# Residential-Energy-Consumption-Prediction

This project is all about using **machine learning** to *predict how much electricity homes will use* each year. It starts with a large dataset of household information, then carefully **prepares the data**, *selects the most important characteristics*, and finally uses different **prediction engines** to make accurate forecasts, which are then *fine-tuned for the best results*.

## Visual Overview

```
flowchart TD
    A0["Project Goal: Residential Energy Prediction
"]
    A1["RECS Dataset
"]
    A2["Data Preprocessing & Feature Engineering
"]
    A3["FeatureSelector Tool
"]
    A4["Machine Learning Models (Prediction Engines)
"]
    A5["Model Optimization & Evaluation
"]
    A6["Dimensionality Reduction (PCA)
"]
    A1 -- "Feeds raw data" → A2
    A2 -- "Prepares data for" → A3
    A2 -- "Prepares data for" → A6
    A6 -- "Benchmarks" → A3
    A3 -- "Selects features for" → A4
    A6 -- "Provides reduced data" → A4
```

```
A5 -- "Optimizes" → A4
A4 -- "Achieves" → A0
```

## Chapters

1. Project Goal: Residential Energy Prediction

2. RECS Dataset

3. Data Preprocessing & Feature Engineering

4. FeatureSelector Tool

5. Dimensionality Reduction (PCA)

6. Machine Learning Models (Prediction Engines)

7. Model Optimization & Evaluation

# Main Takeaway:

This step-by-step guide walks students through cloning, configuring, running, and exploring the Residential Energy Consumption Prediction repository. By following this workflow in Notion, learners will gain hands-on experience with virtual environments, Jupyter Notebooks, feature selection, model training, and interpreting energy-consumption predictions.

# 1. Project Setup

# 1.1 Clone the Repository

1. Open your terminal or command prompt.

2. Execute:

```bash
git clone https://github.com/BytesofSurajm/Residential-Energy-Consumption-Prediction.git
cd Residential-Energy-Consumption-Prediction
```

# 1.2 Create and Activate a Virtual Environment

- **Windows**

```
bashpython -m venv venv
venv\Scripts\activate
```

- **macOS/Linux**

```
bashpython -m venv venv
source venv/bin/activate
```

## 1.3 Install Dependencies

```
bashpip install -r requirements.txt
```

# 2. Explore Project Structure

In Notion, create a **Table of Contents** and break down each folder/file:

- **data/**

  Raw and processed datasets used for training and evaluation.

- **imgs/**

  Plots and visualizations generated by notebooks (e.g., correlation matrices, prediction error plots).

- **models/**

  Serialized model files ( `.pkl` or similar) saved after training.

- **feature_selector.py**

  Script implementing feature-importance ranking and selection routines.

- **1.0-Residential-Energy-Modeling.ipynb**

  Core Jupyter Notebook: data loading, preprocessing, exploratory analysis, model training, and evaluation.

- **requirements.txt**

  Python dependencies list.

- **README.md**, **LICENSE**, **.gitignore**

  Documentation, licensing, and Git configuration.

# 3. Launching Jupyter Notebook

1. Ensure the `venv` environment is activated.

2. Run:

`bashjupyter notebook`

3. In the browser, open **1.0-Residential-Energy-Modeling.ipynb**.

# 4. Step-by-Step in Notion

Create a **Notion Page** with the following sections:

## 4.1 Overview

- **Objective:** Predict residential energy consumption using historical data.
- **Key Concepts:** Regression, feature engineering, model evaluation.

## 4.2 Data Exploration

- Embed screenshots of data tables from the notebook.
- Document questions:
  - What features correlate most strongly with energy usage?
  - Are there missing values or outliers?

## 4.3 Feature Selection

- Summarize the `feature_selector.py` logic:
  1. Load dataset.
  2. Compute feature importances (e.g., via tree-based methods).
  3. Select top $k$ features.
- Add code snippet block:

```python
from feature_selector import select_features
X_selected = select_features(X, y, k=10)
```

## 4.4 Model Training and Evaluation

- List algorithms used (e.g., Random Forest, Linear Regression).
- Embed key performance metrics tables (MAE, RMSE, $R^2$).

- Insert plots from **imgs/** folder using Notion's gallery or image block.

## 4.5 Hands-On Exercises

1. **Modify Hyperparameters:** Change tree depth or learning rate; observe impacts.

2. **New Feature:** Engineer a time-of-day or weather-based feature; retrain model.

3. **Cross-Validation:** Implement k-fold CV and compare results.

## 4.6 Advanced Challenge

- Use `models/` directory: load a pretrained model and predict on a new CSV.

- Guide:

```python
import joblib
model = joblib.load('models/best_model.pkl')
predictions = model.predict(new_data)
```

## 5. Learning Outcomes

- Master environment setup and dependency management.

- Develop skills in exploratory data analysis and visualization.

- Understand feature selection techniques and model evaluation metrics.

- Gain practical experience in deploying and using serialized machine-learning models.

# Chapter 1: Project Goal: Residential Energy Prediction

Welcome to the start of your journey into understanding how we can predict energy usage! In this first chapter, we're going to talk about the main "mission" of our project. Think of it like deciding where you want to go on a trip before you even pack your bags. Knowing your destination helps you plan everything else!

## What's Our Big Goal?

Imagine you're a homeowner trying to save on your electricity bill, or maybe you work for an energy company trying to plan how much power to generate. Both would greatly benefit from knowing how much electricity homes are likely to use.

This project's main mission, or **Project Goal**, is to **predict how much electricity (in Kilowatt-hours, or KWH) a residential home will use in a whole year**. We want to do this by looking at different details, or "characteristics," of each home.

It's like having a crystal ball for energy usage!

## Why is This Important?

Knowing annual electricity consumption is incredibly useful:

- **For Homeowners**: It helps them understand their usage patterns, find ways to save energy, and budget for their electricity bills.

- **For Energy Companies**: It allows them to plan better for electricity supply, manage grids efficiently, and even design programs to encourage energy saving.

- **For Policy Makers**: It can inform decisions about energy efficiency standards and sustainable development.

## The "Ingredients" for Our Prediction

To make these predictions, we use special "ingredients":

1. **Household Characteristics Data**: These are details about homes, like their size, age, number of residents, types of heating/cooling systems, appliances, and even the climate zone they are in. This data comes from a very reliable source: the **U.S. Energy Information Administration (EIA)**, specifically from their 2009 Residential Energy Consumption Survey (RECS). We'll talk more about this dataset in the next chapter: RECS Dataset.

2. **Machine Learning**: This is the "magic" tool we use. Machine learning is a way of teaching computers to learn from data without being explicitly programmed. We feed the computer lots of examples of homes and their energy use, and it learns the patterns to make predictions for new homes.

So, combining these ingredients, our project takes raw data about homes and turns it into useful forecasts of energy use.

## How Does the Prediction Process Work (Conceptually)?

Let's look at a very simplified view of what happens to achieve our goal.

```
sequenceDiagram
    participant EIA as U.S. Energy Info Admin
    participant OurProject as Our Project
    participant MLBrain as Machine Learning Brain
    participant Prediction as Energy Prediction (KWH)

    EIA→>OurProject: Provides raw household data (e.g., home size, AC type)
    OurProject→>MLBrain: Sends data to learn patterns
    MLBrain→>MLBrain: "Learns" from data how characteristics relate to KWH
    MLBrain→>OurProject: Provides a trained prediction model
    OurProject→>Prediction: Generates annual KWH forecast
    Note over Prediction: This is our project goal!
```

1. The **U.S. Energy Information Administration (EIA)** collects a huge amount of data about homes and their energy use.

2. **Our Project** takes this raw data as input.

3. We use the **Machine Learning Brain** (which is essentially our computer code running smart algorithms) to "learn" from this data. It figures out the connections between a home's characteristics and its annual electricity consumption.

4. Once the Machine Learning Brain has learned these patterns, it can then generate an **Energy Prediction (KWH)** for new homes we haven't seen before.

This final prediction of annual KWH is what our entire project aims to achieve! As mentioned in the `README.md` file for this project, the goal is to predict "annual electricity usage (KWH) using machine learning."

## What's Next?

Now that we understand *what* we want to predict (annual electricity consumption), the next step is to understand *where* we get the data to make these predictions. In the next chapter, we will dive into the RECS Dataset, which is the

source of all the household characteristics and energy consumption information we use.

# Chapter 2: RECS Dataset

Welcome back! In <u>Chapter 1: Project Goal: Residential Energy Prediction</u>, we learned that our big mission is to predict how much electricity homes use each year. But to make any prediction, we need information, right? Think of it like a chef wanting to bake a cake – they can't just wish for a cake; they need ingredients!

## Where Do We Get Our "Ingredients"?

This is where the **RECS Dataset** comes in. Imagine the RECS Dataset as our project's giant, well-stocked pantry or a valuable treasure chest. It's the essential raw material that holds all the information we need to make our energy predictions. Without it, our project wouldn't have anything to learn from!

