

Machine Learning Project

Approach for Investigating the Impact of Circular Economy on
Taiwan's Economy and the Energy efficiency

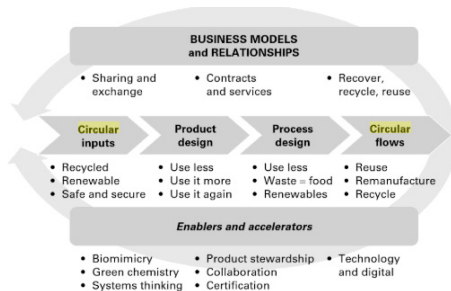
Yuze Tsai

May 17, 2025

Table of Contents

Introduction

- What is Circular Economy?
- Why this topic important?



SOURCE: © Catherine Weetman

References

Variables of interest

- Tanatu, A., et al.(2018): Resource productivity and domestic material consumption has a strong relationship with municipal waste recycling,
- The Ex'TAX project (2022) study: Taxshift can accelerate the transition to circular in construction industry
- Radivojević, V, et al.(2024): A strong positive correlation between circular economy indicators and GDP per capita.

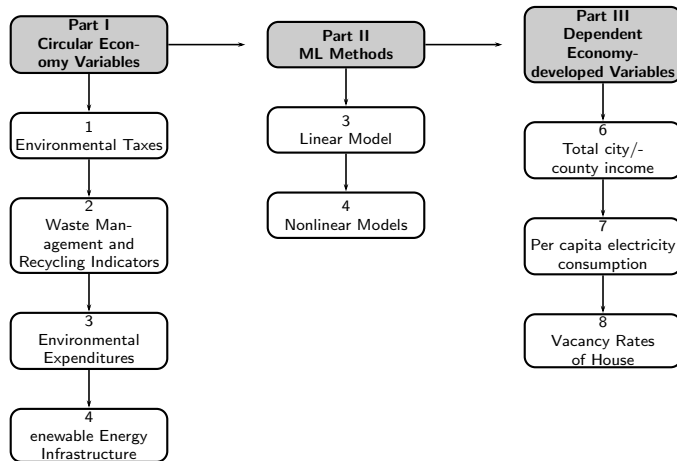
What we have right now...

- **Dependent Variables:** Total city/county income, Per capita electricity consumption, Vacancy Rates of House
- **Independent and Policy-Related Variables** Water pollution environmental tax, Air pollution environmental tax, General waste recycling rate...(total = 11)
- **Year of Interest:** 2014 - 2023 (10 years)

Data Manipulation

```
ML_Project.R  
1 # Data manipulation  
2  
3 library(ggplot2)  
4 library(tidyverse)  
5 library(data.table)  
6 library(readr)  
7 install.packages("arrow")  
8 library(arrow)  
9 gc()  
10  
11 df_elec <- fread("~/Desktop/untitled folder/歷年各縣市再生能源裝置容量(101-11401).csv")  
12 rm(df, data)  
13 df <- fread("~/Desktop/untitled folder/q15019-2140091076.csv")  
14 df2 <- fread("~/Desktop/Result.csv")  
15  
16 df_elec[, year:= 年份+1911]  
17 df_elec[, 1] = NULL  
18 df[, year:= as.numeric(gsub("[^0-9]", "", 統計期)) + 1911]  
19 df[, 1]= NULL  
20 df[, total_income := as.numeric(家庭戶數)*as.numeric(所得收入)]  
21 df_elec <- df_elec[~(287:308), ]  
22 df_elec[, 縣市 := gsub("台北市", "臺北市", 縣市)]  
23 df_elec[, 縣市 := gsub("台中市", "臺中市", 縣市)]  
24 df_elec[, 縣市 := gsub("台南市", "臺南市", 縣市)]  
25 df_elec[, 縣市 := gsub("台東縣", "臺東縣", 縣市)]  
26  
27 result <- merge(df_elec, df, by.x = c("year", "縣市"), by.y = c("year", "地區別"), all = TRUE)  
28 colnames(df2) <- c("year", "縣市", "垃圾焚化量(公噸)",  
29 "垃圾衛生掩埋量(公噸)",  
30 "執行機關資源回收量(公噸)",  
31 "一般廢棄物產生量(公噸)",  
32 "一般廢棄物回收率(%)",  
33 "地方環保單位處決算數(千元)")  
34 df2 <- df2[~(1), ]  
35 df_final <- merge(result, df2, by = c("year", "縣市"), all = T)  
36 rm(df, df_elec, df2)  
37 gc()  
38 df_final[, t_capacity_of_renewable_e := rowSums(.SD, na.rm = TRUE), .SDcols = c("風力", "太陽能電", "其他(含水力)")]  
39  
40 write_csv(df_final, "~/Desktop/untitled folder/output.csv")  
41  
42  
43  
44 people <- fread("~/Desktop/people.csv")  
45 elec <- fread("~/Desktop/untitled folder/elec.csv")  
46 df_transit <- merge(people, elec, by = c("year", "縣市"), all = T)  
47 df_transit <- df_transit[, 1:4]  
48 rm(people, elec)  
49 df_final <- merge(df_final, df_transit, by = c("year", "縣市"), all = T)  
50 df_final <- df_final[~(1:44), ]  
51 df_final <- df_final[~(241:262), ]  
52 write_csv(df_final, "~/Desktop/untitled folder/output.csv")  
53
```


