

Documentation for recreating "Two-dimensional Pareto frontier forecasting for technology planning and roadmapping"

Arshia Feizmohammady

August 8, 2025

Introduction

This document provides thorough documentation for recreating the results of the article written by Ilya Yuskevich, Rob Vingerhoeds, and Alessandro Golkar. It includes an overview of the methodology, code snippets, and explanations to ensure reproducibility and understanding.

Data Gathering

This section provides an overview of the data gathering and cleaning process. It includes detailed steps on how data was collected, and what sources were used.

The authors of the article utilized the following website to gather their data:

Car Specs Database — cars-data.com [WWW Document], n.d., URL: <https://www.cars-data.com/> (accessed 2.9.20), 2020,

This website hosts a comprehensive dataset of all car models. However, it does not provide an option to directly download the data, necessitating web scraping to obtain the required information.

Due to technical challenges associated with web scraping, an alternative approach was taken. A publicly available dataset from Kaggle was used. This dataset can be found at:

Car Specification Dataset 1945-2020, URL: <https://www.kaggle.com/datasets/jahaidulislam/car-specification-dataset-1945-2020>.

It is important to acknowledge the differences between these two datasets, as

they can impact the recreated figures. The dataset used by the original authors consists of 3,289 datapoints and encompasses car models from 32 car brands representing all major automotive companies. In contrast, the dataset obtained from Kaggle contains over 50,000 datapoints, providing a broader range of information.

Data Visualization

This section visualizes the total max power vs. fuel consumption tradeoff of cars between 1975 and 2015. Below is the Python code used for data visualization:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from matplotlib.colors import LinearSegmentedColormap
4
5 # The dataset has the input Fuel Consumption, but we need to
6   # convert that to the output Miles Per Gallon (MPG)
7 KM_PER_100KM_TO_MPG = 235.214583
8
9 # Load and read the data
10 file_path = 'Car Dataset 1945-2020.csv'
11 car_data = pd.read_csv(file_path, low_memory=False)
```

Listing 1: Data Visualization Code

First, we import the pandas and matplotlib.pyplot libraries for data manipulation and plotting, respectively. Additionally, we import LinearSegmentedColormap for creating custom colormaps. We then define a conversion factor from kilometers per 100 kilometers to miles per gallon (MPG) and load the dataset from a CSV file.

```
1 # Exclude rows with NaN or non-numeric data points in critical
2   # columns - not excluding these will create problems with
3   # the code, possibly not allowing it to run
4 car_data = car_data.dropna(subset=['engine_hp', '
5   mixed_fuel_consumption_per_100_km_l'])
6 car_data = car_data[pd.to_numeric(car_data['engine_hp'], errors='
7   coerce').notnull()]
8 car_data = car_data[pd.to_numeric(car_data['
9   mixed_fuel_consumption_per_100_km_l'], errors='coerce').notnull
10  ()]
```

Next, we exclude rows with NaN or non-numeric data points in critical columns such as engine_hp and mixed_fuel_consumption_per_100_km_l to prevent errors during analysis.

```
1 # Remove duplicate data points - In the article they stated that we
2   # should remove duplicate data points
3 car_data = car_data.drop_duplicates(subset=['engine_hp', '
4   mixed_fuel_consumption_per_100_km_l', 'fuel_grade', '
5   transmission', 'Year_from'])
```

We then remove duplicate data points based on several columns (engine_hp, mixed_fuel_consumption_per_100_km_l, fuel_grade, transmission, Year_from).

```

1 # Define a custom colormap - this custom colormap is to visualize
  the year for each dataset - I followed the gradient they used
  in the article
2 cmap = LinearSegmentedColormap.from_list('custom_cmap', ['darkblue',
  , 'lightblue', 'blue', 'lightgreen', 'green', 'orange', 'yellow'
  ])

```

A custom colormap is defined to visualize the data according to the year, following the gradient used in the article.