## What is the RECS Dataset?

RECS stands for **Residential Energy Consumption Survey**. This specific dataset comes from 2009 and was put together by the **U.S. Energy Information Administration (EIA)**. The EIA is a super reliable source that collects tons of data about energy use in the United States.

So, the 2009 RECS Dataset is basically a huge collection of surveys from thousands of homes across the U.S. It asks all sorts of questions about those homes and their energy use.

## What's Inside Our "Treasure Chest"?

The RECS Dataset is packed with two main types of information:

1. **Household Characteristics (Our "Inputs" or "Features")**:

   - These are all the details *about* the homes. Think of them as the ingredients for our prediction recipe.

   - Examples: How old the home is, how many people live there, what type of heating system it has, if it has a swimming pool, where it's located (which climate zone), and so much more!

- The 2009 RECS is incredibly detailed, containing **over 900 different features** for each home! That's a lot of ingredients!

2. **Annual Electricity Usage (Our "Target" or "Output")**:

- This is the actual amount of electricity each home used in a year, measured in **Kilowatt-hours (KWH)**.

- This is the "dish" we want to cook – the specific number we're trying to predict. For every home, the dataset *already* tells us its KWH usage. This is crucial because our machine learning "brain" learns by seeing examples of inputs (household characteristics) and their corresponding outputs (KWH).

So, in simple terms, the RECS Dataset gives us: "If a home has THIS set of characteristics, it used THIS much KWH." Our goal is to learn that "if...then" relationship!

## How Does Our Project Use the RECS Dataset (Conceptually)?

At the start of our project, the first big step is to "open" this treasure chest and get the data ready.

Imagine this simplified process:

```
sequenceDiagram
    participant RECSFile as RECS Data File
    participant ProjectCode as Our Project Code
    participant InMemTable as Data Table (in Computer Memory)

    RECSFile→>ProjectCode: Provides raw RECS data (e.g., from a file)
    ProjectCode→>InMemTable: Reads and organizes data into a table
    Note over InMemTable: Now our data is ready to use!
```

1. **RECS Data File**: The RECS data is stored in a file, like a large spreadsheet or database.

2. **Our Project Code**: Our computer program (the code we write) acts like a librarian. It knows where to find this data file.

3. **Data Table (in Computer Memory)**: The project code reads the data from the file and loads it into a special format in the computer's memory. This format is often called a "DataFrame," which looks a lot like a neat spreadsheet table with rows and columns. Each row is a different home, and each column is a different feature (like `HOMEAGE` or `KWH` ).

This `Data Table` is the raw form of our ingredients, ready for us to start preparing them.

## Peeking Inside with a Little Code

In the project's main notebook ( `1.0-Residential-Energy-Modeling.ipynb` ), one of the first things that happens is loading this data. Python has a fantastic tool called `pandas` that makes working with tables of data super easy.

Here's a very simple idea of how the data is loaded and what it might look like:

```python
import pandas as pd
# In our actual project, the RECS dataset is loaded from a specific file.# We'll use a placeholder name here for demonstration.try:
    # This line is similar to how the project loads the main RECS data.    # It reads the data from a file into something called a 'DataFrame'.    data = pd.read_csv('recs2009_public.csv') # Imaginary RECS data file    print("First few rows of the RECS dataset:")
    print(data.head()) # 'head()' shows the top 5 rows    print(f"\nTotal columns (features) in the dataset: {data.shape[1]}")
except FileNotFoundError:
    print("Dummy data used: Actual RECS file not found for this example.")
    # If you run this without the actual RECS file, we'll create a small dummy table.    data = pd.DataFrame({
        'HOMEAGE': [1980, 2005, 1995, 2010],
        'NHSLDMEM': [3, 2, 4, 1],
        'KWH': [12000, 8500, 15000, 6000]
    })
    print("Using a small dummy dataset for demonstration:")
    print(data.head())
    print(f"\nTotal columns (features) in the dummy dataset: {data.shape[1]}")
```

**What happened here?**

* We used `pd.read_csv()` to "read" the data from a file (which would be our actual RECS dataset in the real project).

* The `data.head()` command then shows us the first few rows of this data. You'd see columns like `HOMEAGE` (the year the home was built), `NHSLDMEM` (number of household members), and importantly, `KWH` (the annual electricity usage).

* `data.shape[1]` tells us how many columns (features) are in our dataset. As mentioned, for the full RECS, this would be over 900!

The `KWH` column is our **target variable** – the exact value we are trying to predict. All the other columns are our **features** (inputs) that help us make that prediction.

| Category | Description | Example Columns (from RECS) |
|---|---|---|
| **Inputs** | These are the "characteristics" of a home that we use to make a prediction. | `HOMEAGE` (year built), `NHSLDMEM` (household members), `CDD30YR` (climate data) |
| **Target** | This is the specific value we want to predict for each home. | `KWH` (annual electricity usage in Kilowatt-hours) |

## What's Next?

Now that we have our raw ingredients from the RECS Dataset, they're not always perfectly ready to use. Some might need to be cleaned, chopped, or combined with others. This process is called "data preprocessing" and "feature engineering."

In the next chapter, we will learn all about Data Preprocessing & Feature Engineering and how we prepare our RECS data to be ready for the machine learning "brain"!

# Chapter 3: Data Preprocessing & Feature Engineering

Welcome back, aspiring energy predictors! In Chapter 1: Project Goal: Residential Energy Prediction, we set our mission: predicting annual electricity usage. Then, in Chapter 2: RECS Dataset, we learned about our "treasure chest" of raw data, the RECS Dataset, which contains all the household characteristics and their actual energy consumption.

Now, imagine you're a chef, and the RECS Dataset is your basket of fresh ingredients – vegetables, fruits, and meats. While they are fresh, they aren't immediately ready to be cooked! You wouldn't throw muddy carrots or tough, unchopped meat directly into your pot, would you? You need to wash them, peel them, chop them, and perhaps even combine some small pieces.

This is exactly what **Data Preprocessing & Feature Engineering** is all about in our project!

## What's Our Goal in This Chapter?

Our goal here is to understand the crucial steps we take to **clean, transform, and prepare** our raw RECS data so that it's perfect for our machine learning "brain" to learn from. Just like a chef prepares ingredients to enhance the final dish's quality, we prepare our data to make our energy predictions more accurate and reliable.

If we feed "dirty" or unprepared data to our machine learning models, the predictions will be like a bad meal – messy and unhelpful!

## Two Big Ideas: Preprocessing and Feature Engineering

Let's break down this concept into two main parts:

## 1. Data Preprocessing (Cleaning the Ingredients)

This is the "cleaning" phase. Raw data, especially from surveys like RECS, often has issues. Think of it as:

- **Muddy Carrots: Handling Missing Values**: Sometimes, for certain homes, information might be missing. For example, a house might not have reported its number of air conditioning units. These are "missing values." If too much data is missing for a particular characteristic (a column), it might be better to remove that characteristic entirely because it won't help our prediction. We also need to get rid of special "Z" columns, which are just flags indicating

missing data for other features – they don't contain real information we can use.

- **Overripe Fruits: Handling Outliers**: Imagine a home that reports using an extremely high amount of electricity, far more than any other home in the dataset. This could be a data entry error or a very unusual situation. These are called "outliers." They can confuse our machine learning model because it might try to learn from these unusual cases, making it less accurate for typical homes. We might remove extreme outliers to keep our data realistic.

- **Redundant Items: Removing Duplicates**: Sometimes, there might be duplicate characteristics or information that doesn't add new value. We want to remove these to keep our data neat and efficient.

## 2. Feature Engineering (Chopping & Combining Ingredients)

This is the "transforming" phase, where we make our clean data even more useful. Think of it as:

- **Small Herb Bits: Merging Infrequent Categorical Levels**: Many characteristics in the RECS data are "categorical," meaning they describe categories rather than numbers. For example, "Type of Heating System" could be "Gas," "Electric," "Oil," "Wood," etc. Sometimes, a category might have only a few homes (e.g., "Wood" heating might be very rare). If a category is too rare, our model might not have enough examples to learn from it properly. In such cases, we might "engineer" a new feature by combining these rare categories into a single "Other" category. This helps simplify the data and makes it easier for the model to learn.

## How Does This Process Work (Conceptually)?

Here's a simplified view of how our raw RECS data goes through these preparation steps:

```
sequenceDiagram
    participant RawData as Raw RECS Data
    participant Preprocessing as Data Preprocessing Steps
    participant FeatureEngineering as Feature Engineering Steps
    participant PreparedData as Prepared Data for ML
```

```
    RawData->>Preprocessing: Input raw, "dirty" data
    Note over Preprocessing: 1. Handle Missing Values (e.g., remove "Z" colum
ns)
    Note over Preprocessing: 2. Drop Extreme Outliers (e.g., KWH > 80k)
    Note over Preprocessing: 3. Remove Redundant Columns
    Preprocessing->>FeatureEngineering: Pass cleaned data
    Note over FeatureEngineering: 1. Merge Infrequent Categories (e.g., "Other"
heating types)
    FeatureEngineering->>PreparedData: Output ready-to-use data
    Note over PreparedData: Our "ingredients" are now perfect for cooking!
```

First, our `Raw RECS Data` goes through the `Data Preprocessing Steps` where we clean it up. Then, the cleaned data moves to `Feature Engineering Steps` where we refine it further. Finally, we get `Prepared Data for ML`, which is now ready to be fed into our machine learning models.

