# Data Worx Ltd

# Spain Electricity ShortFall Project

# TABLE OF CONTENTS

# 01. INTRODUCTION

# ABOUT US

**Services Provided:**

- ❑     **Data collection,**
- ❑     **Data cleaning,**
- ❑     **Data analysis and**
- ❑     **Model building.**

# The Team

1. **Ms. Mandy Rasemphe:** Data Science Team lead
2. **Mr. Karabo Molema:** Data Analyst
3. **Mr. Michael Benjamin:** Machine Learning Engineer
4. **Mr. John Chukwuebuka:** Data Engineer

# 02. Problem Statement

# PROBLEM STATEMENT

❏ **Importance of electricity**

❏ **Inadequate infrastructure**

❏ **Expansion of renewable energy resource infrastructure investments.**

❏ **Model the shortfall between energy generated by means of fossil fuels & various renewable resources**
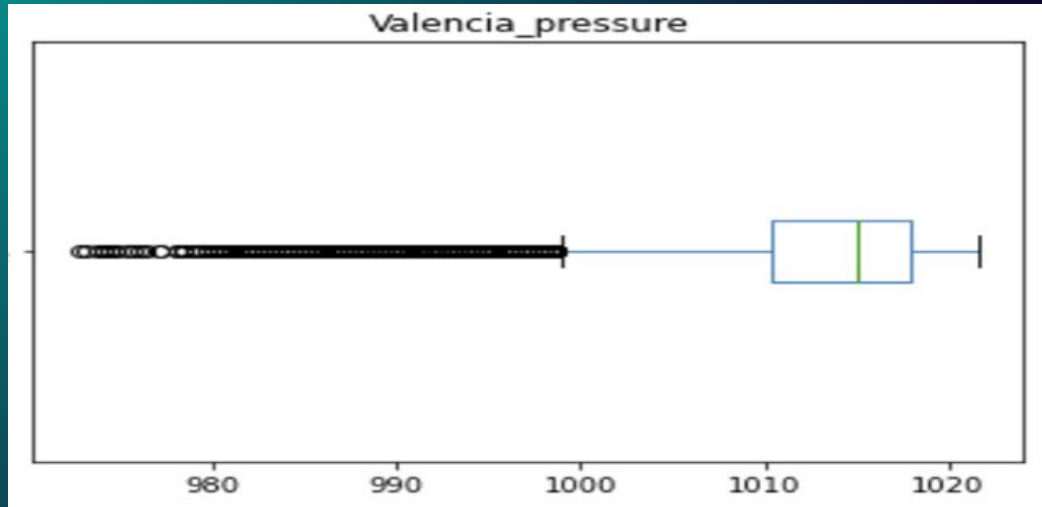
# 03. EDA &
# Feature Engineering

# Investigating Dataset

❏ **Provided Dataset and supporting information explored**

❏ **Contains information on weather conditions.**

❏ **Spain's 5 prominent cities were tracked.**

❏ **Duration is from 1st January 2015 to 31st December 2017**

# Data Issues

## Outliers

- ❑ **28 Features contained outliers**
- ❑ **Keep Outliers**
- ❑ **Use outlier robust techniques.**