```

1 # Filter data to include only petrol and manual transmission cars -
  In the article, they filtered for manual transmission cars
  only
2 # Some petrol car models in the data set aren't defined as petrol,
  but some other keywords, we therefore need to filter for those
  keywords
3 petrol_keywords = ['95', '98', '92', '80', 'Gasoline', 'Gas', '
  Ethanol']
4 filtered_data = car_data[
5     (car_data['fuel_grade'].str.strip().str.lower().isin([kw.lower
6     () for kw in petrol_keywords])) &
7     (car_data['transmission'].str.strip().str.lower() == 'manual'.
8     lower()) &
9     (car_data['Year_from'] >= 1970) &
10    (car_data['Year_from'] <= 2015)
11 ].copy()

```

The data is filtered to include only petrol cars with manual transmissions, produced between 1970 and 2015. Specific keywords are used to identify petrol cars.

```

1 # Print first few rows of filtered data to check filtering - this
  is an optional step, mainly for debugging, I have commented it
  out for now
2
3 ##print(filtered_data.head())

```

An optional step to print the first few rows of filtered data for debugging purposes is included but commented out.

```

1 # Function to preprocess data and convert fuel consumption -
  converting input to output
2 def preprocess_data(data):
3     data['mixed_fuel_consumption_mpg'] = KM_PER_100KM_TO_MPG / data
4     ['mixed_fuel_consumption_per_100_km_l']
5     return data
6
7 # Preprocess data for filtered data
8 filtered_data = preprocess_data(filtered_data)

```

A function preprocess_data is defined to convert fuel consumption to MPG, and the filtered data is preprocessed accordingly.

```

1 # Print the number of data points - the article had around 5000
  datapoints, the dataset I'm using has around 9000
2 # This could result in some discrepancies later on
3 print(f"Number of data points: {len(filtered_data)}")

```

The number of data points in the filtered dataset is printed for comparison with the original article.

```

1 # Function to create scatter plots with actual release years
2 def plot_car_data(data, title):
3     plt.figure(figsize=(12, 6))
4     scatter = plt.scatter(data['engine_hp'], data['
5         mixed_fuel_consumption_mpg'], c=data['Year_from'], cmap=
6         cmap)
7     cbar = plt.colorbar(scatter, label='Year')
8     plt.title(title)
9     plt.xlabel('Max Engine Power (HP)')
10    plt.ylabel('Average Fuel Consumption (MPG)')
11    plt.grid(True)
12    plt.show()
13
14 # Plot the car data for all models
15 plot_car_data(filtered_data, 'Total max power vs Fuel Consumption
16     tradeoff')

```

Finally, a function `plot_car_data` is defined to create scatter plots of engine power vs. fuel consumption, with the color representing the year. The plot is displayed using `matplotlib`.

Comparison of Results - Total Max Power vs. Fuel Consumption Tradeoff

In this subsection, we will compare the generated plot with the corresponding plot from the original article. The comparison will help in understanding the discrepancies and similarities between the recreated results and the article's results.

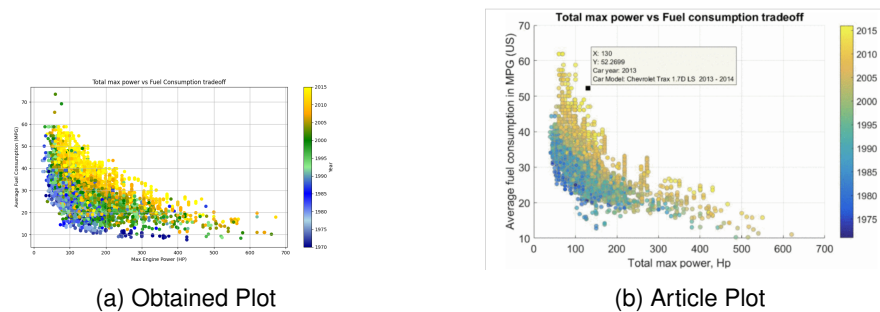


Figure 1: Comparison of Total Max Power vs. Fuel Consumption Tradeoff

By comparing the two images, we can see that while they are not identical, they are quite similar, suggesting that we chose a good dataset to use for recreating the results of this article.

Pareto Frontier

In this section, we will try to find the best responses of the data and estimate a Pareto frontier based on that. Below is the Python code used to achieve this:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from matplotlib.colors import LinearSegmentedColormap
4 import numpy as np
5 from scipy.optimize import curve_fit
6
7 # The dataset has the input Fuel Consumption, but we need to
8   # convert that to the output Miles Per Gallon (MPG)
9 KM_PER_100KM_TO_MPG = 235.214583
10
11 # Load and read the data
12 file_path = 'Car Dataset 1945-2020.csv'
13 car_data = pd.read_csv(file_path, low_memory=False)
```

