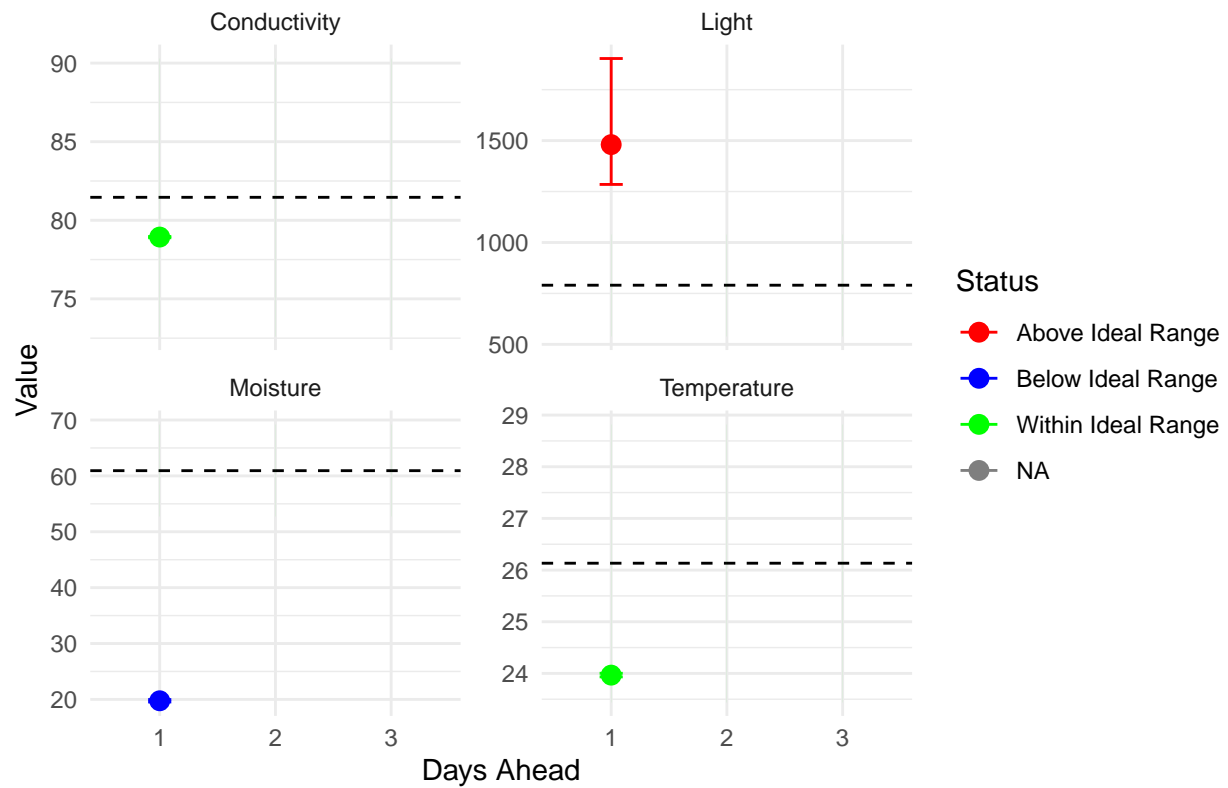
```
##      Variable Ideal_Mean Ideal_Min Ideal_Max
## 1 Temperature      26.13298      23.44616      28.81979
## 2 Moisture      60.93523      52.65403      69.21643
## 3 Light      790.28120      540.38063      1040.18176
## 4 Conductivity      81.46064      72.62890      90.29238
```

Table 2: Forecast Comparison to Ideal Conditions

Variable	Day	Mean	Ideal_Mean	Status
Conductivity	1	78.93	81.46	Within Ideal Range
Conductivity	2	NA	81.46	NA
Conductivity	3	NA	81.46	NA
Light	1	1480.36	790.28	Above Ideal Range
Light	2	NA	790.28	NA
Light	3	NA	790.28	NA
Moisture	1	19.75	60.94	Below Ideal Range
Moisture	2	NA	60.94	NA
Moisture	3	NA	60.94	NA
Temperature	1	23.97	26.13	Within Ideal Range
Temperature	2	NA	26.13	NA
Temperature	3	NA	26.13	NA

```
## Warning: Removed 8 rows containing missing values or values outside the scale range
## ('geom_point()').
```


3-Day Forecast vs. Ideal Growing Conditions



Stage 5: Discussion