## A Peek at the Code (Simplified Examples)

In our project, these steps are performed using Python code, especially with the `pandas` library, which is excellent for working with data tables.

Let's look at how we might conceptually do some of these steps. Remember, the actual code in the project is more complex and robust, but these examples show the core idea!

## 1. Dropping Extreme Outliers (KWH)

As mentioned in the `README.md` file under "EDA and Data Transformation," one of the first things we do is drop extreme outliers, especially for our target variable `KWH`.

```python
import pandas as pd
# Assume 'data' is our DataFrame loaded from the RECS dataset# For demons
tration, let's create a tiny dummy DataFramedata = pd.DataFrame({
    'HOMEAGE': [1980, 2005, 1995, 2010, 1970],
    'KWH': [12000, 8500, 15000, 6000, 90000] # 90000 is an extreme outlier})
print("Original data (first few rows):")
```

```
print(data)
# Remove rows where KWH is greater than 80,000 (an example threshold)# T
his removes the "overripe fruit" that could mess up our model.data_cleaned_o
utliers = data[data['KWH'] < 80000]
print("\nData after dropping KWH outliers:")
print(data_cleaned_outliers)
```

**What happened here?** We used a simple rule ( `data['KWH'] < 80000` ) to filter our data.
Any row where the `KWH` value was too high (like 90,000 in our example) was
removed. This ensures our model learns from more typical and reliable energy
consumption figures.

## 2. Removing Redundant Flag Columns ( `Z` columns)

The RECS dataset includes many columns that end with `z` . These columns are not
actual data but rather "flags" that indicate if the corresponding non-Z column has
missing information. They are like notes saying "this information was missing,"
which isn't useful for prediction.

```
# Assume 'data' is our DataFrame with many columns# Let's add a dummy 'Z'
column for demonstrationdata_with_z = pd.DataFrame({
    'HOMEAGE': [1980, 2005],
    'KWH': [12000, 8500],
    'ACBLZ': [0, 1] # Example of a 'Z' column})
print("Data before removing 'Z' columns:")
print(data_with_z)
# Find and remove columns that end with 'Z' (case-insensitive)columns_to_dr
op_z = [col for col in data_with_z.columns if col.endswith('Z') or col.endswith
('z')]
data_no_z_columns = data_with_z.drop(columns=columns_to_drop_z)
print("\nData after removing 'Z' columns:")
print(data_no_z_columns)
```

**What happened here?** We identified all column names that end with `z` (like
`ACBLZ` ) and then used the `.drop()` function to remove them from our dataset. This
cleans up our data by getting rid of unnecessary information.

## 3. Merging Infrequent Categorical Levels

While showing the exact code for merging infrequent categories can be a bit more complex for a beginner chapter, the idea is simple:

Imagine a column like `TYPEHEAT` (Type of Heating). It might have many categories: 'Gas', 'Electric', 'Oil', 'Wood', 'Solar', 'Geothermal'. If 'Wood', 'Solar', and 'Geothermal' only appear for a handful of homes out of thousands, our model won't learn much from them individually.

Instead, we might combine them:
* 'Gas' → 'Gas'
* 'Electric' → 'Electric'
* 'Oil' → 'Oil'
* 'Wood' → 'Other'
* 'Solar' → 'Other'
* 'Geothermal' → 'Other'

This simplifies the `TYPEHEAT` feature, making it more digestible for the machine learning model. The project code identifies categories that appear very rarely and groups them into a single "Other" category.

## Why is this "Preparation" so Important?

Think about it this way:

| Without Data Preprocessing & Feature Engineering | With Data Preprocessing & Feature Engineering |
| --- | --- |
| **Garbage In, Garbage Out**: Unclean data leads to unreliable predictions. | **Clean Input, Reliable Output**: Well-prepared data leads to accurate and trustworthy predictions. |
| **Confused Models**: Missing values, outliers, and too many rare categories can confuse the machine learning model. | **Happy Models**: Models learn patterns more easily and efficiently from clean, well-structured data. |
| **Inefficient Learning**: Models take longer to train and might use up more computer resources trying to handle messy data. | **Efficient Learning**: Models train faster and perform better, as they focus only on meaningful information. |

| Without Data Preprocessing & Feature Engineering | With Data Preprocessing & Feature Engineering |
|---|---|
| **Misleading Insights**: Patterns derived from raw, noisy data might be incorrect, leading to bad decisions. | **Clear Insights**: Clear patterns emerge, providing valuable understanding of energy consumption. |

## What's Next?

By now, our RECS dataset has been cleaned, outliers removed, and features like categories are tidied up. Our "ingredients" are washed, chopped, and neatly organized. But remember, the RECS dataset still has **over 900 features**! That's a lot of characteristics for our machine learning "brain" to process.

Even after cleaning, we might have many features that are still redundant, not very important, or too similar to each other. Too many features can overwhelm our models, slow them down, and sometimes even hurt their performance.

This leads us to the next crucial step: **Feature Selection**. In the next chapter, we will introduce a special tool called the FeatureSelector Tool that helps us intelligently pick the most important ingredients for our prediction "recipe"!

# Chapter 4: FeatureSelector Tool

Welcome back, energy prediction enthusiasts! In Chapter 3: Data Preprocessing & Feature Engineering, we learned how to clean and prepare our raw RECS data, like washing and chopping ingredients for a perfect meal. We handled missing values, outliers, and merged infrequent categories.

However, even after all that preparation, our RECS dataset is still enormous, containing **over 900 different features** (characteristics about homes)! Imagine having a recipe with 900 ingredients – that's a lot to manage, and many might be redundant or not truly essential.

## Why Do We Need a "Smart Filter"?

Having too many features can be a problem for our machine learning "brain" (models):

- **Slower Training**: More features mean more calculations, making our models take a very long time to learn.

- **Confused Models**: Some features might be very similar, or some might just be noise, confusing the model and making it harder for it to find the real patterns.

- **Overfitting**: With too many features, a model might start memorizing the training data too well, instead of learning general rules. This leads to bad predictions on new, unseen homes.

This is where the **FeatureSelector Tool** comes in! Think of it as your project's personal "smart filter" or a super-efficient "chef's assistant" that helps you pick *only the most relevant ingredients* for your prediction recipe. It systematically goes through all 900+ features and identifies which ones are truly important and which can be safely removed.

Our goal in this chapter is to understand how this smart tool works and how we use it to significantly reduce the number of features in our dataset, making our predictions more efficient and accurate.

## What is the FeatureSelector Tool?

The `FeatureSelector` is a custom Python class (a blueprint for creating special objects) designed specifically for our project. It implements several intelligent strategies to identify and remove "unnecessary" features. It's like having a set of rules to decide if an ingredient is good enough to stay in the recipe.

Here are the main strategies our `FeatureSelector` uses:

## 1. Finding Features with Too Many Missing Values (The "Half-Empty Jars")

Imagine a jar of spice that's almost empty. If most of the information for a particular feature (like "number of fireplaces") is missing for many homes, that feature won't be very useful for making predictions. Our `FeatureSelector` can identify and recommend removing features that have more than a certain percentage of missing values.

## 2. Finding Features with Only One Unique Value (The "Always the Same" Ingredient)

Consider a feature like "Does the house have a roof?" If every single house in our dataset *always* has a roof (and this is always recorded as "Yes"), then this feature doesn't help us distinguish between homes or predict their energy use. It provides no new information. The `FeatureSelector` identifies features where all entries are the same, or almost all, and suggests removing them.

## 3. Finding Highly Similar (Correlated) Features (The "Duplicate Spices")

Sometimes, two different features might tell us almost the same thing. For example, "total square footage of the home" and "number of large rooms" might be highly related. If two features are extremely similar (highly correlated), we only need to keep one of them, as the other is redundant. The `FeatureSelector` finds these "duplicate spices" and marks one for removal.

## 4. Finding Features with Zero Importance (The "Flavorless" Ingredient)

After training a small machine learning model, the `FeatureSelector` can ask the model: "Which features did you find absolutely useless for making predictions?" Some features might simply not contribute *any* predictive power. These "flavorless" ingredients can be removed.

## 5. Finding Features with Very Low Importance (The "Tiny Sprinkles")

Even if a feature isn't entirely "flavorless," it might contribute very, very little to the overall prediction. The `FeatureSelector` can also identify features that, when combined, only contribute a tiny fraction of the total importance, allowing us to remove them and simplify our model without losing much accuracy.

