# Rural Municipality Crop Yield Prediction

Ethan Lam

# Agenda

Introduction

The Problem

Solution

# The problems

I want to find the best place in Saskatchewan to grow Barley

I want to know if my investment is good in terms of the ROI or I want to forecast the yields in a specific Rural Municipality of a crop
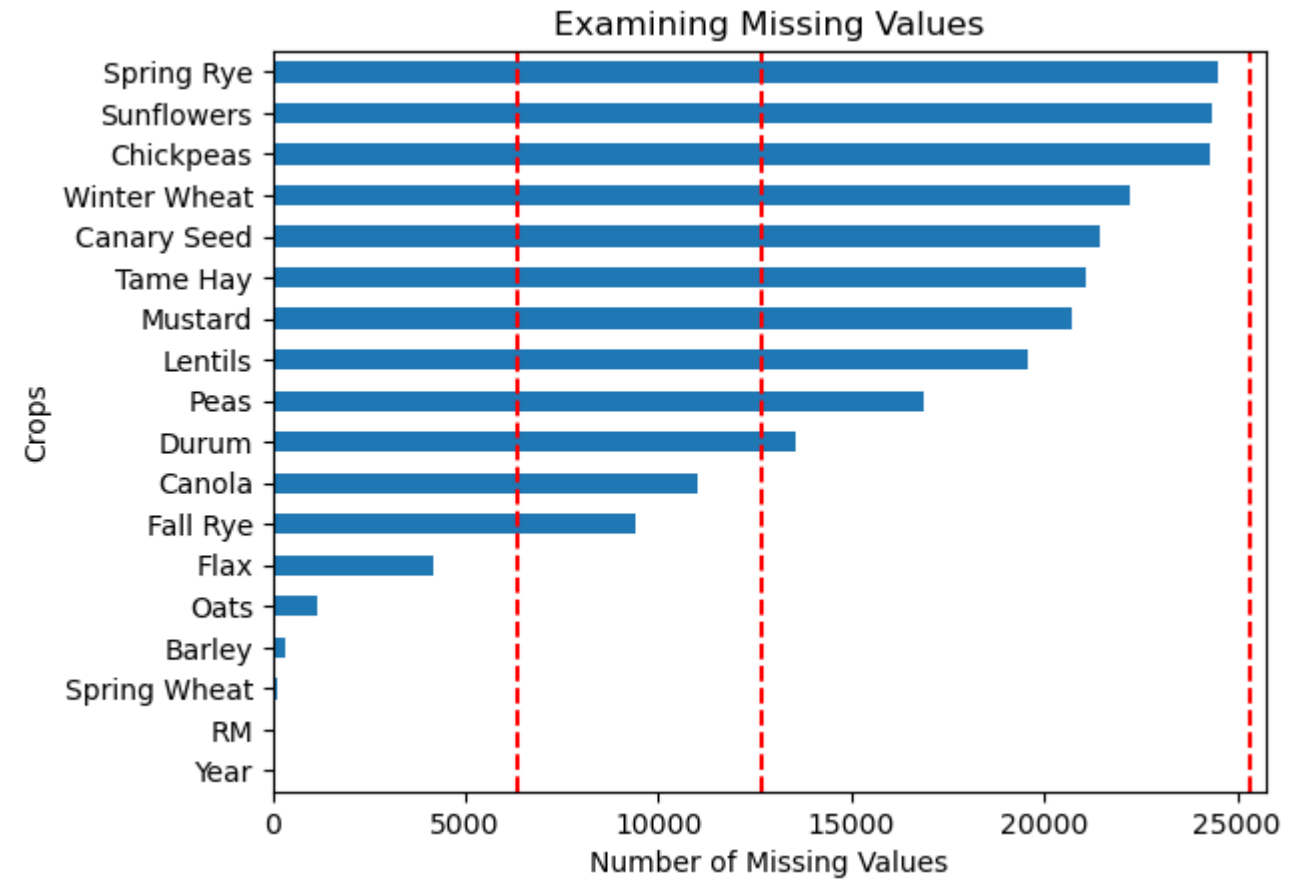
# Data Collection and Preprocessing

- Crop yield data:
  - Crop yields by Rural Municipality (RM) are produced annually from the Ministry of Saskatchewan Crop Report and Saskatchewan Crop Insurance Corporation
  - Data provided from 1938 to 2022.