Listing 2: Pareto Frontier Code

First, we import the necessary libraries: pandas for data manipulation, matplotlib.pyplot for plotting, LinearSegmentedColormap for creating custom colormaps, numpy for numerical operations, and curve_fit from scipy.optimize for fitting a logarithmic curve. We define a conversion factor from kilometers per 100 kilometers to miles per gallon (MPG) and load the dataset from a CSV file.

```
1 # Exclude rows with NaN or non-numeric data points in critical
   # columns
2 car_data = car_data.dropna(subset=['engine_hp', '
   mixed_fuel_consumption_per_100_km_l'])
3 car_data = car_data[pd.to_numeric(car_data['engine_hp'], errors='
   coerce').notnull()]
4 car_data = car_data[pd.to_numeric(car_data['
   mixed_fuel_consumption_per_100_km_l'], errors='coerce').notnull
   ()]
5
6 # Remove duplicate data points
7 car_data = car_data.drop_duplicates(subset=['engine_hp', '
   mixed_fuel_consumption_per_100_km_l', 'fuel_grade', '
   transmission', 'Year_from'])
```

Next, we exclude rows with NaN or non-numeric data points in critical columns such as engine_hp and mixed_fuel_consumption_per_100_km_l to prevent errors during analysis. We also remove duplicate data points based on several columns (engine_hp, mixed_fuel_consumption_per_100_km_l, fuel_grade, transmission, Year_from).

```
1 # Define a custom colormap
2 cmap = LinearSegmentedColormap.from_list('custom_cmap', ['darkblue'
   , 'lightblue', 'blue', 'lightgreen', 'green', 'orange', 'yellow'
   ])
3
4 # Filter data to include only petrol and manual transmission cars
   # from 1970
5 petrol_keywords = ['95', '98', '92', '80', 'Gasoline', 'Gas', '
   Ethanol']
```

```

6 filtered_data = car_data[
7     (car_data['fuel_grade'].str.strip().str.lower().isin([kw.lower
8         () for kw in petrol_keywords])) &
9     (car_data['transmission'].str.strip().str.lower() == 'manual'.
10        lower()) &
11     (car_data['Year_from'] == 1970)
12 ].copy()

```

We define a custom colormap to visualize the data according to the year and filter the data to include only petrol cars with manual transmissions from the year 1970. Specific keywords are used to identify petrol cars.

```

1 # Function to preprocess data and convert fuel consumption
2 def preprocess_data(data):
3     data['mixed_fuel_consumption_mpg'] = KM_PER_100KM_TO_MPG / data
4     data['mixed_fuel_consumption_per_100_km_l']
5     return data
6 # Preprocess data for filtered data
7 filtered_data = preprocess_data(filtered_data)

```

A function preprocess_data is defined to convert fuel consumption to MPG, and the filtered data is preprocessed accordingly.

```

1 # Print the number of data points for 1970
2 print(f"Number of data points in 1970: {len(filtered_data)}")

```

The number of data points in the filtered dataset for the year 1970 is printed.

```

1 # Function to identify maximum points along the x-axis
2 def identify_max_points_along_x(df, step=1, max_x=800):
3     x_values = np.arange(0, max_x+1, step)
4     max_points = []
5     for x in x_values:
6         subset = df[df['engine_hp'] == x]
7         if not subset.empty:
8             max_y = subset['mixed_fuel_consumption_mpg'].max()
9             max_points.append((x, max_y))
10    return pd.DataFrame(max_points, columns=['engine_hp', '
11        mixed_fuel_consumption_mpg'])
12 # Identify maximum points along the x-axis for 1970
13 max_points_1970 = identify_max_points_along_x(filtered_data)

```

We define a function identify_max_points_along_x to identify the maximum points along the x-axis (engine horsepower). This function helps in determining the Pareto frontier. The maximum points for the year 1970 are identified using this function.

```

1 # Function to fit a logarithmic curve
2 def logarithmic_fit(x, a, b, c):
3     return a * np.log(b * (x + 1)) + c
4
5 def fit_logarithmic_curve(max_points):
6     x = max_points['engine_hp'].values
7     y = max_points['mixed_fuel_consumption_mpg'].values
8
9     # Fit the logarithmic curve