## How Do We Use the FeatureSelector Tool?

Using the `FeatureSelector` in our project is like following a few simple steps. First, we "create" the tool, then we tell it to "identify" the features to remove using its strategies, and finally, we tell it to "remove" those features from our dataset.

Here's a simplified look at how it works in our project's code:

```python
import pandas as pd
from feature_selector import FeatureSelector # This is our custom tool!# Assume 'data_prepared' is our DataFrame after Chapter 3's preprocessing# For demonstration, let's create a small dummy DataFrame.# In the real project, data_prepared has over 900 features and includes 'KWH' (our target).data_prepared = pd.DataFrame({
    'feature_A': [1, 2, 3, 4, None],      # Has a missing value    'feature_B': [10, 10, 10, 10, 10],     # Single unique value    'feature_C': [5, 6, 7, 8, 9],
    'feature_D': [5.1, 6.2, 7.3, 8.4, 9.5], # Highly correlated with feature_C    'feature_E': [100, 200, 150, 250, 120], # KWH (Our target variable)    'feature_Z': [0, 0, 0, 0, 0]         # A 'Z' column example})
# Separate features (X) from the target (y)# In our project, 'KWH' is the target variableX_features = data_prepared.drop(columns=['feature_E', 'feature_Z']) # drop target and dummy 'Z'y_labels = data_prepared['feature_E']
print("Original number of features:", X_features.shape[1])
# Step 1: Create an instance of the FeatureSelector# We give it our features (X_features) and our target (y_labels) for importance calculations.fs = FeatureSelector(data = X_features, labels = y_labels)
# Step 2: Tell the FeatureSelector to identify features for removal# We define thresholds for missing values, correlation, and cumulative importance.# 'task' is 'regression' because we are predicting a continuous KWH value.selection_parameters = {
    'missing_threshold': 0.1,  # If >10% missing, consider removal    'correlation_threshold': 0.95, # If correlation >0.95, consider removal    'eval_metric': 'l2',
# Metric for regression model (L2 = squared error)    'task': 'regression',       # Predicting a continuous number (KWH)    'cumulative_importance': 0.99 # Keep features that explain 99% of total importance}
fs.identify_all(selection_parameters)
# Step 3: Get the list of features to be removedfeatures_to_remove = fs.check_removal()
print("\nFeatures identified for removal:", features_to_remove)
# Step 4: Remove the identified features from our datasetdata_final = fs.remove(methods='all') # 'all' means use all identified strategiesprint("\nNumber of features after removal:", data_final.shape[1])
```

```
print("\nFirst few rows of the final dataset:")
print(data_final.head())
```

**What happened here?**

1. We imported the `FeatureSelector` class.

2. We created an object `fs` from this class, giving it our prepared data ( `X_features` ) and our target ( `y_labels` ).

3. We called `fs.identify_all()` to run all its smart filter strategies. This method needs some parameters, like `missing_threshold` (e.g., if more than 10% of values are missing, mark it for removal) and `correlation_threshold` (e.g., if two features are 95% similar, mark one for removal).

4. The `fs.check_removal()` function gave us a list of all the features it decided to remove.

5. Finally, `fs.remove(methods='all')` actually removed those features from our dataset, giving us a much cleaner and more focused set of "ingredients" called `data_final` .

In our actual project, the `FeatureSelector` helps reduce the **428 features** (after Chapter 3's cleaning) down to about **216 features**! That's a huge reduction, making our subsequent steps much more manageable.

## Behind the Scenes: How the FeatureSelector Works Its Magic

Let's peek under the hood to understand how `FeatureSelector` orchestrates these different identification methods.

Imagine our `FeatureSelector` as a conductor of an orchestra. When you call `identify_all()` , it calls on different sections (methods) to perform their specific tasks.

```
sequenceDiagram
    participant You as You (User)
    participant FS as FeatureSelector Tool
    participant IdentifyMissing as identify_missing()
    participant IdentifySingleUnique as identify_single_unique()
    participant IdentifyCollinear as identify_collinear()
    participant IdentifyImportance as identify_zero_importance() / low_importan
```

```
ce()

    You→>FS: fs = FeatureSelector(data, labels)
    You→>FS: fs.identify_all(params)
    FS→>IdentifyMissing: Check for missing values
    IdentifyMissing—>>FS: List of 'half-empty jar' features
    FS→>IdentifySingleUnique: Check for single unique values
    IdentifySingleUnique—>>FS: List of 'always same' features
    FS→>IdentifyCollinear: Check for high correlations
    IdentifyCollinear—>>FS: List of 'duplicate spice' features
    FS→>IdentifyImportance: Train mini-model for importance
    IdentifyImportance—>>FS: List of 'flavorless/tiny sprinkle' features
    FS—>>You: (Internally collects all lists)
    You→>FS: features_to_remove = fs.check_removal()
    You→>FS: data_final = fs.remove(methods='all')
    FS—>>You: Cleaned Dataframe (data_final)
```

Each `identify_` method does its specific calculation:

## 1. Checking for Missing Values (Simplified `identify_missing` logic)

The `identify_missing` method calculates the percentage of missing data in each column.

```
# Inside FeatureSelector's identify_missing method:# self.data is the DataFrame given to FeatureSelectormissing_series = self.data.isnull().sum() / self.data.shape[0]
# Filter to find columns where missing_percentage > thresholdrecord_missing = missing_series[missing_series > self.missing_threshold]
# Convert the result to a list of column names to dropto_drop = list(record_missing.index)
# Store these identified featuresself.ops['missing'] = to_drop
```

**What this code does:** It counts how many "empty slots" ( `isnull().sum()` ) there are in each column, divides by the total number of rows ( `data.shape[0]` ) to get a

percentage, and then checks if this percentage is above the `missing_threshold` you set. If it is, that feature is marked for removal.

## 2. Checking for Highly Correlated Features (Simplified `identify_collinear` logic)

The `identify_collinear` method calculates how similar each pair of features is using something called the "correlation coefficient."

```
# Inside FeatureSelector's identify_collinear method:# self.data is the DataFramecorr_matrix = self.data.corr() # Calculate correlation between all features# Get the upper triangle of the correlation matrix# This avoids checking the same pair twice (e.g., A-B is same as B-A)# and avoids checking a feature against itself (correlation is always 1.0)upper = corr_matrix.where(np.triu(np.ones(corr_matrix.shape), k = 1).astype(bool))
# Find columns where any correlation in the upper triangle is above the threshold# We use .abs() because strong negative correlation is also strong similarityto_drop = [column for column in upper.columns if any(upper[column].abs() > self.correlation_threshold)]
# Store these identified featuresself.ops['collinear'] = to_drop
```

**What this code does:** It creates a table (correlation matrix) showing how strongly related every pair of features is. If "Feature X" and "Feature Y" are very highly correlated (e.g., 0.98), the tool marks one of them (`Feature X` in this simplified example) for removal, because `Feature Y` already gives us almost the same information.

## Why is Feature Selection So Important?

| Without Feature Selection | With Feature Selection |
|---|---|
| **Bloated Models**: Too many features make models slow, complex, and harder to understand. | **Lean Models**: Fewer features mean faster training, simpler models, and better understanding. |
| **Noisy Data**: Unnecessary features add "noise" that can confuse the model. | **Clean Signals**: By removing noise, the model focuses on truly important patterns. |
| **Higher Risk of Overfitting**: Model might learn irrelevant details from noise. | **Reduced Overfitting**: Model generalizes better to new data, leading to more reliable |

| Without Feature Selection | With Feature Selection |
| --- | --- |
| | predictions. |
| **Poor Performance**: Predictions can be less accurate due to complexity and noise. | **Improved Performance**: Models are more efficient and often achieve higher prediction accuracy. |

By using the `FeatureSelector` , we turn our enormous list of ingredients into a focused, powerful set of core components, leading to a much better "energy prediction recipe."

## What's Next?

We've now significantly reduced the number of features using our `FeatureSelector` tool. We've gone from over 900 down to around 200 features! This is a great improvement.

However, sometimes even 200 features can be a lot. What if we could represent most of the information in these 200 features using an even smaller, more abstract set of "super-features"? This is the idea behind **Dimensionality Reduction**.

In the next chapter, we will explore Dimensionality Reduction (PCA), a technique that can compress the information from many features into fewer, more meaningful components, preparing our data even further for our machine learning models!

# Chapter 5: Dimensionality Reduction (PCA)

Welcome back, energy prediction enthusiasts! In Chapter 4: FeatureSelector Tool, we became master chefs who smartly filtered our huge list of over 900 ingredients (features) down to a more manageable ~216. We removed features with too many missing values, redundant information, or those that weren't very important.

While 216 features are much better than 900, it's still quite a lot for our machine learning "brain" to process efficiently. Imagine if you were given 216 pieces of information to make a decision – your brain might feel a bit overwhelmed! Also, some of these 216 features might still be indirectly related or convey similar underlying "information" in different ways.

## Why Do We Need Even More Simplification?

Sometimes, many features describe different aspects of the same underlying "thing." For example, imagine describing a car using its engine size, horsepower, and torque. These are three different numbers, but they all relate to the car's "power." What if we could combine these into a single "power score" that captures most of their collective information?