- Geospatial:
  The shapefile from Government of Saskatchewan

- First look at the data

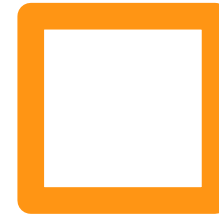| Year | RM | Winter Wheat | Canola | Spring Wheat | Mustard | Durum | Sunflowers | Oats | Lentils | Peas | Barley | Fall Rye | Canary Seed | Spring Rye | Tame Hay | Flax | Chickpeas |
|------|----|--------------|--------|--------------|---------|-------|------------|------|---------|------|--------|----------|-------------|-----------|----------|------|-----------|
| 1938 | 1 | NaN | NaN | 4 | NaN | NaN | NaN | 1 | NaN | NaN | 1 | NaN | NaN | NaN | NaN | 0 | NaN |
| 1939 | 1 | NaN | NaN | 9 | NaN | NaN | NaN | 16 | NaN | NaN | 16 | NaN | NaN | NaN | NaN | 0 | NaN |
| 1940 | 1 | NaN | NaN | 12 | NaN | NaN | NaN | 23 | NaN | NaN | 19 | NaN | NaN | NaN | NaN | 8 | NaN |
| 1941 | 1 | NaN | NaN | 18 | NaN | NaN | NaN | 32 | NaN | NaN | 28 | NaN | NaN | NaN | NaN | 5 | NaN |
| 1942 | 1 | NaN | NaN | 20 | NaN | NaN | NaN | 35 | NaN | NaN | 28 | 14 | NaN | NaN | NaN | 5 | NaN |

# Data Quality Check

- Check for NULL/Missing values
  - Many missing values
- Check for duplicate
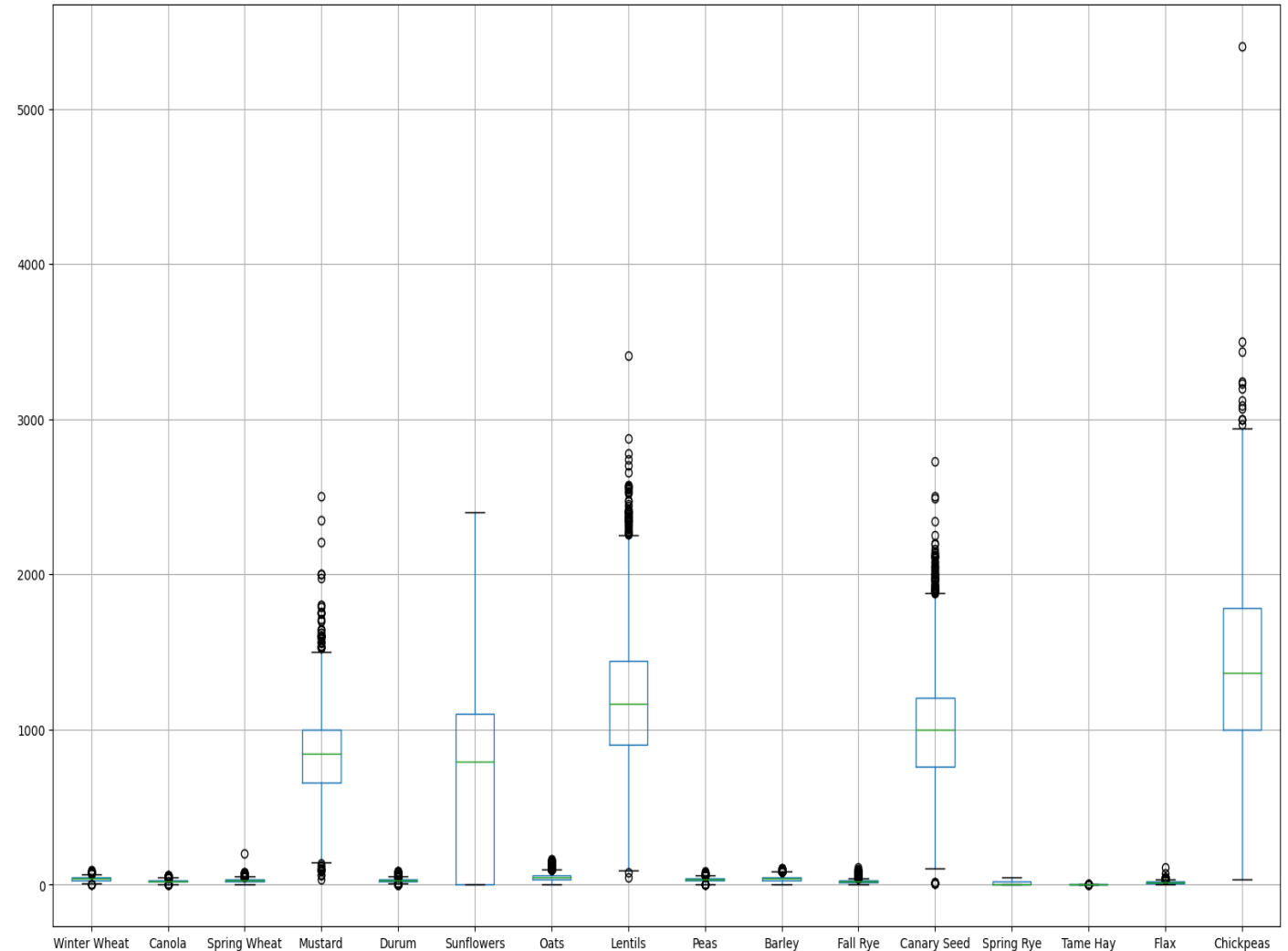  - No duplicate was found