Valencia_pressure

# Data Issues

## Missing Values

- One feature contained Null values: Valencia_pressure
- Replaced with median
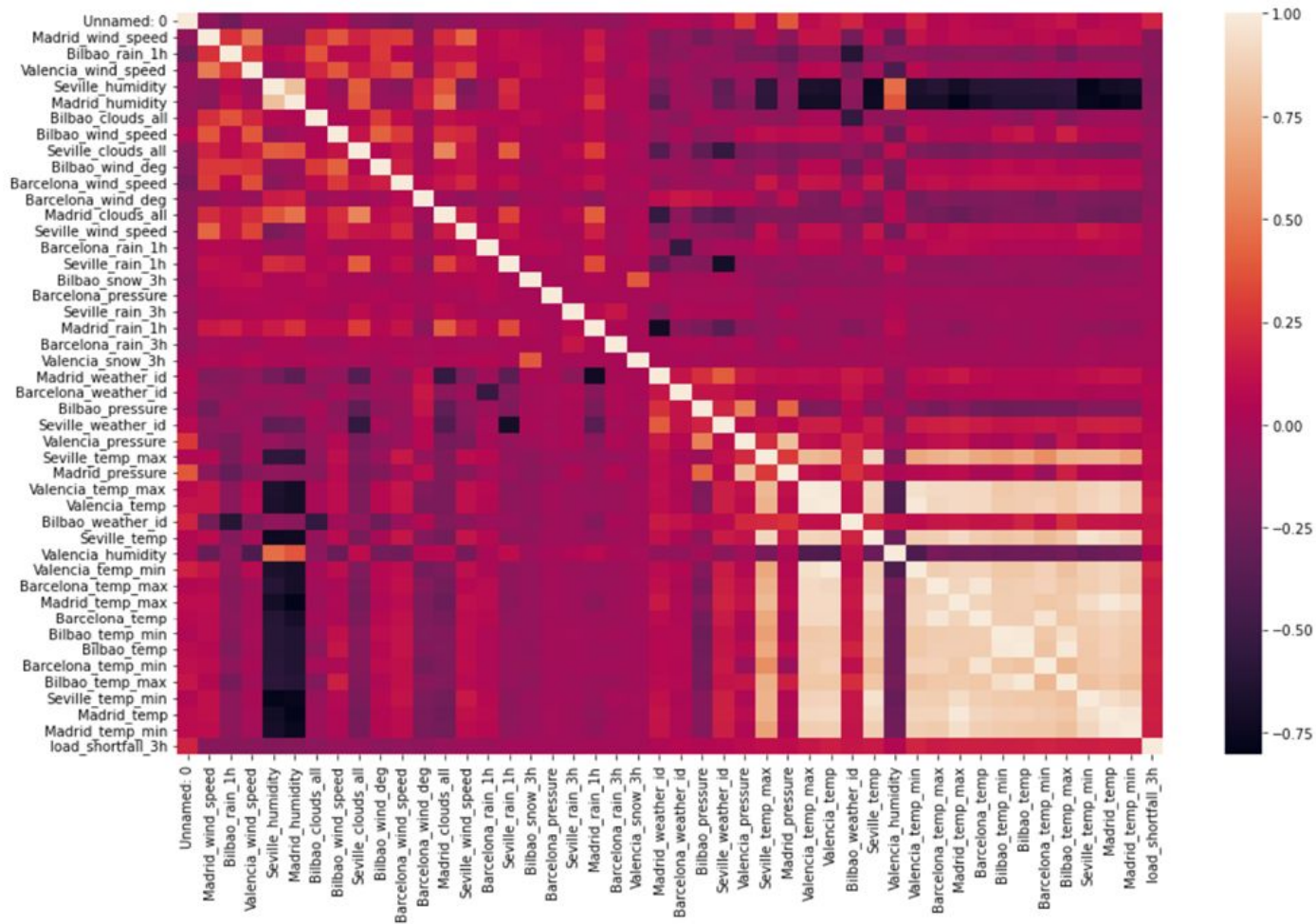
# Data Issues

**Incorrect Feature Data types**

- ❑ Some features were of object data types
- ❑ Features with incorrect data types were converted

