# Predicting Solar Energy Production

This project aims to predict solar energy production using machine learning models based on various factors. Solar energy is a key renewable resource for sustainable power generation. The goal is to optimize solar power generation, improve grid efficiency, and support better decision-making in solar energy investments.

by Rishita Shah

# The Importance of Solar Energy

## Objective :

The main objective of predicting solar energy production is to **enhance energy efficiency, optimize resource planning, and improve grid integration** using data-driven insights.

## Dataset :

The dataset contains **218,115** records with **17 columns** .Key observations:
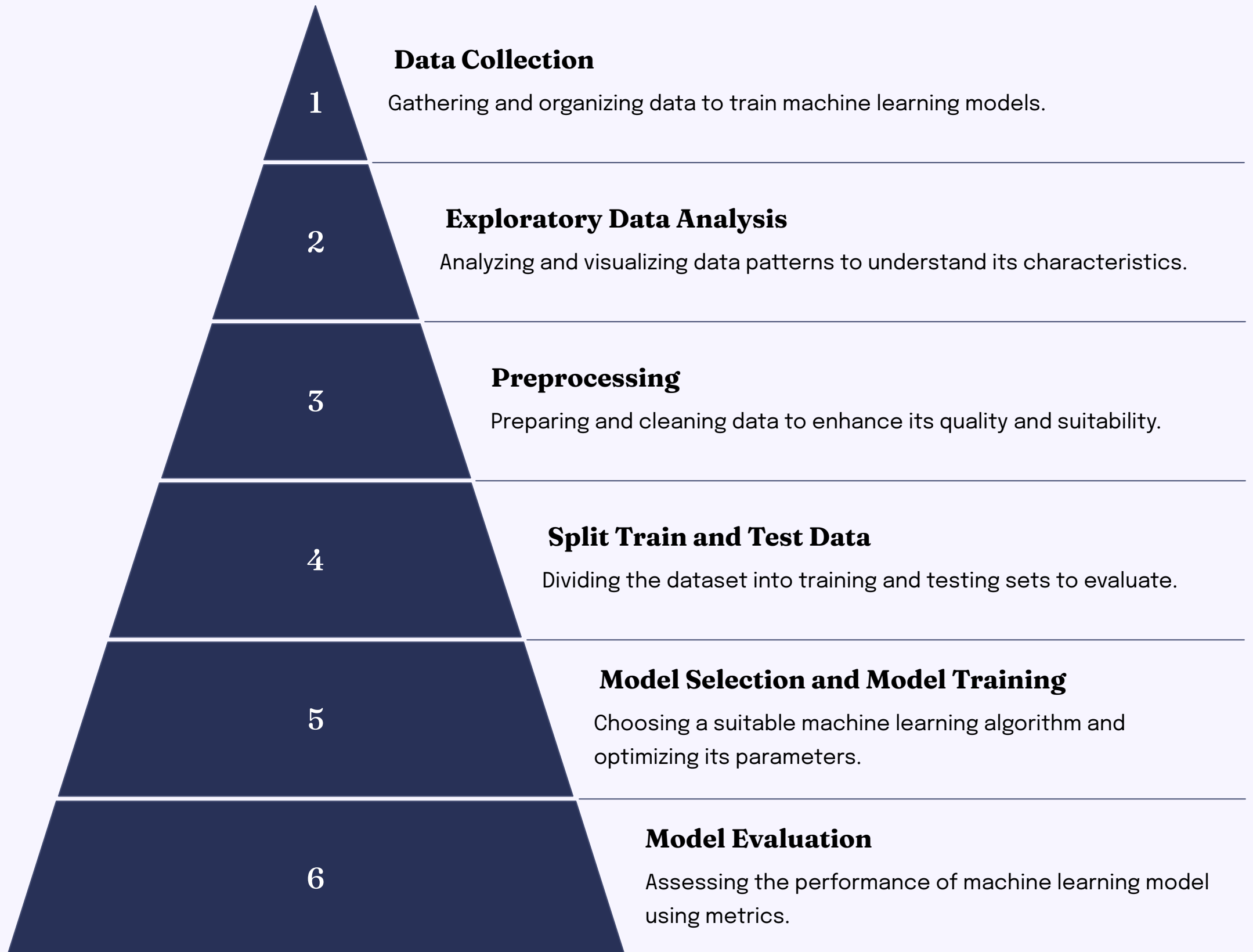
✅ **Target Variable:**

- Estimated Annual PV Energy Production (kWh) → This is the **dependent variable** we aim to predict.