Examining Missing Values

# EXPLORATORY DATA ANALYSIS

- Understand that there are 25312 entries in the dataset

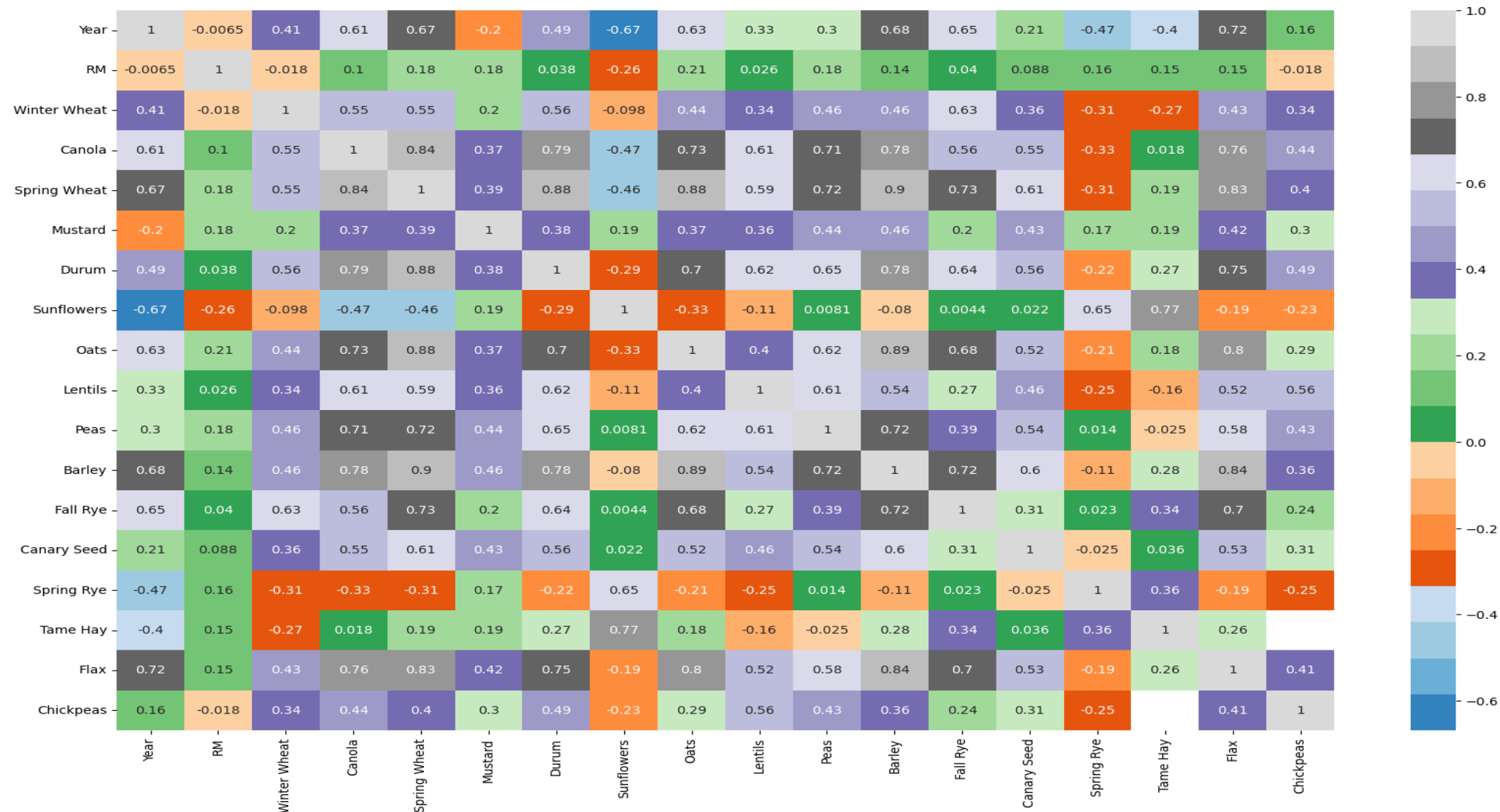- There are 18 columns in the dataset

- There are 299 RMs