```

```

10     popt, _ = curve_fit(logarithmic_fit, x, y, maxfev=10000)
11
12     return popt
13
14 # Fit the logarithmic curve for the maximum points
15 log_params = fit_logarithmic_curve(max_points_1970)

```

A logarithmic function `logarithmic_fit` is defined, and another function `fit_logarithmic_curve` is used to fit this curve to the maximum points identified. The parameters of the fitted logarithmic curve are obtained.

```

1 # Determine the adjusted limits for the plot
2 max_y_pareto = max_points_1970['mixed_fuel_consumption_mpg'].max()
3 max_x_pareto = max_points_1970['engine_hp'].max()
4 y_limit = np.ceil(max_y_pareto / 10) * 10
5 x_limit = np.ceil(max_x_pareto / 50) * 50

```

The adjusted limits for the plot are determined based on the maximum points along the y-axis and x-axis.

```

1 # Function to create scatter plots with actual release years and
  # the Pareto frontier
2 def plot_car_data_with_pareto_frontier(data, max_points, log_params,
  # title):
3     plt.figure(figsize=(12, 6))
4     scatter = plt.scatter(data['engine_hp'], data['
        mixed_fuel_consumption_mpg'], c=data['Year_from'], cmap=
        cmap)
5     cbar = plt.colorbar(scatter, label='Year')
6     plt.title(title)
7     plt.xlabel('Max Engine Power (HP)')
8     plt.ylabel('Average Fuel Consumption (MPG)')
9     plt.grid(True)
10
11 # Plot the Pareto frontier
12 plt.plot(max_points['engine_hp'], max_points['
    mixed_fuel_consumption_mpg'], 'ro-', label='Pareto Frontier
    ')
13
14 # Plot the fitted logarithmic curve within the specified limits
15 x_vals = np.linspace(0, x_limit, 1000)
16 y_vals = logarithmic_fit(x_vals, *log_params)
17 y_vals = np.minimum(y_vals, y_limit)
18 plt.plot(x_vals, y_vals, 'g--', label='Logarithmic Fit')
19
20 plt.legend()
21 plt.xlim(0, x_limit)
22 plt.ylim(0, y_limit)
23 plt.show()
24
25 # Plot the car data for 1970 with Pareto frontier and logarithmic
  # fit
26 plot_car_data_with_pareto_frontier(filtered_data, max_points_1970,
    log_params, '1970: Max Engine Power vs Fuel Consumption
    tradeoff with Pareto Frontier')

```

Finally, a function `plot_car_data_with_pareto_frontier` is defined to create scatter plots of engine power vs. fuel consumption with the actual release

years and the Pareto frontier. The fitted logarithmic curve is also plotted within the specified limits. The car data for 1970 is plotted with the Pareto frontier and the logarithmic fit.

Result

In this subsection, we present the results of the Pareto frontier analysis. The following image shows the Pareto frontier for the year 1970 along with the fitted logarithmic curve.

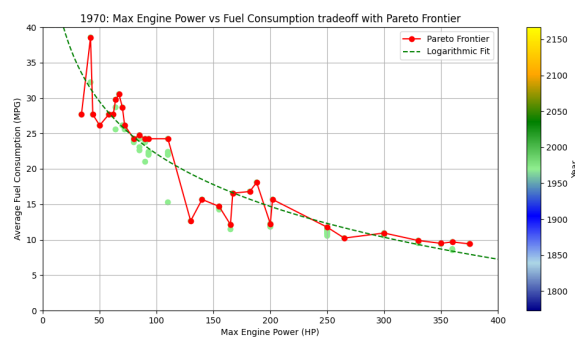


Figure 2: 1970: Max Engine Power vs Fuel Consumption Tradeoff with Pareto Frontier

We can see from the result that this is a good start. However, the logarithmic fit will need to be adjusted in our future code so that no data points are above it.

Pareto Frontiers Backtesting

In this section, we will utilize our Pareto frontier estimation method to estimate the Pareto frontiers until 2050, using both a complete (1970-2015) and incomplete (1970-1995) dataset. We will then compare those two results to backtest our method.

