# Forest Fires

Areeb Shafqat

22nd November 2021

# Main Objectives

✓ Employ a variety of linear regression approaches to predict forest fire (wildfire) area accurately (Prediction focused) with meteorological data based in Portugal

✓ With the data variables or features given such as humidity, wind, rain, moisture, drought, month (Jan, Feb etc., ), day, location coordinates of the fire area, and initial spread index of fire, there is not much benefit in the interpreting or discovering features of high importance for a regression model.

✓ What matters is predictability or fast detection of a wildfire to avoid severe harm to nature, wildlife and humans and enhance firefighting resource management

# Data Summary

## All feature variables:

- X–axis spatial coordinate
- Y-axis spatial coordinate
- Month of the Year
- Day of the week
- Fine Fuel Moisture Code (Moisture Conditions)
- Duff Moisture Code (DMC)
- Drought Code (DC)
- Initial Spread Index (ISI)
- Temperature (Celsius Degrees)
- Relative Humidity
- Wind Speed (Km/h)
- Outside Rain (mm/m2)

- Data is available at https://archive.ics.uci.edu/ml/datasets/Forest+Fires

- There are 517 observations and 13 columns (2 categorical which are Month & Day and the rest numerical)

- No null and duplicated observations

- Area burnt (Target variable) is measured in hectares from 0.00 to 1090.84

- The target is very much skewed towards 0.0 and hence a logarithm transform will be performed later.
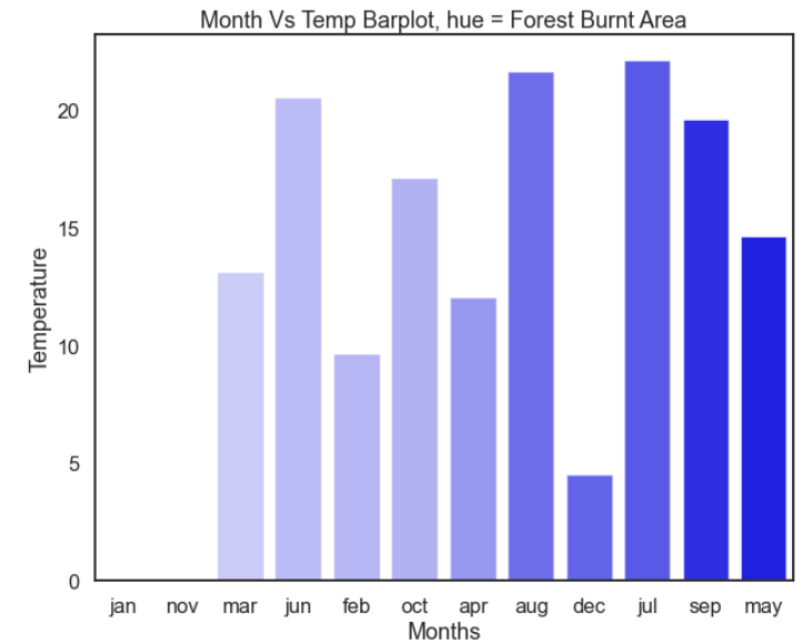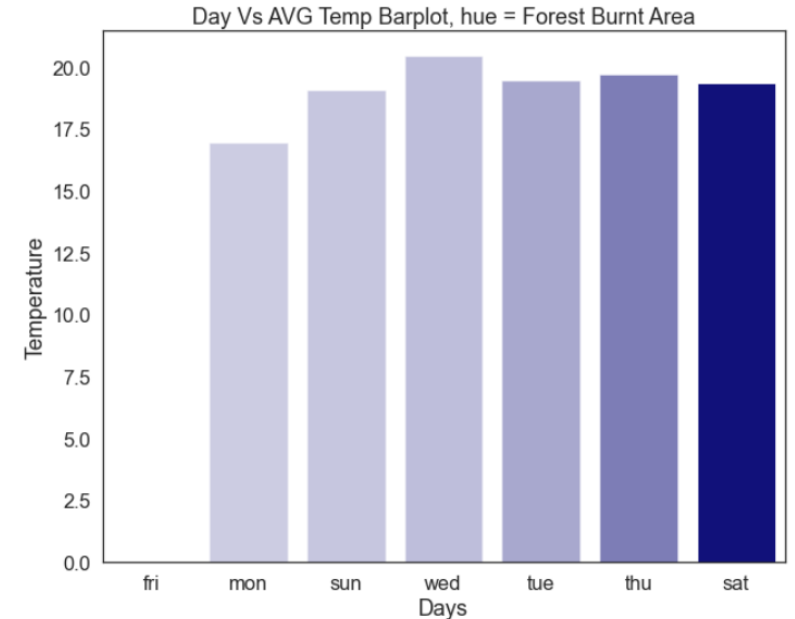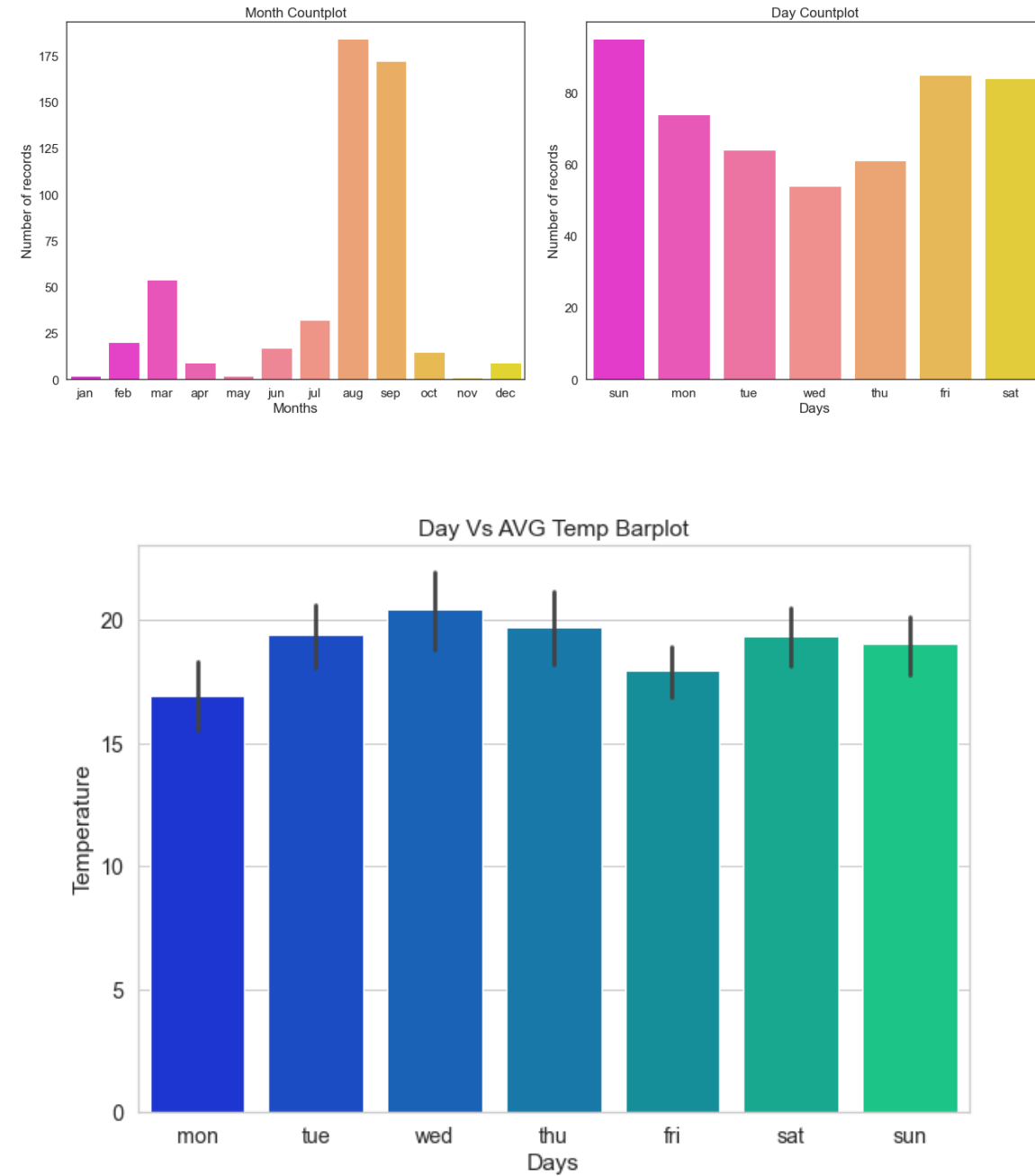
# Data Exploration and Engineering