# Check Outliers

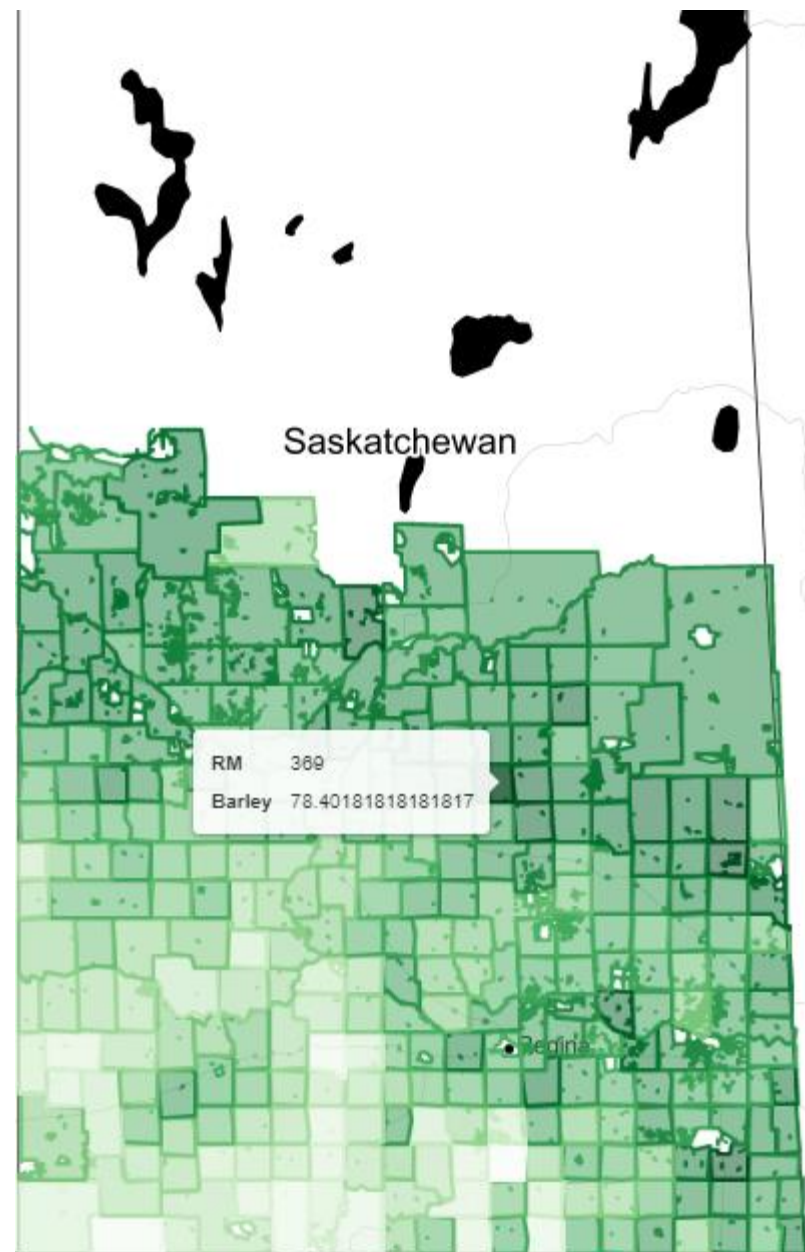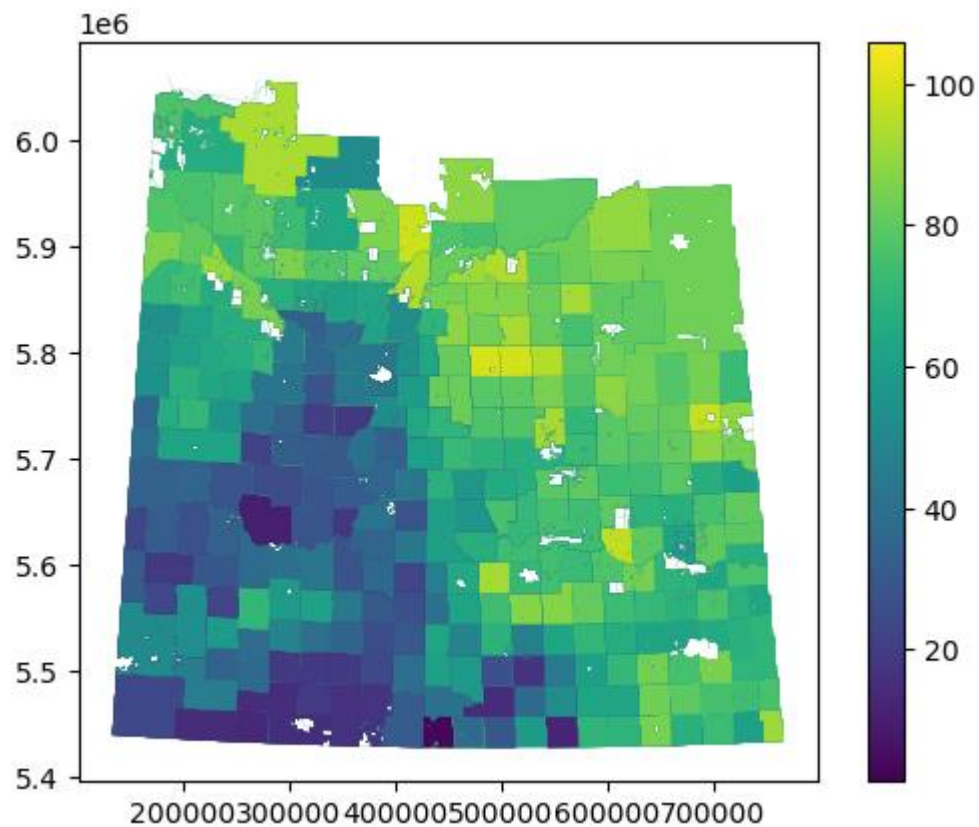There are many outliers in the data of Mustard, Lentils, Canary Seeds, and Chickpeas
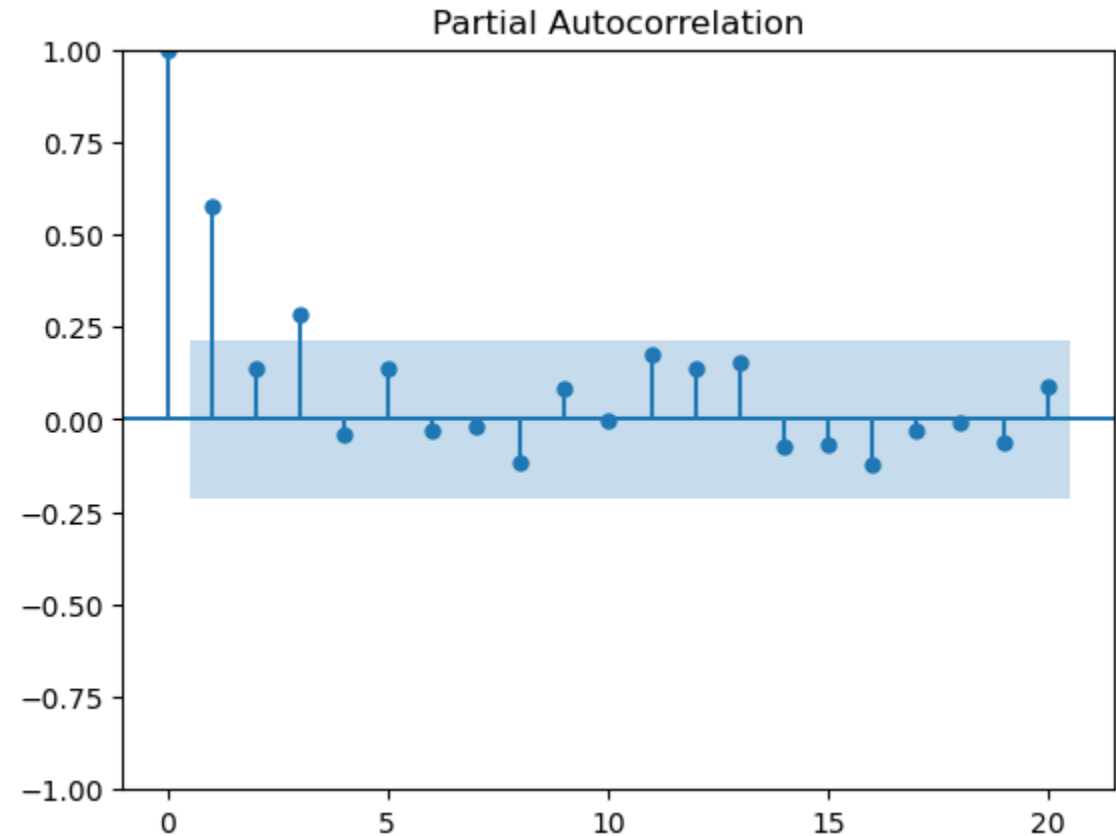
# Correlation Matrix