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from matplotlib.colors import LinearSegmentedColormap, Normalize
4 import numpy as np
5 from scipy.optimize import curve_fit
6
7 # Constants
8 KM_PER_100KM_TO_MPG = 235.214583
9 FILE_PATH = 'Car Dataset 1945-2020.csv'
10
11 # Define a custom colormap
```



```

12 cmap = LinearSegmentedColormap.from_list('custom_cmap', ['darkblue'
13         , 'lightblue', 'blue', 'lightgreen', 'green', 'orange', 'yellow'
14         , 'red'])
15 norm = Normalize(vmin=1970, vmax=2060)
16
17 # Load and preprocess the data
18 def load_and_preprocess_data(file_path, start_year, end_year):
19     car_data = pd.read_csv(file_path, low_memory=False)
20     car_data = car_data.dropna(subset=['engine_hp', '
21         mixed_fuel_consumption_per_100_km_l'])
22     car_data = car_data[pd.to_numeric(car_data['engine_hp'], errors
23         ='coerce').notnull()]
24     car_data = car_data[pd.to_numeric(car_data['
25         mixed_fuel_consumption_per_100_km_l'], errors='coerce').
26         notnull()]
27     car_data = car_data.drop_duplicates(subset=['engine_hp', '
28         mixed_fuel_consumption_per_100_km_l', 'fuel_grade', '
29         transmission', 'Year_from'])
30
31     petrol_keywords = ['95', '98', '92', '80', 'Gasoline', 'Gas', '
32         Ethanol']
33     filtered_data = car_data[
34         (car_data['fuel_grade'].str.strip().str.lower().isin([kw.
35             lower() for kw in petrol_keywords])) &
36         (car_data['transmission'].str.strip().str.lower() == '
37             manual'.lower()) &
38         (car_data['Year_from'] >= start_year) &
39         (car_data['Year_from'] <= end_year)
40     ].copy()
41
42     filtered_data['mixed_fuel_consumption_mpg'] =
43         KM_PER_100KM_TO_MPG / filtered_data['
44             mixed_fuel_consumption_per_100_km_l']
45     return filtered_data

```

Listing 3: Pareto Frontiers Backtesting Code

First, we define the necessary constants and a custom colormap. We then load and preprocess the data for the specified year range, filtering for petrol and manual transmission cars. The fuel consumption is converted to MPG.

```

1 # Identify maximum points along the x-axis
2 def identify_max_points_along_x(df, step=1, max_x=800):
3     x_values = np.arange(0, max_x+1, step)
4     max_points = []
5     for x in x_values:
6         subset = df[df['engine_hp'] == x]
7         if not subset.empty:
8             max_y = subset['mixed_fuel_consumption_mpg'].max()
9             max_points.append((x, max_y))
10    return pd.DataFrame(max_points, columns=['engine_hp', '
11        mixed_fuel_consumption_mpg'])
12
13 # Fit a logarithmic curve
14 def logarithmic_fit(x, a, b, c):
15     return a * np.log(b * (x + 1)) + c
16
17 def fit_logarithmic_curve(max_points):

```

```

17 x = max_points['engine_hp'].values
18 y = max_points['mixed_fuel_consumption_mpg'].values
19 popt, _ = curve_fit(logarithmic_fit, x, y, maxfev=10000)
20 return popt

```

We define functions to identify the maximum points along the x-axis (engine horsepower) and fit a logarithmic curve to these points.

```

1 # Adjust the fitted curve upwards
2 def adjust_curve_upwards(data, log_params):
3     x = data['engine_hp'].values
4     y = data['mixed_fuel_consumption_mpg'].values
5     y_fit = logarithmic_fit(x, *log_params)
6     max_diff = np.max(y - y_fit)
7     if max_diff > 0:
8         log_params[2] += max_diff + 1 # Adjust the curve upwards
9     return log_params

```

A function is defined to adjust the fitted logarithmic curve upwards based on the data.