This is exactly what **Dimensionality Reduction** aims to do. It's like taking a very detailed, high-resolution photograph and compressing it into a smaller file. You still see the main picture clearly, but the file is much easier to store, share, and process.

Our goal in this chapter is to understand a powerful technique called **Principal Component Analysis (PCA)**, which helps us **compress** the information from our ~216 features into a much smaller set of **new, combined features** without losing too much of the original "detail" or information that's useful for prediction.

## What is Principal Component Analysis (PCA)?

**Principal Component Analysis (PCA)** is a statistical technique that helps us simplify complex datasets. Here's the core idea:

- **It creates "Principal Components"**: Instead of just picking and dropping original features (like the FeatureSelector), PCA creates entirely *new* features. These new features, called **Principal Components (PCs)**, are combinations of the original ones.

- **It captures "Variance"**: Think of "variance" as the spread or diversity in our data. If all homes had the exact same size, "home size" would have zero variance and wouldn't help predict energy use. The more variation a feature has, the more information it likely contains. PCA tries to find directions in your data where the information (variance) is maximized.

- **The Goal**: PCA's main job is to reduce the number of features while trying to keep as much of the original information (variance) as possible. It's like summarizing a long story into a few key bullet points, ensuring you don't miss the main plot.

**Why is this helpful for our energy prediction project?**

1. **Managing Complexity**: Instead of dealing with ~216 individual features, we might only need to work with, say, 20-30 Principal Components. This simplifies our data a lot.

2. **Faster Models**: Machine learning models often train much faster when they have fewer features to process.

3. **Reduced Noise**: Sometimes, less important features can be like "noise" in the data. PCA can help focus on the main "signals" that truly matter.

4. **Benchmark**: PCA also gives us a good benchmark. If we can get good predictions using a certain number of PCA components, we can compare that to how well our models perform with other feature selection methods. As mentioned in the `README.md` for this project, we found that "**~200 components explain >95% of the variance**." This means we can capture most of the valuable information from our 216 features using a slightly smaller number of PCA components.

## How PCA Works (Conceptually)

Let's imagine our data points are scattered in a 2D graph. PCA works by finding new "directions" or "axes" in the data that capture the most spread (variance).

1. **First Principal Component (PC1)**: PCA finds the direction along which the data varies the most. This is our PC1. It captures the largest amount of information.

2. **Second Principal Component (PC2)**: Then, it finds another direction, completely perpendicular (at a 90-degree angle) to PC1, where the *remaining* variance is highest. This is PC2.

3. **And so on…**: It continues to find more components, each perpendicular to the previous ones, until it has as many components as original features (or fewer, if we specify).

Each Principal Component is a weighted combination of the original features. For example, PC1 might be `0.6 * HOMEAGE + 0.3 * SQFT + 0.1 * NUM_ROOMS + …` . The important thing is that these components are new, abstract features that summarize the information.

The `README.md` file for our project shows a plot called "**PCA**" which illustrates how much variance (information) each component explains. You'll see that a few

components explain a lot, and then it tapers off. We can decide to keep only enough components to explain a high percentage, like 95% of the total variance.

PCA

This plot helps us see that we don't need all 216 features; we can capture a lot of the essential information with fewer components.

## How to Use PCA in Our Project

In our project, we use a powerful library called `scikit-learn` (often shortened to `sklearn` ) which has a built-in PCA tool.

Here's how we'd typically use it:

```python
from sklearn.decomposition import PCA
import pandas as pd
import numpy as np # Needed for dummy data# Assume 'data_selected' is our DataFrame after FeatureSelector,# containing about 216 features.# For demonstration, let's create a small dummy DataFrame with 4 features.data_selected = pd.DataFrame({
    'home_size_sqft': [1500, 2000, 1800, 2200, 1600],
    'num_bedrooms': [3, 4, 3, 5, 3],
    'house_age_years': [30, 10, 25, 5, 40],
    'num_appliances': [5, 7, 6, 8, 5]
})
print("Original number of features:", data_selected.shape[1])
print("First few rows of original data:")
print(data_selected.head(2))
# Step 1: Initialize PCA. We tell it to keep enough components to explain 95% of the variance.pca = PCA(n_components=0.95)
# Step 2: Fit PCA to our data and transform it.# 'fit' learns the best directions (components).# 'transform' converts our original data into these new components.data_pca = pca.fit_transform(data_selected)
print("\nNumber of components after PCA:", data_pca.shape[1])
print("First few rows of PCA data (first 2 components only, can be more):")
```

```
# We only print the first two columns (components) for simplicityprint(data_pc
a[:2, :])
```

**What happened in this code?**

1. We imported the `PCA` tool from `sklearn.decomposition` .

2. We created a dummy `data_selected` DataFrame. In our actual project, this would be the DataFrame that the `FeatureSelector` tool cleaned and reduced to ~216 features.

3. We created a `PCA` object by saying `pca = PCA(n_components=0.95)` . This is a crucial step! It tells PCA, "Please find enough principal components to capture at least 95% of the total information (variance) in my data."

4. Then, `data_pca = pca.fit_transform(data_selected)` is where the magic happens.

   - `fit(data_selected)` makes PCA "learn" from our `data_selected` which directions (components) are most important.

   - `transform(data_selected)` then takes our original data and "projects" it onto these new component directions, giving us `data_pca` , which is our new, compressed dataset.

5. Finally, we print the number of features before and after PCA, and show a glimpse of the new `data_pca` . You'll notice `data_pca` has fewer columns (features) than `data_selected` , achieving our goal of dimensionality reduction! In our project, this step helps reduce features from ~216 to about 200 components to capture 95% of the variance, as highlighted in the `README.md` .

## Behind the Scenes: How PCA Works Its Magic

The `PCA` tool in `scikit-learn` performs complex mathematical calculations very efficiently. Here's a simplified conceptual walkthrough of what happens:

```
sequenceDiagram
    participant YourData as Your Cleaned Data (e.g., 216 features)
    participant SklearnPCA as sklearn.decomposition.PCA
    participant InternalMath as Internal Math Engine
    participant PrincipalComponents as Principal Components (New Features)
```

```
    YourData→>SklearnPCA: Send data (for .fit_transform())
    SklearnPCA→>InternalMath: "Calculate best summary directions"
    Note over InternalMath: Finds directions where data spreads most
    InternalMath→>SklearnPCA: Returns information about these directions
    SklearnPCA→>PrincipalComponents: Transform original data onto new dire
ctions
    Note over PrincipalComponents: Result is fewer columns, but retains key inf
o
```

When you call `pca.fit_transform(data)`, the `sklearn.decomposition.PCA` object:

1. **Analyzes the spread**: It looks at how all your original features vary and how they vary together.

2. **Identifies Principal Components**: Using advanced math (specifically, eigenvectors and eigenvalues), it identifies the new "directions" or "axes" (our Principal Components) that capture the most variance. The first PC captures the most, the second the next most, and so on.

3. **Determines Number of Components**: Since we set `n_components=0.95`, PCA checks how many of these principal components are needed to explain 95% of the total variance in the original data.

4. **Transforms the Data**: It then converts your original data points from their original feature space into this new, smaller space defined by the chosen Principal Components. The result is a new dataset where each column is one of these Principal Components.

## Why is Dimensionality Reduction (PCA) Important?

| Without PCA | With PCA (Dimensionality Reduction) |
|---|---|
| **Complex Models**: More features mean more complicated relationships for models to learn. | **Simpler Models**: Fewer features lead to simpler, easier-to-understand models. |
| **Slower Training**: More calculations mean models take longer to train, especially with large datasets. | **Faster Training**: Models train much quicker on fewer, more concentrated features. |
| **Memory Intensive**: Storing and processing data with many features uses more computer | **Memory Efficient**: Reduced features require less memory. |

| Without PCA | With PCA (Dimensionality Reduction) |
|---|---|
| memory. | |
| **Hard to Visualize**: Difficult to plot or understand data with hundreds of features. | **Easier to Visualize**: If reduced to 2 or 3 components, data can be plotted and understood. |
| **Potential Overfitting**: Models might try to learn from noise in less important features. | **Reduced Overfitting**: By focusing on main information, models generalize better to new data. |

By using PCA, we make our data even more compact and efficient, ensuring that our subsequent machine learning models can learn the most important patterns quickly and accurately for predicting residential energy consumption.

## What's Next?

We've now gone from a massive raw dataset with over 900 features, cleaned it up, used the FeatureSelector Tool to reduce it to ~216 important features, and finally used PCA to condense that information into an even smaller, powerful set of Principal Components. Our data "ingredients" are now perfectly prepared!

The next exciting step is to actually build our **Machine Learning Models** – our "prediction engines" – that will learn from this refined data and start making those valuable annual electricity usage predictions!

Let's move on to Chapter 6: Machine Learning Models (Prediction Engines).