# GIS Analysis

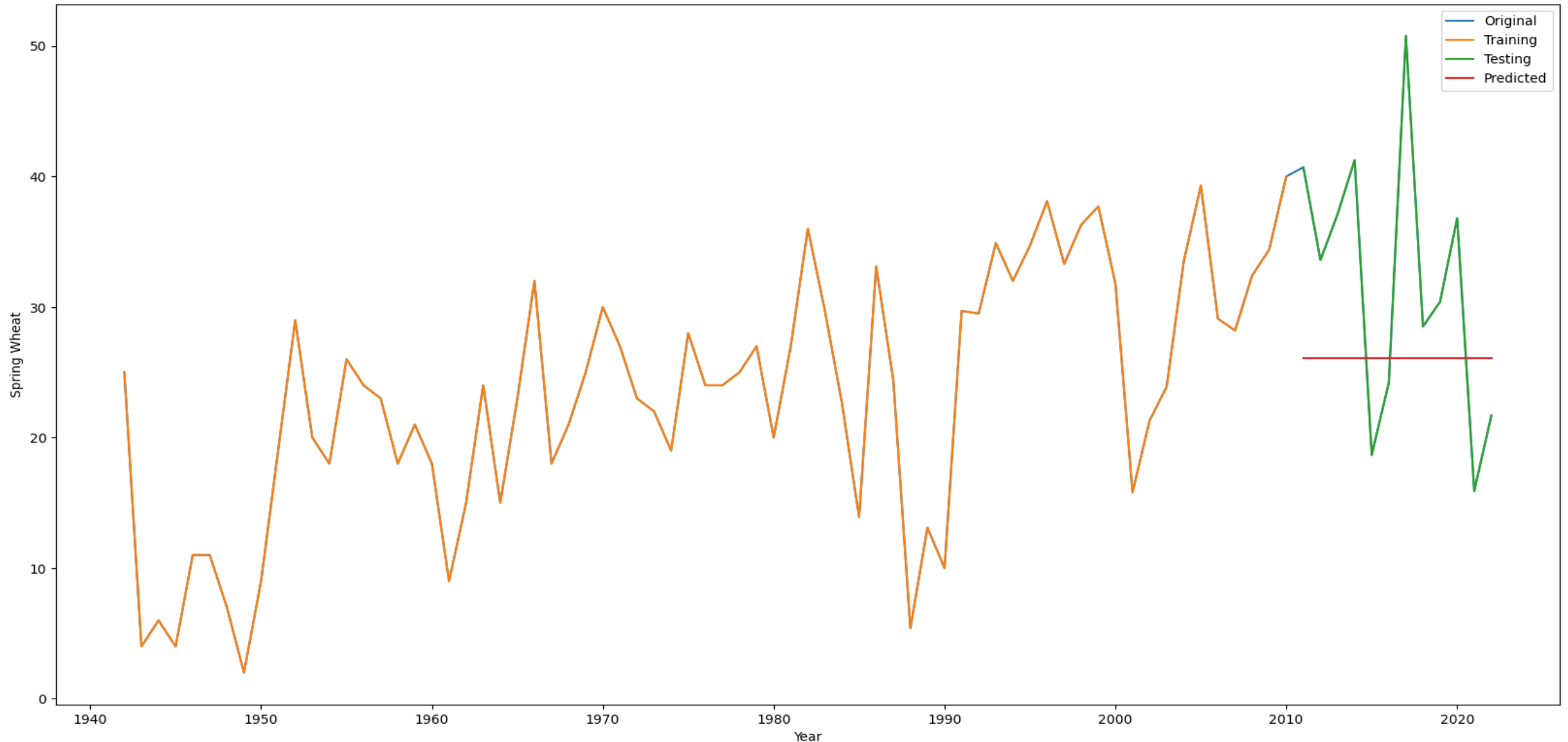The place has the highest average Barley yield in the last 10 year is RM: 369