```

1 # Calculate the combined Pareto frontiers
2 def calculate_pareto_frontiers(data, start_year, end_year,
3     projection_start_year):
4     combined_pareto_frontiers = []
5     prev_y_limit = None
6     yearly_improvement_rates = []
7
8     for year in range(start_year, projection_start_year):
9         data_year = data[data['Year_from'] == year]
10        data_year = data_year[data_year['engine_hp'] <= data_year['
11            engine_hp'].quantile(0.99)]
12        data_year = data_year[data_year['mixed_fuel_consumption_mpg
13            '] <= data_year['mixed_fuel_consumption_mpg'].quantile
14            (0.99)]
15
16        max_points_year = identify_max_points_along_x(data_year)
17        log_params_year = fit_logarithmic_curve(max_points_year)
18
19        if not max_points_year.empty:
20            max_y_pareto = max_points_year['
21                mixed_fuel_consumption_mpg'].max()
22            if prev_y_limit is not None:
23                y_limit = max(np.ceil(max_y_pareto / 10) * 10,
24                    prev_y_limit)
25            else:
26                y_limit = np.ceil(max_y_pareto / 10) * 10
27
28            log_params_year = adjust_curve_upwards(data_year,
29                log_params_year)
30            year_color = cmap(norm(year))
31
32            x_vals = np.linspace(0, 900, 1000)
33            y_vals = logarithmic_fit(x_vals, *log_params_year)
34            y_vals = y_vals[x_vals >= data_year['engine_hp'].min()]
35            x_vals = x_vals[x_vals >= data_year['engine_hp'].min()]
36            y_vals = np.minimum(y_vals, y_limit)

```

```

31         combined_pareto_frontiers.append((x_vals, y_vals, year)
32         )
33         if prev_y_limit is not None:
34             yearly_improvement_rates.append(y_limit -
35             prev_y_limit)
36
37         prev_y_limit = y_limit
38
39         average_improvement_rate = np.mean(yearly_improvement_rates)
40         total_y_increase_needed = 99 - prev_y_limit
41         years_remaining = 2060 - (projection_start_year - 1)
42         adjusted_yearly_improvement_rate = total_y_increase_needed /
43         years_remaining
44         current_y_limit = prev_y_limit
45
46         for year in range(projection_start_year, 2061):
47             current_y_limit += adjusted_yearly_improvement_rate
48             prev_pareto_frontier = combined_pareto_frontiers[-1]
49             x_vals_prev, y_vals_prev, _ = prev_pareto_frontier
50             y_vals_new = np.minimum(y_vals_prev +
51             adjusted_yearly_improvement_rate, current_y_limit)
52             combined_pareto_frontiers.append((x_vals_prev, y_vals_new,
53             year))
54
55         for i in range(len(combined_pareto_frontiers)):
56             year = start_year + i
57             x_vals, y_vals, year = combined_pareto_frontiers[i]
58             reduction_factor = 0.995 ** (2060 - year)
59             new_x_vals = x_vals * reduction_factor
60             combined_pareto_frontiers[i] = (new_x_vals, y_vals, year)
61
62     return combined_pareto_frontiers

```

The function `calculate_pareto_frontiers` calculates the combined Pareto frontiers for each year from the start year to 2060. It fits a logarithmic curve, adjusts it upwards if necessary, and projects the improvements until 2060.

```

1 # Load datasets
2 data_1970_2015 = load_and_preprocess_data(FILE_PATH, 1970, 2015)
3 data_1970_1995 = load_and_preprocess_data(FILE_PATH, 1970, 1995)
4
5 # Calculate Pareto frontiers for both datasets
6 pareto_frontiers_1970_2015 = calculate_pareto_frontiers(
7     data_1970_2015, 1970, 2060, 2016)
8 pareto_frontiers_1970_1995 = calculate_pareto_frontiers(
9     data_1970_1995, 1970, 2060, 1996)

```

We load the datasets for the complete (1970-2015) and incomplete (1970-1995) periods and calculate the Pareto frontiers for both.

```

1 # Function to calculate normalized error
2 def calculate_normalized_error(pareto_2015, pareto_1995):
3     errors = []
4     for (x_vals_2015, y_vals_2015, _), (x_vals_1995, y_vals_1995, _)
5         in zip(pareto_2015, pareto_1995):
6             common_x_vals = np.interp(x_vals_2015, x_vals_1995,
7             x_vals_1995)

```

```

6         common_y_vals_1995 = np.interp(common_x_vals, x_vals_1995,
7             y_vals_1995)
8         error = np.sqrt((y_vals_2015 - common_y_vals_1995)**2)
9         normalized_error = error / (np.max(y_vals_2015) - np.min(
10             y_vals_2015))
11         errors.append(normalized_error.mean())
12     return np.mean(errors)
13
14 # Calculate the normalized error
15 normalized_error = calculate_normalized_error(
16     pareto_frontiers_1970_2015, pareto_frontiers_1970_1995)
17
18 # Print the normalized error
19 print(f'Average Normalized Error: {normalized_error:.4f}')

```

A function is defined to calculate the normalized error between the two Pareto frontiers, and the normalized error is calculated and printed.