# Chapter 6: Machine Learning Models (Prediction Engines)

Welcome back, energy prediction explorers! In our previous chapters, we've been busy preparing our data. We started with the raw RECS Dataset in Chapter 2: RECS Dataset, then cleaned and refined it in Chapter 3: Data Preprocessing & Feature Engineering. After that, we used the clever Chapter 4: FeatureSelector Tool to pick only the most important features. Finally, in Chapter 5: Dimensionality Reduction (PCA), we compressed all that valuable information into a powerful, smaller set of "super-features" (Principal Components).

Our data "ingredients" are now perfectly prepared, organized, and optimized! But what's next? We have all the refined information about homes, but we still haven't actually built anything that can *predict* energy usage.

## What Are Our "Prediction Engines"?

This is where **Machine Learning Models** come into play! Think of these models as the "engines" or "brains" of our prediction system. They are the core algorithms that take our perfectly prepared data, learn patterns from it, and then use those patterns to forecast annual electricity consumption (KWH) for new homes.

Without these models, our project would be like having a perfectly designed car body and a full tank of fuel, but no engine to make it move!

Our goal in this chapter is to understand:
* What these "prediction engines" are.
* The different types of engines our project explores.
* How they learn from data (training).
* How they make predictions.

## What Do Machine Learning Models Do? (The Core Idea)

At their heart, machine learning models are computer programs designed to **find hidden rules and relationships** within data. Instead of us telling them explicit rules (like "if home is big, then use more energy"), we show them lots and lots of examples.

For our project:
* We feed them thousands of examples of homes, each with its **characteristics** (our refined features from PCA) and its **actual annual KWH usage**.
* The model "learns" from these examples: it figures out how different characteristics influence KWH usage.
* Once it has learned these patterns, we can give it the characteristics of a *new* home (one it has never seen before), and it will use its learned patterns to **predict** that home's KWH usage!

## Different Types of "Engines" for Prediction

Just like cars can have different types of engines (e.g., gasoline, diesel, electric), machine learning has various types of models, each with its own strengths and

ways of learning. Our project explores a few key types:

# 1. Simpler Linear Models (The "Straightforward" Engines)

These models look for **direct, straight-line relationships** between the home's characteristics and its energy consumption. They are like simple, efficient engines that get the job done without too many complex parts.

- **Linear Regression**: This is the simplest. It tries to draw the "best fit" line through our data points to show how features relate to KWH.
- **Ridge / Lasso Regression**: These are like Linear Regression but with extra "regularization" features. Imagine them as having a built-in "tune-up" mechanism that helps them prevent learning too much from noise in the data, making them more stable.

# 2. More Complex Ensemble Methods (The "Teamwork" Engines)

These models are much more sophisticated. Instead of one single decision process, they work by combining the "decisions" of many simpler "sub-models" (often called "decision trees"). They are like powerful engines that achieve great performance by having many components working together.

- **Decision Trees**: Imagine a tree where each branch is a "yes/no" question about a home's characteristic (e.g., "Is the home built before 1990?"). By asking a series of these questions, it arrives at a prediction.
- **Random Forest**: This is a "forest" of many individual Decision Trees. Each tree makes its own prediction, and then the Random Forest combines all their predictions (e.g., by averaging them) to get a more accurate and robust final prediction. It's like having a committee of experts, all giving their opinion.
- **XGBoost (eXtreme Gradient Boosting)**: This is a highly powerful and popular ensemble method. It also builds many decision trees, but it does so in a very clever way: each new tree tries to correct the mistakes made by the previous trees. It's like a team that learns from its past errors to get progressively better.

# How Do We Use These "Prediction Engines"? (Training and Predicting)

Using a machine learning model generally involves two main steps:

1. **Training (Learning)**: We show the model our historical data (the prepared RECS data with known KWH values). The model learns the patterns from this data. This is like teaching a student using many examples.

2. **Predicting (Applying)**: Once the model is trained, we can give it new data (characteristics of homes it hasn't seen before) and ask it to predict their KWH usage. This is like the student applying what they've learned to solve a new problem.

In our project, we use `scikit-learn` (often called `sklearn`), a fantastic Python library that makes working with machine learning models very straightforward.

Let's look at a simplified example using **Linear Regression**, one of our baseline models:

```python
from sklearn.linear_model import LinearRegression
import pandas as pd
import numpy as np # For creating dummy data# --- Step 1: Prepare Dummy Data (Similar to our PCA-transformed data) ---# In our actual project, X_train_pca would be the result of Chapter 5 (PCA)# and y_train would be the actual KWH values for those homes.X_train_pca = pd.DataFrame(np.random.rand(100, 20)) # 100 homes, 20 PCA featuresy_train = pd.Series(np.random.rand(100) * 10000)   # 100 KWH values (0 to 10000)print("Shape of our training features (X):", X_train_pca.shape)
print("Shape of our training target (y):", y_train.shape)
# --- Step 2: Choose and Create our "Engine" (Model) ---# We select the Linear Regression model.linear_model = LinearRegression()
print("\nOur Linear Regression engine is ready!")
# --- Step 3: Train the "Engine" (Fit the Model to Data) ---# This is where the model learns patterns from our data.print("Training the Linear Regression engine...")
linear_model.fit(X_train_pca, y_train)
print("Engine trained! It has learned the patterns.")
# --- Step 4: Make Predictions ---# Imagine we have new, unseen homes (X_test_pca)X_test_pca = pd.DataFrame(np.random.rand(5, 20)) # 5 new homes, 20 PCA featuresprint("\nMaking predictions for new homes...")
predicted_kwh = linear_model.predict(X_test_pca)
```

```
print("Predicted KWH for the first 5 new homes:")
print(predicted_kwh)
```

**What happened here?**

1. We imported

`LinearRegression` from `sklearn.linear_model` .

2. We created some

**dummy** `X_train_pca` (our features, which would be the PCA components from the previous chapter) and `y_train` (the actual KWH values for those homes). In the real project, these are our carefully prepared ingredients!

3. We created our "prediction engine":

`linear_model = LinearRegression()` .

4. We "trained" our engine using

`linear_model.fit(X_train_pca, y_train)` . This is the most crucial step – the model looks at the `X_train_pca` (home characteristics) and `y_train` (actual KWH) and figures out the relationship between them.

5. Finally, we simulated new homes with

`X_test_pca` and used `linear_model.predict(X_test_pca)` to get the **predicted** `KWH` **values**. This is the ultimate output of our project!

## Behind the Scenes: How an "Engine" Learns and Predicts

Let's simplify how a machine learning model, like our Linear Regression example, works internally.

```
sequenceDiagram
    participant RawData as Prepared Data (X_train, y_train)
    participant MLModel as Machine Learning Model (e.g., LinearRegression)
    participant LearnedPatterns as Learned Patterns (Internal Rules)
    participant NewData as New Data (X_test)
    participant Prediction as KWH Prediction

    RawData→>MLModel: Provide data for training (fit())
    MLModel→>LearnedPatterns: "Learns" relationships between features and
KWH
    Note over LearnedPatterns: Model adjusts its internal rules (e.g., line equati
```

on) to fit data
    LearnedPatterns→>MLModel: Stores learned rules
    NewData→>MLModel: Provide new data for prediction (predict())
    MLModel→>Prediction: Applies learned rules to new data to make KWH pre
diction

When you call `model.fit(X_train, y_train)` :
* The
`MLModel` takes your `X_train` (the features) and `y_train` (the actual KWH values).
* It performs calculations to find the best possible
`Learned Patterns` (e.g., for Linear Regression, it finds the best line that minimizes the
error between its predictions and the actual `y_train` ).
* These
`Learned Patterns` are stored inside the `MLModel` .

When you call `model.predict(X_test)` :
* The
`MLModel` takes the `NewData` ( `X_test` , which are just features of unseen homes).
* It applies its
`Learned Patterns` (the rules it figured out during training) to these new features.
* It then outputs the
`KWH Prediction` for each of those new homes.

## Simplified Internal Look (Linear Regression)

The magic behind `linear_model.fit(X, y)` and `linear_model.predict(X_new)` is essentially finding
an equation.

Imagine a super-simplified case where KWH depends only on `home_size` . Linear
Regression tries to find values for `m` and `b` in the equation:

`KWH = m * home_size + b`

When you call `.fit()` :

```
# Conceptual idea inside .fit() for Linear Regression:# It calculates the best
'm' (slope) and 'b' (intercept)# based on all the 'home_size' and 'KWH' examp
les you give it.# For example, it might find:# linear_model.coef_ = 10  # This is
'm'# linear_model.intercept_ = 1000 # This is 'b'
```

The model learns that for every 1 unit increase in `home_size` , `KWH` increases by 10, plus a base usage of 1000. These `coef_` (coefficients) and `intercept_` are the `Learned Patterns` .