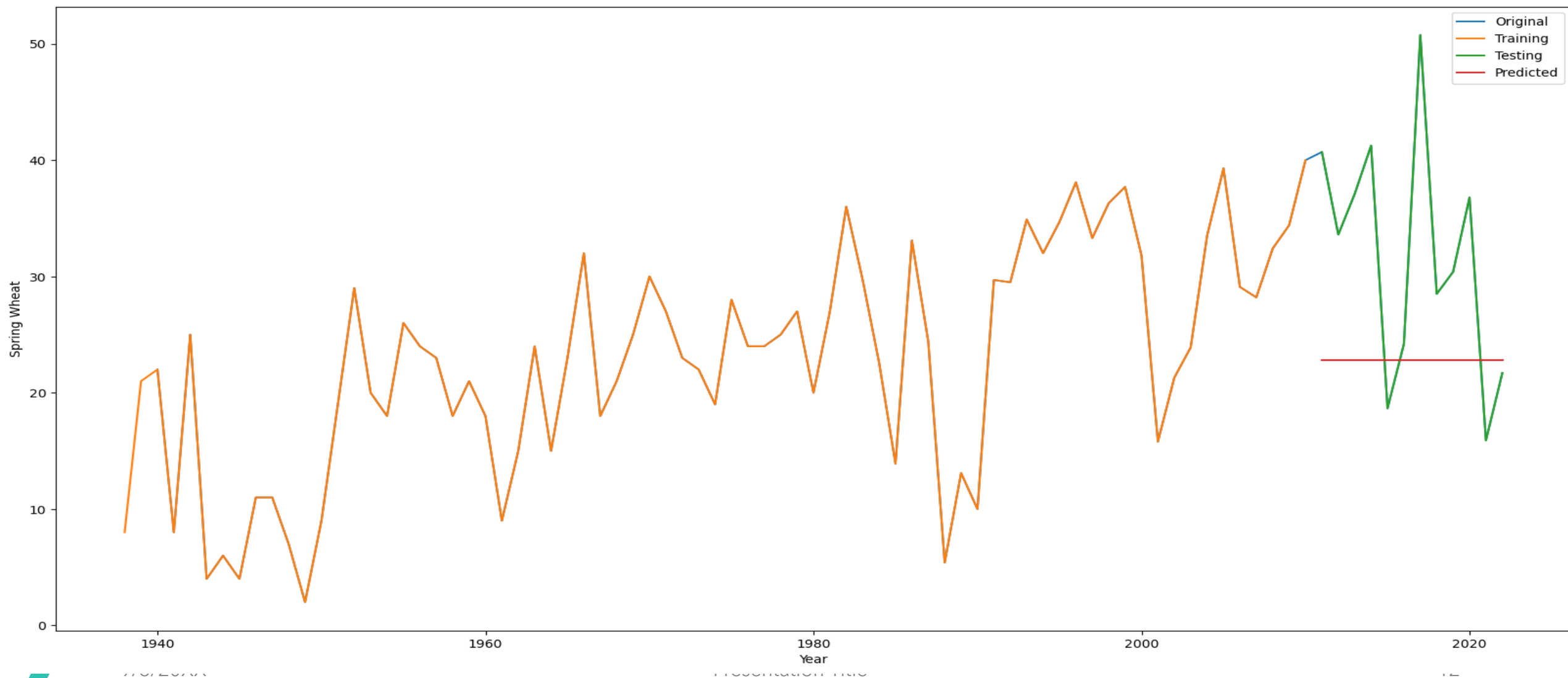# Time Series: Partial Autocorrelation (Spring Wheat)

- When analyzing the plot, we can see that the first lag has a very strong correlation to our future value.