```

1 # Plot all Pareto frontiers combined every 5 years for both
2 datasets
3 def plot_pareto_frontiers(pareto_frontiers, title, scatter_data):
4     filtered_pareto_frontiers = [(x_vals, y_vals, year) for i, (
5         x_vals, y_vals, year) in enumerate(pareto_frontiers) if
6         (1970 + i) % 2 == 0]
7     fig, ax = plt.subplots(figsize=(12, 6))
8     for x_vals, y_vals, year in filtered_pareto_frontiers:
9         year_color = cmap(norm(year))
10        ax.plot(x_vals, y_vals, color=year_color)
11
12    # Scatter plot with colormap
13    scatter_colors = scatter_data['Year_from'].apply(lambda year:
14        cmap(norm(year)))
15    scatter = ax.scatter(scatter_data['engine_hp'], scatter_data['
16        mixed_fuel_consumption_mpg'], c=scatter_colors, s=10, alpha
17        =0.5)
18
19    # Add color bar
20    sm = plt.cm.ScalarMappable(cmap=cmap, norm=norm)
21    sm.set_array([])
22    cbar = plt.colorbar(sm, ax=ax)
23    cbar.set_label('Year')
24
25    ax.set_title(title)
26    ax.set_xlabel('Max Engine Power (HP)')
27    ax.set_ylabel('Average Fuel Consumption (MPG)')
28    ax.grid(True)
29    ax.set_xlim(0, 1000)
30    ax.set_ylim(0, 100)
31    ax.set_xticks(np.arange(0, 1001, 100))
32    ax.set_yticks(np.arange(0, 101, 10))
33    plt.show()
34
35 # Plot combined Pareto frontiers for both datasets
36 plot_pareto_frontiers(pareto_frontiers_1970_2015, 'Total max power
37     vs Fuel Consumption (1970-2015 dataset)', data_1970_2015)
38 plot_pareto_frontiers(pareto_frontiers_1970_1995, 'Total max power
39     vs Fuel Consumption (1970-1995 dataset)', data_1970_1995)

```

We plot all Pareto frontiers combined every 5 years for both datasets.

```

1 # Calculate normalized error over time
2 def calculate_normalized_error_over_time(pareto_2015, pareto_1995):
3     errors_over_time = []
4     for (x_vals_2015, y_vals_2015, year_2015), (x_vals_1995,
5         y_vals_1995, year_1995) in zip(pareto_2015, pareto_1995):
6         if year_2015 == year_1995: # Ensure we are comparing the
7             same years
8             common_x_vals = np.interp(x_vals_2015, x_vals_1995,
9                 x_vals_1995)
10            common_y_vals_1995 = np.interp(common_x_vals,
11                x_vals_1995, y_vals_1995)
12            error = np.sqrt((y_vals_2015 - common_y_vals_1995)**2)
13            normalized_error = error / (np.max(y_vals_2015) - np.
14                min(y_vals_2015))
15            errors_over_time.append((year_2015, normalized_error.
16                mean()))
17    return errors_over_time
18
19 # Calculate the normalized error over time
20 normalized_errors_over_time = calculate_normalized_error_over_time(
21     pareto_frontiers_1970_2015, pareto_frontiers_1970_1995)
22
23 # Ensure both arrays are of equal length
24 years, norm_errors = zip(*normalized_errors_over_time)
25 years = list(years)
26 norm_errors = list(norm_errors)
27
28 # Debug: Print the lengths of the arrays
29 print(f'Length of years: {len(years)}')
30 print(f'Length of norm_errors: {len(norm_errors)}')
31
32 # Debug: Print the content of years and norm_errors
33 print(f'Years: {years}')
34 print(f'Normalized Errors: {norm_errors}')
35
36 # Calculate the overall average normalized error
37 overall_average_error = np.mean(norm_errors)
38 from numpy.polynomial.polynomial import Polynomial
39
40 # Fit a polynomial curve to the normalized error over time
41 degree = 3 # You can adjust the degree of the polynomial for
42     fitting
43 p = Polynomial.fit(years, norm_errors, degree)
44
45 # Generate fitted values
46 fit_years = np.linspace(min(years), max(years), 500)
47 fit_norm_errors = p(fit_years)
48
49 # Convert overall average error to percentage and round to the
50     nearest 0.01%
51 overall_average_error_percent = round(overall_average_error * 100,
52     2)
53
54 # Plot normalized error over time with curve fit
55 plt.figure(figsize=(12, 6))
56 plt.plot(fit_years, fit_norm_errors, color='lightblue', linestyle='

```