## Importance :

- Improves energy efficiency and grid stability.
- Enhances decision-making for power distribution.
- Helps in cost reduction and resource optimization.

# Steps :

**Data Collection**

1 Gathering and organizing data to train machine learning models.

**Exploratory Data Analysis**

2 Analyzing and visualizing data patterns to understand its characteristics.

**Preprocessing**

3 Preparing and cleaning data to enhance its quality and suitability.

**Split Train and Test Data**

4 Dividing the dataset into training and testing sets to evaluate.

**Model Selection and Model Training**

5 Choosing a suitable machine learning algorithm and optimizing its parameters.

**Model Evaluation**

6 Assessing the performance of machine learning model using metrics.

# Exploratory Data Analysis

**1** **Descriptive Statistics :**

This step involves calculating basic statistics like mean, median, and standard deviation of key features, to understand data characteristics.
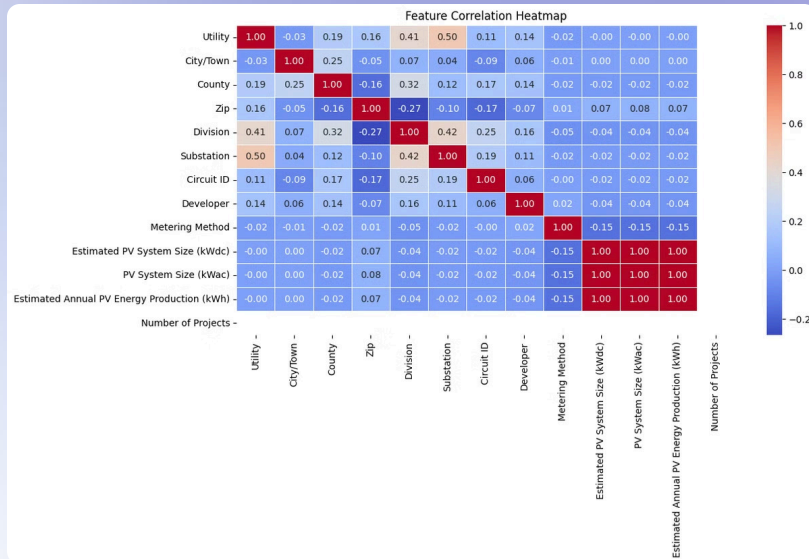
**2** **Correlation Analysis :**

Examining correlations between features reveals potential relationships and helps identify which features are most relevant for segmentation.

**3** **Visualizations :**

Creating visualizations like histograms, scatter plots, and box plots helps understand data patterns and outliers in the data.

# Feature Correlation Heatmap.



Feature Correlation Heatmap

- The color bar on the right shows the correlation scale from **-1 (blue, strong negative correlation)** to **1 (red, strong positive correlation)**.

- The diagonal values are all **1.00** (deep red) because each variable is perfectly correlated with itself.

- **Estimated PV System Size (kWdc), PV System Size (kWac), and Estimated Annual PV Energy Production (kWh)** are **highly correlated** (values close to 1), which makes sense as larger PV systems generate more energy.

- **Substation & Division** show a moderate positive correlation (0.42), indicating some level of dependency.

- **Zip Code & Division have a weak negative correlation (-0.27)**, meaning geographical areas may have different energy profiles.

# Histogram for Numerical Variables.

- These histograms help in understanding how each numerical variable is distributed.