When you call `.predict()` :

```
# Conceptual idea inside .predict():# It takes a new home_size (e.g., 2000 sqft)# and plugs it into the learned equation:# predicted_kwh = linear_model.coef_[0] * new_home_size + linear_model.intercept_# predicted_kwh = 10 * 2000 + 1000# predicted_kwh = 20000 + 1000 = 21000 KWH
```

The model simply applies the learned `m` and `b` to the `home_size` of the new home to make a prediction. Of course, in our project, we have many more features, so the equation is much more complex, but the idea is the same!

## Why Are These "Prediction Engines" Important?

| Without Machine Learning Models | With Machine Learning Models (Prediction Engines) |
|---|---|
| **No Prediction Capability**: We have data, but no way to forecast energy usage. | **Forecast Power**: Enables accurate prediction of annual KWH for any home. |
| **Manual Rules**: Would need to write countless "if-then" rules, which is impossible for complex data. | **Automated Learning**: Models learn complex patterns automatically from data. |
| **Guesses Only**: Decisions based on intuition or rough estimates. | **Data-Driven Decisions**: Predictions are based on historical data and learned relationships. |
| **Inefficient for New Data**: Impossible to scale predictions to millions of homes. | **Scalable**: Once trained, models can quickly predict for vast amounts of new data. |

By building and training these machine learning models, we transform our prepared data into a powerful tool that can give us concrete, data-driven answers to our project's main goal: predicting residential energy consumption.

## What's Next?

We now have our trained "prediction engines" ready to make forecasts! But how do we know if these engines are good? Are their predictions accurate? Which engine is the best?

In the next crucial chapter, we will learn all about **Model Optimization & Evaluation**. We'll discover how to measure our models' performance and fine-tune them to make the best possible energy predictions.

Let's move on to Chapter 7: Model Optimization & Evaluation.

# Chapter 7: Model Optimization & Evaluation

Welcome back, intrepid energy prediction explorers! In Chapter 6: Machine Learning Models (Prediction Engines), we finally built and trained our "prediction engines" – our machine learning models. These models have learned patterns from our highly prepared RECS data and are now capable of forecasting annual electricity usage (KWH) for homes.

But here's a crucial question: **How good are these predictions?** Are they accurate? Which of our "engines" (models) performs best? And how can we make them even *better*?

Imagine you're a mechanic. You've just built a custom engine for a race car. You can't just put it in the car and hope for the best, right? You need to:
1.
**Test it rigorously**: Measure its horsepower, fuel efficiency, and acceleration under various conditions. This is **Model Evaluation**.
2.
**Fine-tune it**: Adjust its carburetors, spark plugs, and other settings to squeeze out every bit of extra performance. This is **Model Optimization**.

Our goal in this chapter is to understand these two vital steps: **evaluating** our machine learning models to know how well they perform and **optimizing** them to achieve their peak prediction accuracy for residential energy consumption.

## Part 1: Model Evaluation (Testing Our Engines)

First, let's talk about how we measure if our predictions are good.

# 1. What Metrics Do We Use? (Measuring Engine Performance)

When our model predicts a `KWH` value for a home, and we know the *actual* `KWH` value, we can calculate how "wrong" the prediction was. We use special numbers, called **metrics**, to summarize this "wrongness" or "goodness" across many predictions.

For our project, since we're predicting a continuous number (KWH), two common and very useful metrics are:

- **Root Mean Squared Error (RMSE)**:

  - **What it means**: Think of RMSE as the **average size of the prediction errors** our model makes. If a home actually uses 10,000 KWH, and our model predicts 10,500 KWH, the error is 500 KWH. RMSE collects all these errors, squares them (to make positive and penalize larger errors more), averages them, and then takes the square root.

  - **Goal**: We want RMSE to be **as low as possible**. A lower RMSE means our predictions are, on average, closer to the actual values.

- **R-squared ($R^2$)**:

  - **What it means**: R-squared tells us **how well our model explains the variation in energy consumption**. If all homes used the exact same KWH, there would be no variation. But homes use different amounts, and R-squared tells us how much of that difference our model can account for using the home's features.

  - **Value Range**: It ranges from 0 to 1 (or 0% to 100%).

    - An R-squared of **1 (or 100%)** means our model perfectly explains all the variation – a perfect prediction (rare in real life!).

    - An R-squared of **0 (or 0%)** means our model is no better than just guessing the average KWH for all homes.

  - **Goal**: We want R-squared to be **as high as possible**. A higher R-squared means our model is capturing more of the real-world patterns.

# 2. Cross-Validation (Rigorously Testing Our Engines)

When we train a machine learning model, it learns patterns from the data it *sees*. If we then test it on the *exact same data* it learned from, it's like a student taking a test with the answers already provided – they'll get a perfect score, but it doesn't show if they truly *learned* the material.

To get a reliable estimate of how well our model will perform on *new, unseen* homes, we use **Cross-Validation**.

**The Big Idea of Cross-Validation:**
Instead of just splitting our data once into training and testing sets, we split it *multiple times* in different ways. This gives us several training-testing cycles, and then we average the performance metrics (like RMSE and R-squared) from all those cycles.

Think of it like testing our race car engine on different tracks, or having multiple independent reviewers test our energy prediction software. It gives us a much more stable and trustworthy performance estimate.

The most common type of cross-validation is **K-Fold Cross-Validation**.

**How K-Fold Cross-Validation Works (Conceptually):**

Imagine our entire dataset of homes. We divide it into `K` equal parts (or "folds"). Let's say `K=5` (a common choice).

```
sequenceDiagram
    participant Data as Our Full Dataset (Homes)
    participant Splitter as K-Fold Splitter
    participant Model as ML Model (e.g., XGBoost)
    participant Metrics as Performance Metrics (RMSE, R2)

    Data→>Splitter: Full Dataset
    Splitter→>Splitter: Divides into 5 Folds (e.g., Fold 1, Fold 2, ..., Fold 5)

    Note over Splitter: --- Round 1 ---
    Splitter→>Model: Training Data (Folds 2,3,4,5)
    Splitter→>Model: Testing Data (Fold 1)
    Model→>Model: Trains on (Folds 2,3,4,5)
    Model→>Model: Predicts on (Fold 1)
    Model→>Metrics: Calculate RMSE, R2 for Round 1
```

```
    Note over Splitter: --- Round 2 ---
    Splitter→>Model: Training Data (Folds 1,3,4,5)
    Splitter→>Model: Testing Data (Fold 2)
    Model→>Model: Trains on (Folds 1,3,4,5)
    Model→>Model: Predicts on (Fold 2)
    Model→>Metrics: Calculate RMSE, R2 for Round 2

    Note over Splitter: (Continues for 5 Rounds, each fold gets to be the test set
    once)

    Metrics→>Model: Collects 5 sets of RMSE, R2
    Model→>Model: Calculates Average RMSE, Average R2
    Model→>Data: Reports stable performance estimate
```

In each "round":
1. One fold is set aside as the
**testing data**.
2. The remaining
`K-1` folds are used as the **training data**.
3. The model is
**trained** on the training data.
4. The trained model
**predicts** on the testing data.
5. We
**calculate metrics** (like RMSE, $R^2$) for that round.

This process repeats `K` times, with each fold getting a chance to be the testing data exactly once. Finally, we average the metrics from all `K` rounds. This average is a much more robust and trustworthy measure of our model's performance. The `README.md` notes that we "Used cross-validation to get stable performance estimates."

## A Peek at the Evaluation Code (Simplified)

`scikit-learn` makes this incredibly easy.

```
from sklearn.model_selection import KFold # For splitting data into foldsfrom s
klearn.linear_model import LinearRegression # Our example modelfrom sklear
n.metrics import mean_squared_error, r2_score # For RMSE and R-squaredim
port numpy as np
import pandas as pd
# --- Dummy Data (like our PCA features and KWH) ---X_pca = pd.DataFrame
(np.random.rand(100, 20)) # 100 homes, 20 PCA featuresy_kwh = pd.Series(n
p.random.rand(100) * 20000 + 5000) # 100 KWH values# --- Set up K-Fold Cr
oss-Validation ---kf = KFold(n_splits=5, shuffle=True, random_state=42) # 5 f
olds, shuffle dataall_rmse_scores = []
all_r2_scores = []
fold_num = 1for train_index, test_index in kf.split(X_pca):
    print(f"\n--- Fold {fold_num} ---")
    # Split data for current fold    X_train, X_test = X_pca.iloc[train_index], X_pc
a.iloc[test_index]
    y_train, y_test = y_kwh.iloc[train_index], y_kwh.iloc[test_index]
    # Create and train our model for this fold    model = LinearRegression()
    model.fit(X_train, y_train)
    # Make predictions on the test set for this fold    y_pred = model.predict(X_t
est)
    # Calculate metrics    rmse = np.sqrt(mean_squared_error(y_test, y_pred))
    r2 = r2_score(y_test, y_pred)
    print(f"  RMSE: {rmse:.2f}")
    print(f"  R-squared: {r2:.2f}")
    all_rmse_scores.append(rmse)
    all_r2_scores.append(r2)
    fold_num += 1print("\n--- Cross-Validation Results ---")
print(f"Average RMSE across all folds: {np.mean(all_rmse_scores):.2f}")
print(f"Average R-squared across all folds: {np.mean(all_r2_scores):.2f}")
```

**What happened here?**
1. We set up
`KFold` to split our data 5 times.
2. In a loop, for each fold:
* We split our dummy

`X_pca` and `y_kwh` into `X_train`, `X_test`, `y_train`, `y_test`.
* We created a
`LinearRegression` model and `fit` (trained) it on the `X_train` and `y_train` for *that specific fold*.
* We made predictions (
`y_pred`) on the `X_test` for *that specific fold*.
* We calculated
`RMSE` and `R-squared` by comparing `y_test` (actual values) with `y_pred` (predicted values).
* We stored these scores.
3. Finally, we printed the
**average** `RMSE` and `R-squared` across all 5 folds. This is our reliable performance estimate! The `README.md` shows a "Baseline Models" plot, which is the result of such cross-validation for different models.