# Generated Insights

- From the first plot, we can infer that Saturday and Thursday are days where there is most likely to be a larger forest wildfire relative to other days of the week.

- Based on both the plots of temperature against months and days, it is evident that a higher temperature does not necessarily lead to larger area forest fires where for example, contrary to what one might expect, December has one of the highest forest burnt area on average (The color concentration represents the mean area of forest fire for a month or day).

- The bottom plot indicates that September, July and May are the months to look out for when preparing for wildfires.



Day Vs AVG Temp Barplot, hue = Forest Burnt Area



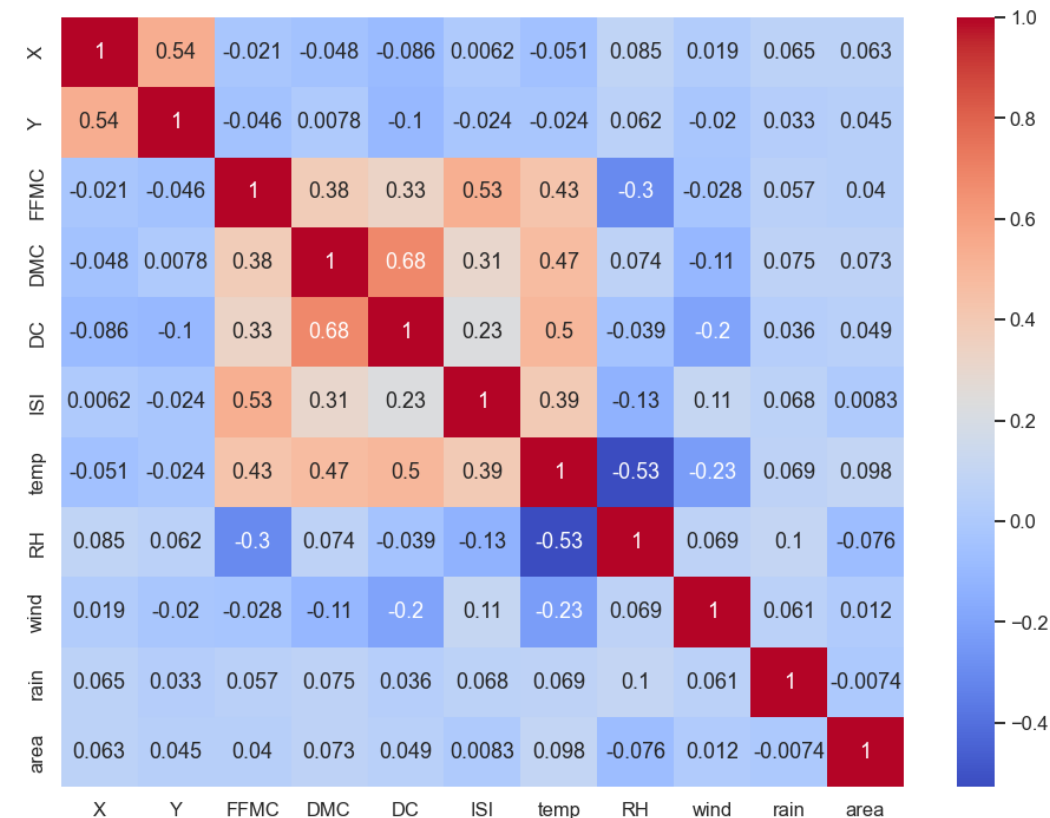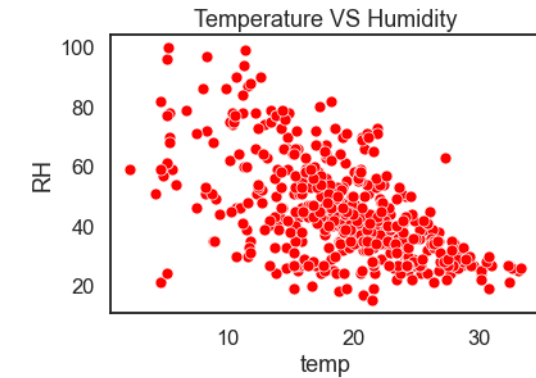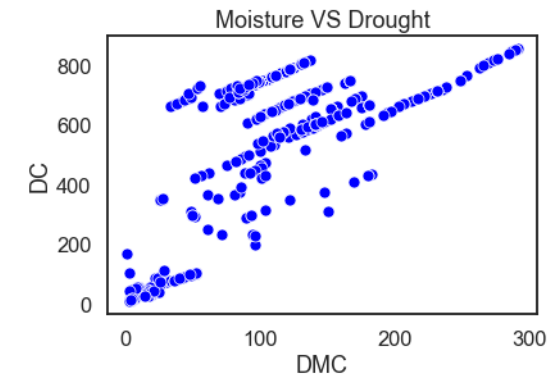Month Vs Temp Barplot, hue = Forest Burnt Area