Flow Chart of Study Structure



Linear Regression Method – Panel Data Regression

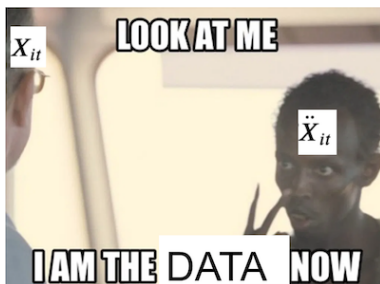
$$y_{it} = a + X'_{it}\beta + c_i + e_{it} \quad (1)$$

where we define:

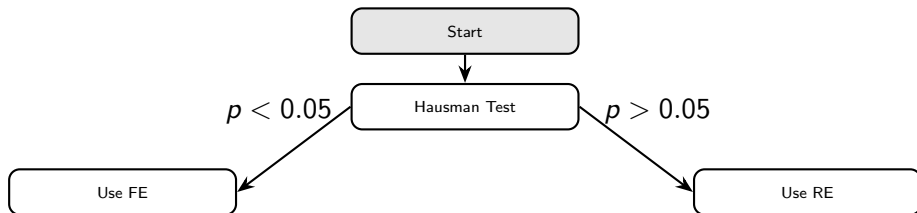
- $i = 1, \dots, N$, where i represents a county/city, and we have N cross-sectional data ($N = 22$).
- $t = 2013, \dots, 2023$.
- y_{it} is the dependent variable.
- X'_{it} are the independent variables of K -dimension.
- a is the intercept.
- β is the coefficient vector of K -dimension.
- c_i is the individual effect.
- e_{it} is the error term.

Fixed Effects vs. Random Effects

Feature	Fixed Effects (FE)	Random Effects (RE)
Individual Effect c_i	Fixed (correlated with X_{it})	Random (uncorrelated)
Main Assumption	Controls for time-invariant differences	Individual effects are uncorrelated
Interpretation	Within-group variation	Between and within variation
When to Use	If c_i is correlated with X_{it}	If c_i is random and uncorrelated



Panel Data Model Selection



Nonlinear regression method – XGBoost

XGBoost (Extreme Gradient Boosting) is an ameliorated gradient boosting method that efficiently handles large datasets and has performed exceptionally well in many competitions.

- Uses **gradient boosting** to iteratively minimize errors.

Mathematical Model:

$$\text{Loss} = \sum_{i=1}^N l(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

where:

- The first term represents prediction error.
- The second term is a regularization term controlling model complexity.

Pros and Cons:

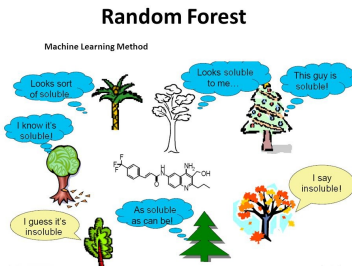
- **Pros:** High accuracy, fast training.
- **Cons:** Complex parameter tuning.