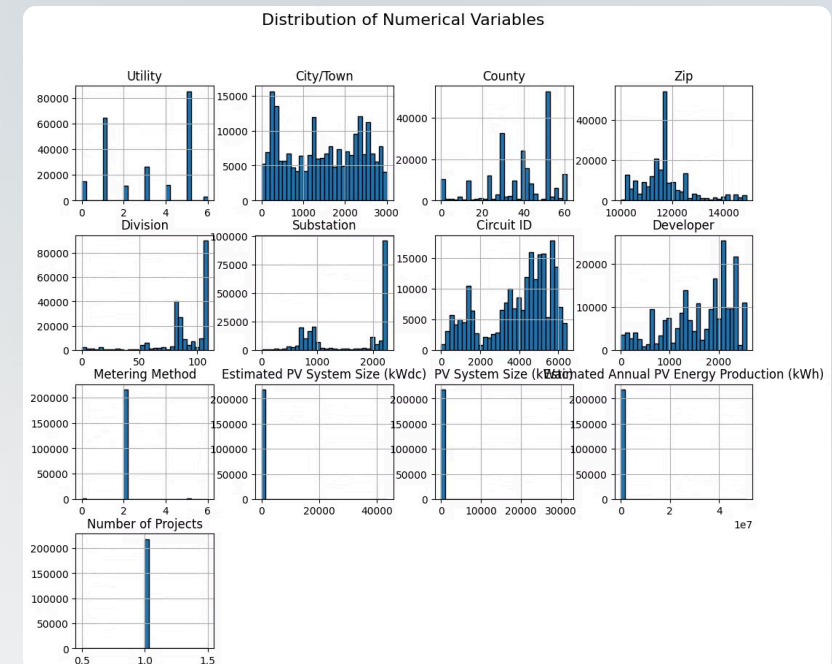- It also provides insights into possible **outliers, skewness, and data concentration**.

## Analysis of Specific Variables:

1. ## Categorical-Like Numerical Variables:

- Some variables like Utility, Division, Metering Method, and Number of Projects have very few unique values, suggesting they are categorical (even though they are represented numerically).

- Metering Method and Number of Projects seem to have only one or very few distinct values, indicating low variance.

2. ## Geographical Variables:

- Variables like City/Town, County, Zip, and Substation show a **multi-modal** distribution, meaning multiple peaks exist. This makes sense, as different locations have different frequencies of solar energy projects.
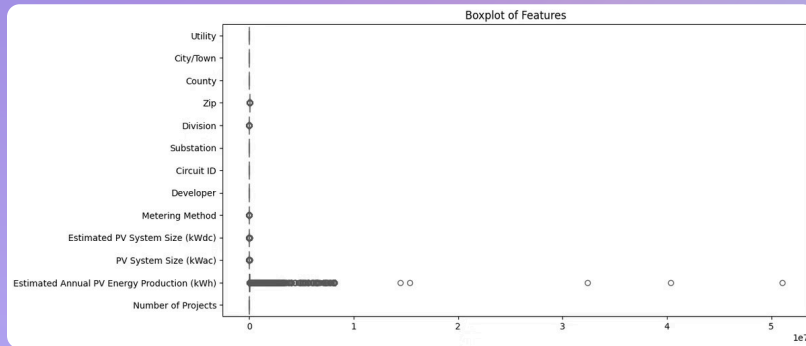


Distribution of Numerical Variables

# Boxplot to Detect Outliers.


Boxplot of Features

## Purpose:

- This boxplot helps visualize the distribution and detect outliers in numerical variables.

## Outliers Observed:

- The **Estimated Annual PV Energy Production (kWh)** has extreme outliers, suggesting the presence of very high-value solar energy projects.

- Other features like **PV System Size (kWdc, kWac)** also exhibit some outliers.

- **Categorical-like variables** (e.g., Utility, Division, and Substation) show small variances but some dispersed values.

Made with Gamma

# Data Collection and Preprocessing

## 1 Data Sources

- The data includes **project details, system specifications, and energy production estimates**.
- The dataset appears to be collected from **solar energy project records**, likely from utility companies or government agencies.

## 2 Data Cleaning

- Remove inconsistencies, missing values, and outliers to ensure data quality and accuracy.

## 3 Feature Engineering

- Create new features from existing ones to improve model performance.

# The Future of Solar Energy Prediction

Advancements in technology, AI, and data analysis are enabling more accurate and sophisticated solar energy prediction models. This will lead to a more reliable and efficient solar energy grid.

# Thank You !!