- Lag 5 is the last lag the clearly goes above the green threshold line. As such, we now know to use 5 lags to create our auto regression model
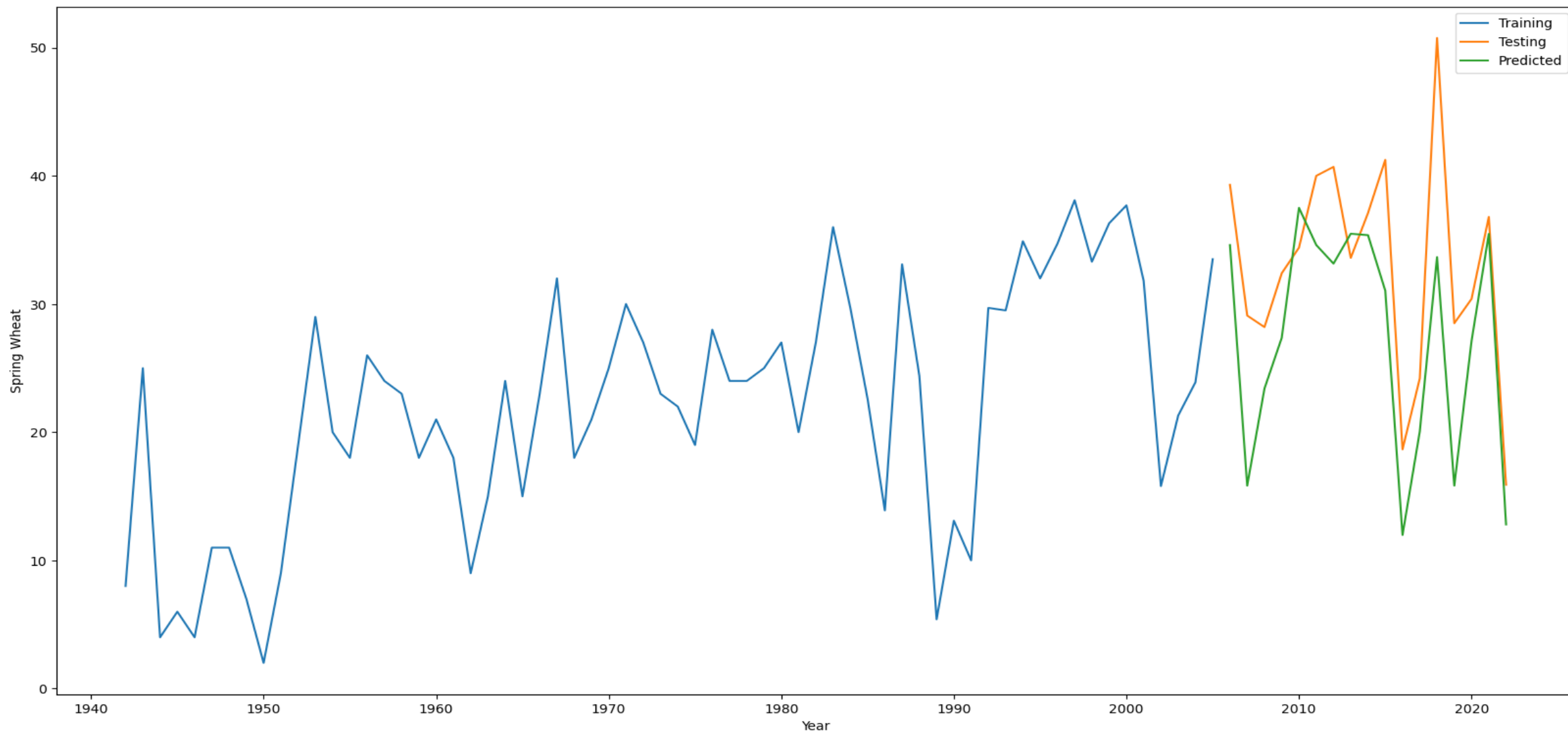