# Generated Insights

- The plots shown here indicate the number of times a forest fire event was recorded and the temperature variation between days

- The greatest number of wildfires recorded were on Sundays, Saturdays and Fridays

- The most frequent months in which forest fires occurred were August and September although this does not mean that these fires were the most severe which was discussed on the previous slide.

# Generated Insights



- Looking at relationships between variables, we find that temperature and humidity have a negative correlation which is intuitive since higher temperatures result in lower humidity's.

- As for Duff Moisture (DMC) and Drought, these have a strong positive correlation where DMC indicates fire extinguishment

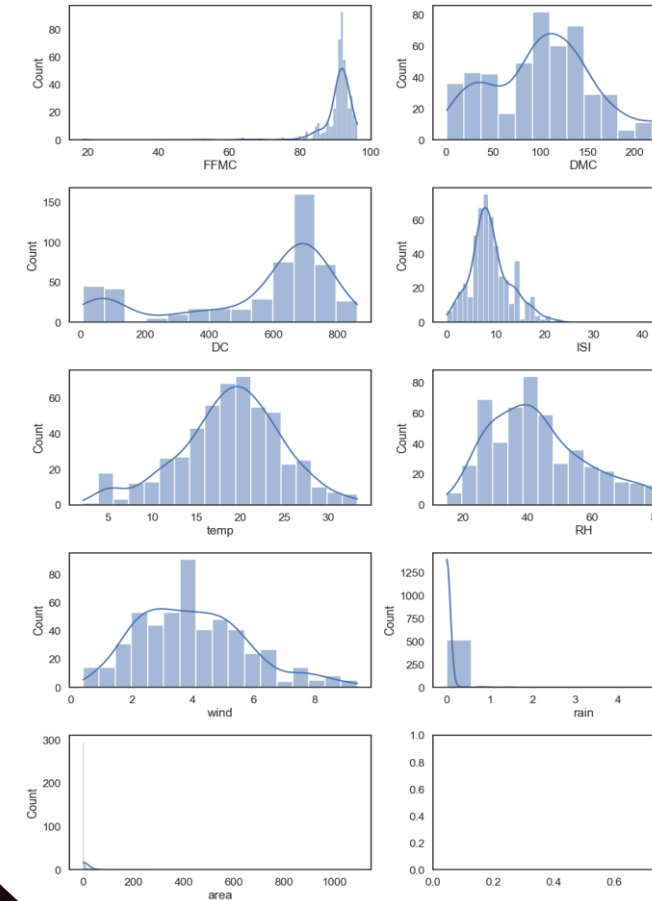- Other correlations include between the X and Y coordinate, temperature and DC, and ISI and FFMC
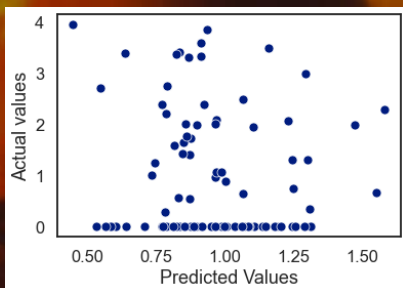
# Data Cleaning Steps