## Part 2: Model Optimization (Fine-Tuning Our Engines)

Once we know how well our models perform, the next step is to make them even better. This is where **Hyperparameter Tuning** comes in.

## What Are "Hyperparameters"? (Engine Knobs and Settings)

Think back to our race car engine. It has various settings: fuel mixture, spark timing, valve clearance, etc. These settings aren't learned *by* the engine; they are set *for* the engine *before* it starts running. Changing these settings affects how the engine performs.

In machine learning, models also have such "settings" called **hyperparameters**.
* For a
`RandomForest` model, hyperparameters include:
*
`n_estimators`: How many individual decision trees to build (number of team members).
*
`max_depth`: How "deep" each decision tree can grow (how many questions it can ask).
* For

`XGBoost` , there are many more, like `learning_rate` (how fast the model learns) or `subsample` (what percentage of data to use for each tree).

The optimal values for these hyperparameters are **not learned from the data**; we have to choose them ourselves. Choosing the right hyperparameters can drastically improve a model's performance.

## How to Find the Best Settings: GridSearchCV (Systematically Turning Knobs)

Trying to find the best combination of hyperparameters by hand can be like trying to tune an engine by randomly turning knobs – it's tedious and inefficient.

**GridSearchCV** (Grid Search Cross-Validation) is a smart way to find the best hyperparameters. It systematically tries out **every possible combination** of hyperparameter values that you specify, using cross-validation to evaluate each combination!

**How GridSearchCV Works (Conceptually):**

```
sequenceDiagram
    participant Model as ML Model (e.g., XGBoost)
    participant Params as Hyperparameter Combinations (Grid)
    participant GridSearchCV as GridSearchCV Tool
    participant CrossVal as Cross-Validation Process
    participant Results as Performance Scores

    Model→>GridSearchCV: Give Model to Tune
    Params→>GridSearchCV: Give List of Settings to Try (e.g., {'n_estimators': [100, 200], 'max_depth': [3, 5]})

    GridSearchCV→>GridSearchCV: For each combination of settings:
    Note over GridSearchCV: e.g., First Try: n_estimators=100, max_depth=3

    GridSearchCV→>CrossVal: Evaluate this combination using Cross-Validation
    CrossVal→>GridSearchCV: Returns average performance (e.g., Average RMSE)
```

> Note over GridSearchCV: (Repeats for all combinations)
> GridSearchCV→>Results: Collects performance for all combinations
> GridSearchCV→>Model: Identifies the BEST combination of settings
> GridSearchCV→>Model: Returns the BEST tuned model

GridSearchCV automates this process:
1. You provide a model (e.g., `XGBoost`).
2. You provide a "grid" of hyperparameters: a dictionary where keys are hyperparameter names and values are lists of settings to try (e.g., `{'n_estimators': [100, 200], 'max_depth': [3, 5]}` means it will try 100/3, 100/5, 200/3, 200/5).
3. For
*each* combination, `GridSearchCV` trains and evaluates the model using **cross-validation** (just like we learned in Part 1!).
4. It then identifies the combination that resulted in the
**best average performance** (e.g., lowest RMSE).
5. Finally, it gives you the model trained with these best settings.

The `README.md` clearly states: "Used `GridSearchCV` across all models. Best performing model: **XGBoost**, based on lowest RMSE." This is how we found our champion model!

## A Peek at the Optimization Code (Simplified)

```
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestRegressor # Our example ensemble modelimport pandas as pd
import numpy as np
# --- Dummy Data (like our PCA features and KWH) ---X_pca = pd.DataFrame(np.random.rand(100, 20)) # 100 homes, 20 PCA featuresy_kwh = pd.Series(np.random.rand(100) * 20000 + 5000) # 100 KWH values# --- Step 1: Choose a model ---base_model = RandomForestRegressor(random_state=42) # A simple Random Forest# --- Step 2: Define the "grid" of hyperparameters to try ---param_grid = {
    'n_estimators': [50, 100],  # Try building 50 trees and 100 trees    'max_dept
```

```
h': [5, 10]      # Try trees with max depth of 5 and 10}
# This will test 2 * 2 = 4 combinations.print("Starting GridSearchCV to find th
e best hyperparameters...")
# --- Step 3: Set up GridSearchCV ---# 'cv=3' means use 3-fold cross-validat
ion for each combination# 'scoring='neg_mean_squared_error'' means it will f
ind the lowest RMSE (neg because GridSearchCV aims to maximize score)grid
_search = GridSearchCV(
    estimator=base_model,
    param_grid=param_grid,
    cv=3,
    scoring='neg_mean_squared_error',
    n_jobs=-1, # Use all available CPU cores for faster processing    verbose=0
# Don't print too much detail during search)
# --- Step 4: Run the search! ---grid_search.fit(X_pca, y_kwh)
print("\nGrid Search Complete!")
print(f"Best RMSE found (lower is better): {-grid_search.best_score_:.2f} (Not
e: neg_mean_squared_error is used, so we negate for RMSE)")
print(f"Best Hyperparameters: {grid_search.best_params_}")
# --- Step 5: Get the best model ---best_tuned_model = grid_search.best_esti
mator_
print("\nOur best-tuned Random Forest engine is ready!")
# You can now use best_tuned_model to make predictions# y_pred_tuned = b
est_tuned_model.predict(X_test)
```

**What happened here?**

1. We selected

`RandomForestRegressor` as our model to tune.

2. We defined

`param_grid`, a dictionary listing the hyperparameter names ( `n_estimators` , `max_depth` ) and the specific values we want to try for each.

3. We created a

`GridSearchCV` object, giving it our `base_model` , the `param_grid` , and specifying `cv=3` (to use 3-fold cross-validation for evaluation). `scoring='neg_mean_squared_error'` tells it to try and minimize RMSE (lower is better).

4. We called

`grid_search.fit(X_pca, y_kwh)` to start the entire tuning process. This command runs all the

combinations and performs cross-validation for each.

5. After the search, `grid_search.best_score_` gave us the best RMSE (it's negative because GridSearchCV always tries to maximize, so we negate it back to positive RMSE).

6. `grid_search.best_params_` showed us the exact combination of `n_estimators` and `max_depth` that achieved that best score.

7. Finally, `grid_search.best_estimator_` gave us the actual model object, already trained and optimized with the best hyperparameters. This is the ultimate "tuned engine"!

## Why is Model Optimization & Evaluation So Important?

| Without Evaluation & Optimization | With Evaluation & Optimization |
|---|---|
| **Unknown Performance**: Don't know how accurate predictions are. | **Quantified Performance**: Know exactly how well models predict. |
| **Unreliable Results**: Predictions might work on training data, but fail on new data. | **Reliable Predictions**: Models are thoroughly tested and perform well on unseen data. |
| **Suboptimal Models**: Using default settings might lead to poor accuracy. | **Peak Performance**: Models are fine-tuned to achieve the best possible accuracy. |
| **No Comparison**: Can't compare different models fairly. | **Informed Model Selection**: Can confidently choose the best model for the job. |
| **Wasted Effort**: Building models that aren't truly effective. | **Efficient Project**: Focus on models that deliver real value and accuracy. |

By rigorously evaluating and carefully optimizing our machine learning models, we ensure that our project delivers the most accurate and reliable predictions of residential energy consumption possible. This allows homeowners, energy companies, and policymakers to make better, data-driven decisions.

## Project Complete!

Congratulations! You've reached the end of our journey through the "Residential Energy Consumption Prediction" project concepts.

We started by defining our goal: predicting annual KWH usage. Then, we explored the RECS Dataset, learned about Data Preprocessing & Feature Engineering to clean our data, used the clever FeatureSelector Tool to pick essential features,

and applied <u>Dimensionality Reduction (PCA)</u> to compress our data even further. Finally, we built our <u>Machine Learning Models (Prediction Engines)</u> and, in this chapter, learned how to **optimize and evaluate** them for peak performance.

You now have a foundational understanding of the entire machine learning pipeline used in this project, from raw data to highly accurate energy predictions. Keep exploring, keep learning, and happy coding!