# Correlation

❏ **Little to no correlation between predictors and response feature**

❏ **Multicollinearity amongst predictor features**
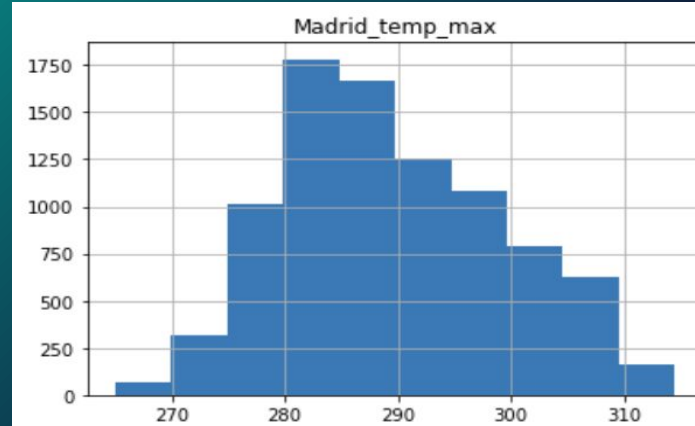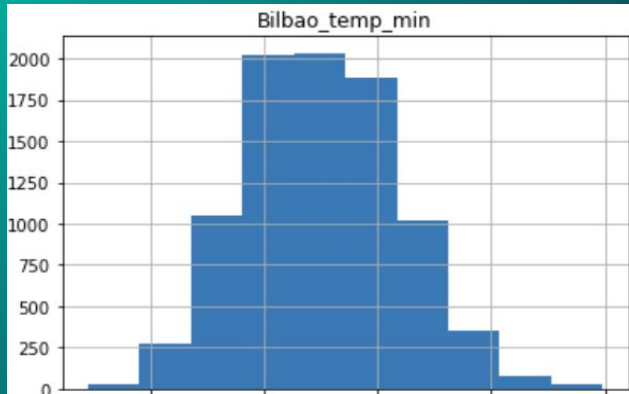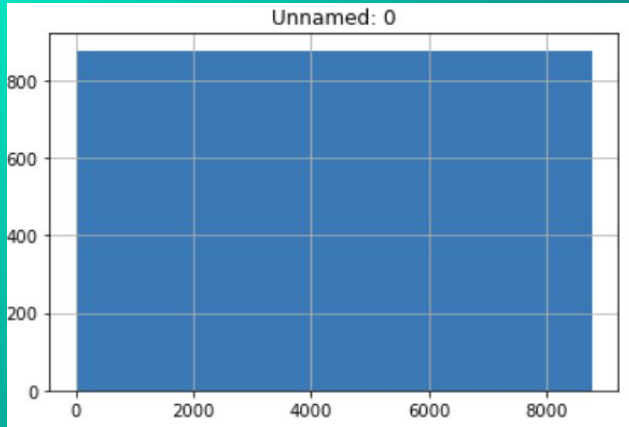
❏ **Drop Multicollinear features**

# Distribution of data

❑ **Distributions of all features**

❑ **Importance of distribution**

❑ **Too much skewness may affect prediction models**

# Distribution of data

# Date time feature

❏ **Split  time feature**

❏ **Get different information such as Year, Month, day, time.**

❏ **Useful in capturing seasonal patterns in a dataset**

# 04. Model Building

# TYPES OF MODELS

- ❏ **Multiple linear regression**

- ❏ **Ridge regression**

- ❏ **Lasso regression**

- ❏ **Decision Tree**

- ❏ **Random Forest**

# 05. Model Evaluation

# Evaluation Metrics

❑ **ROOT MEAN SQUARED ERROR**



**RMSE COMPARISON**

| Model | RMSE |
|---|---|
| MULTIPLE LINEAR REGRESSION | 4846 |
| RIDGE REGRESSION | 4844 |
| LASSO REGRESSION | 4844 |
| DECISION TREE | 3713 |
| RANDOM FOREST | 2920 |

RMSE

Models

# MEAN ABSOLUTE ERROR
# &
# R squared

| | MEAN ABSOLUTE ERROR | R Squared |
|---|---|---|
| Linear Regression | 3857 | 0.17 |
| Ridge | 3858 | 0.17 |
| Lasso | 3858 | 0.17 |
| Decision Tree | 2626 | 0.51 |
| Random Forest | 2285 | 0.70 |

# BUSINESS VALUE OF PROJECT

❏ **Our robust analysis has taken most factors into account.**

❏ **Model will assist in making accurate future predictions.**

❏ **Trends and patterns identified between fossils and renewables.**

❏ **This will assist government make informed business decisions.**

23

# Conclusion

❏ **Show the extent of the shortfall between producing electricity from fossil fuel and renewable source.**

❏ **Random Forest best modelled the data**

❏ **Limitations of project included: Outliers**

❏ **Improvement for future purposes could be further investigation into where the outliers originates.**

❏ **Findings of the project comprehensively highlights the energy deficit between the two systems.**

# QUESTIONS???