Nonlinear regression method – Random Forest

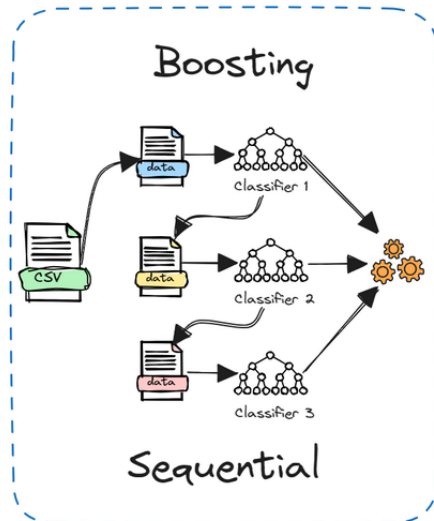
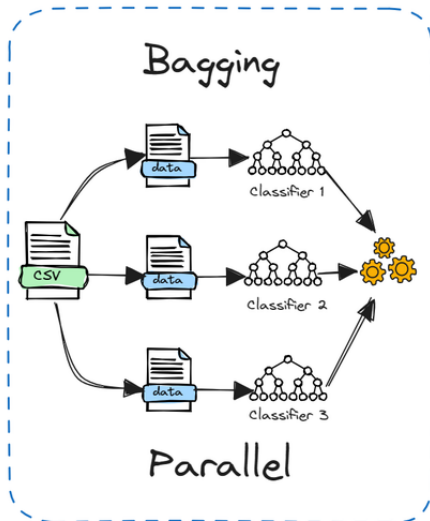
Random Forest is an **ensemble learning** method consisting of multiple decision trees. The final prediction is based on **voting (classification)** or **averaging (regression)**.

Pros and Cons:

- **Pros:** Handles high-dimensional data, reduces overfitting.
- **Cons:** Long training time, less interpretable.



Boosting vs Bagging



Nonlinear regression method – Support Vector Machine (SVM)

SVM is a supervised learning algorithm primarily used for classification.

- It finds the **optimal hyperplane** that best separates different classes.
- Uses **kernel functions** to handle non-linear classification.

Mathematical Formulation:

$$\max_{\mathbf{w}, b} \quad \frac{1}{\|\mathbf{w}\|}$$

$$\text{subject to } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad \forall i$$

Pros and Cons:

- **Pros:** Effective in small sample, high-dimensional data.
- **Cons:** Computationally expensive, slow for large datasets.

Comparison of ML Methods

Feature	SVM	Random Forest	XGBoost
Speed	Slow	Medium	Fast
Accuracy	High	High	Very High
Interpretability	Low	Medium	Low
Handles Large Data	No	Yes	Yes

Intuitive Clairvoyance...

- Environmental budgets and recycling efforts have a positive relationship with total income of a city.
- Cities that invested more in renewable energy infrastructure will generate some critical improvements in energy efficiency.
- Higher levels of incineration and landfill usage are associated with lower energy efficiency, possibly due to poor recycling practices.

Hausman Test for Model Selection

- Null Hypothesis: $H_0 : \mathbb{E}(c_i|X_{it}) = 0 \Rightarrow$ Random Effects (RE)
- Alternative Hypothesis: $H_1 : \mathbb{E}(c_i|X_{it}) \neq 0 \Rightarrow$ Fixed Effects (FE)
- Based on the Hausman test, FE model is chosen if $p < 0.05$, suggesting correlation between individual effects and regressors.
- In this study, FE models are generally preferred due to practical data nature.

Model Evaluation and Comparison

- **Performance Metrics:**

- Mean Squared Error (MSE)
- R-Squared (R^2)
- Cross-validation (10-fold)

- **Findings:**

- Nonlinear models generally outperform linear ones.
- CatBoost achieved the highest R^2 for household income and vacancy rate.
- SVM showed poor predictive performance; excluded from SHAP analysis.

SHAP (SHapley Additive Explanations)

- SHAP assigns each feature an importance value for prediction.
- Based on cooperative game theory—accounts for all possible feature combinations.
- **Advantages:**
 - Global and local interpretability
 - Indicates both direction (+/-) and magnitude of effect
 - Helps identify nonlinear thresholds and interactions
- **Applied to:** XGBoost, CatBoost, LightGBM, Random Forest