Partial Autocorrelation

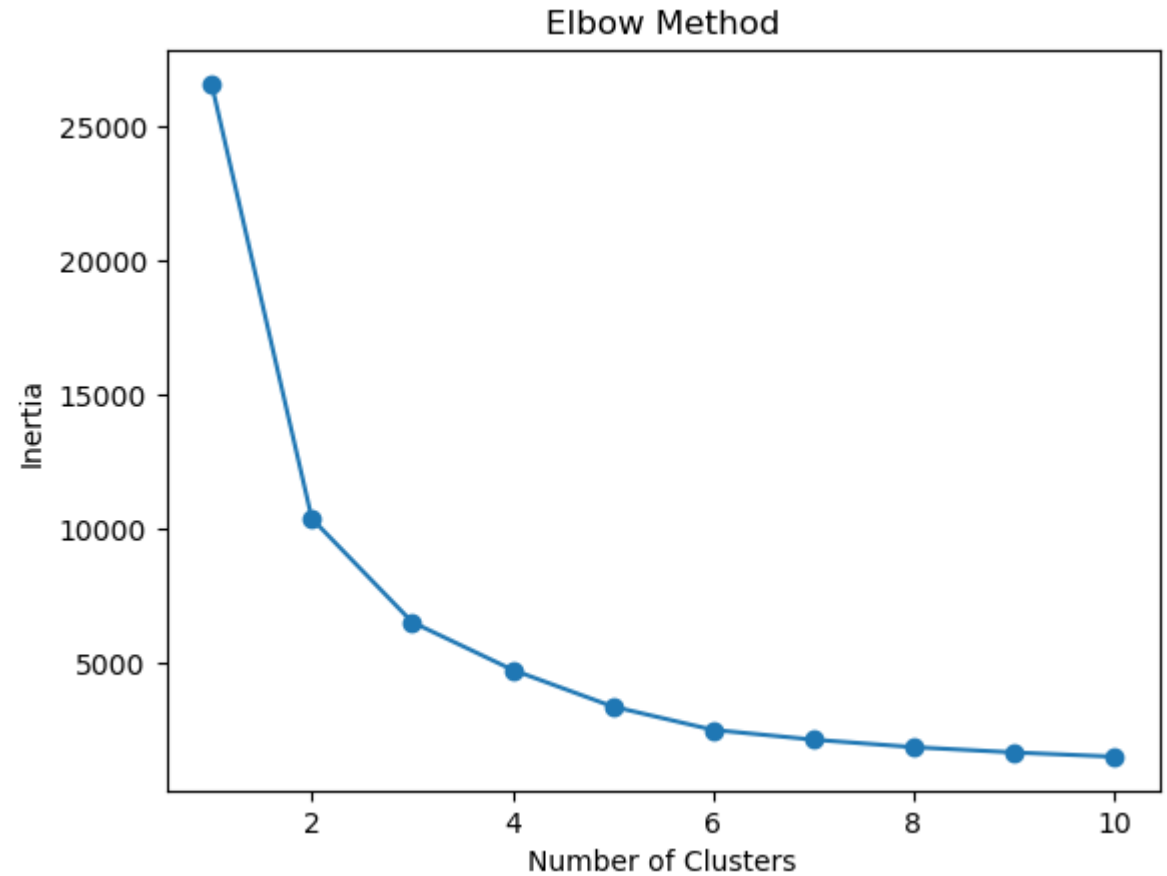# Time Series (Autoregressive model)

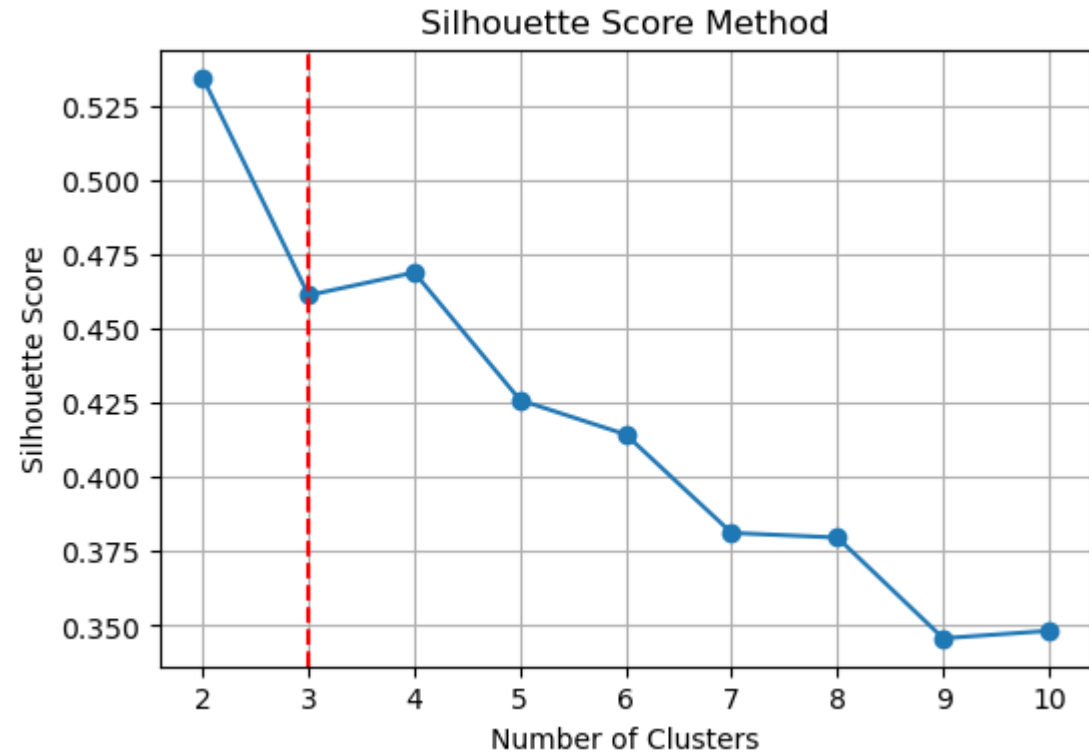# Time Series: ARIMA

# Time Series: XGBoost

# Unsupervised Learning – k number

- The image show that the k=5 is not a bad choice



Elbow Method

# Silhouette Analysis

- Used to determine the degree of separation between clusters



Silhouette Score Method

# Thank you