Removal of all outliers with a z-score or standard deviation above 4

Convert categorical columns like Months and Days to numerical e.g., (Jan – 0, Feb – 1,…) , (Sun – 0, Mon - 1…)

To remove much of the skewness (>0.75) of all variables including area, a natural logarithm transform was performed on all the data (np.log1(data))

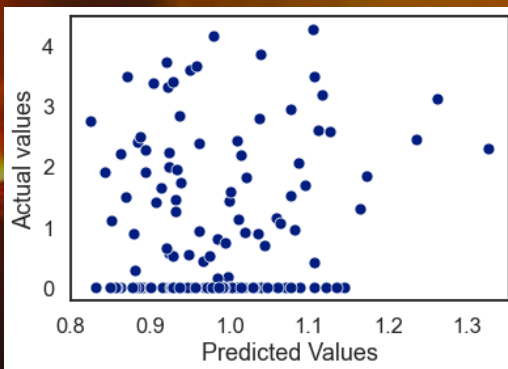# Linear Regressions



Polynomials Features = 2    Standard Scalar    Same Train Test Splits

## Vanilla

| Metric | Value |
|--------|-------|
| RMSE | 1.247 |
| R2 - score | -0.061 |

| Feat | Importance |
|------|-----------|
| FFMC | -3.694 |
| Month | 0.499 |

## Ridge ✔

### RidgeCV

| Metric | Value |
|--------|-------|
| RMSE | 1.2195 |
| R2 - score | 0.00913 |



## Lasso

### LassoCV

| Metric | Value |
|--------|-------|
| RMSE | 1.226 |
| R2 - score | -0.000862 |



## ElasticNet

### ElasticNetCV

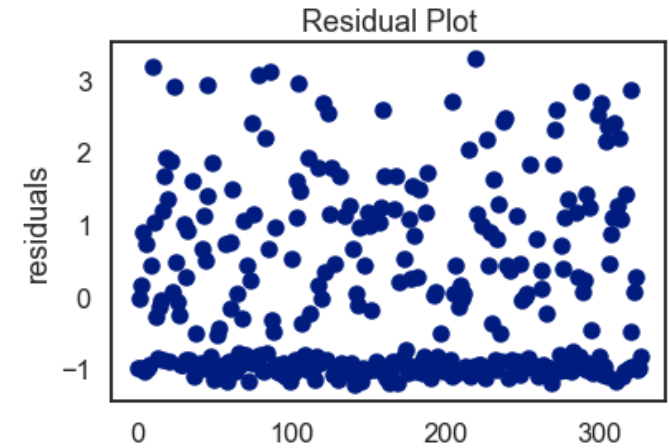| Metric | Value |
|--------|-------|
| RMSE | 1.226 |
| R2 - score | -0.000862 |

# Flaws and Recommendations

- All R squared terms are extremely low which indicates that the independent variables do not explain much of the variation in the area of wildfire variable (dependent)

- This is mainly because the residuals or the errors between predictions and actual values do not follow a normal distribution

- The above result is found through using D'Agostino K^2 Test where an output p-value lower than 0.05 would indicate the distribution of residuals is not normally distributed and hence a linear regression cannot be used for this dataset. The plot of residuals from Ridge Regression, which performs the best on metrics such as RMSE and R_2 squared is shown on the right:

- Using models such as Random Forests and Support Vector Machines would be the ideal choice to go with in this case along with potentially adding data from other countries which could enhance the model tenfold



Residual Plot

P-value = 5.64376948e-09

Thank You!