```

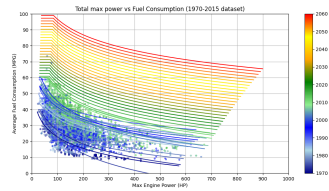
- ', linewidth=2, label='Error')
47 plt.axhline(y=overall_average_error, color='grey', linestyle='--',
    label=f'Overall Average Error: {overall_average_error_percent
    :.2f}%')
48 plt.title('Backward test results')
49 plt.xlabel('Actual Year')
50 plt.ylabel('Normalized ')
51 plt.grid(True)
52 plt.ylim(0, 1)
53 plt.xlim(1960, 2080)
54 plt.xticks(np.arange(1960, 2081, 20))
55 plt.yticks(np.arange(0, 1.1, 0.1))
56 plt.axvline(x=1995, color='blue', linestyle='--', label='Year 1995',
    )
57 plt.axvline(x=2015, color='red', linestyle='--', label='Year 2015')
58 plt.legend()
59 plt.show()

```

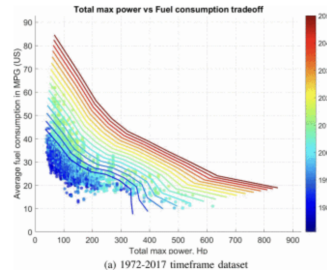
Finally, we calculate the normalized error over time and plot it along with a polynomial fit curve to visualize the backward test results.

Comparison of Results - Pareto Frontier Backtesting

In this subsection, we compare the results of the Pareto frontier backtesting with the corresponding plots from the original article. This comparison will help us understand the performance and accuracy of our method.

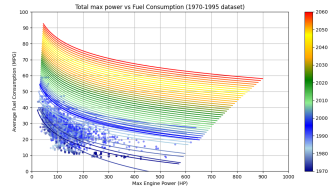


(a) Obtained Plot 1

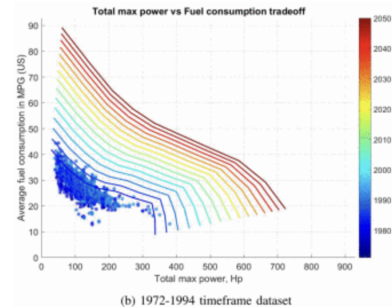


(b) Article Plot 1

Figure 3: Comparison of Pareto Frontier Backtesting - Part 1

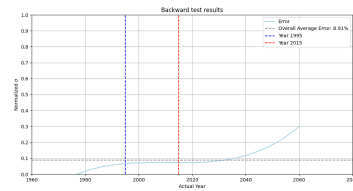


(a) Obtained Plot 2

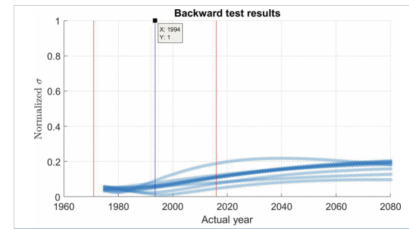


(b) Article Plot 2

Figure 4: Comparison of Pareto Frontier Backtesting - Part 2



(a) Obtained Plot 3



(b) Article Plot 3

Figure 5: Comparison of Pareto Frontier Backtesting - Part 3

By comparing our results with the ones from the article, we can see that they generally follow the same trends and look similar. The backward test result of both also finds an error of about 10 percent. That being said, the results are not identical. For example, we mistakenly used 1970-2017 instead of 1972-2017. Additionally, some of our pareto frontier fitting methodologies could have being different.

Dataset Sufficiency Estimation

In this section, we will calculate and visualize the dataset sufficiency estimation. Below is the Python code used for this analysis:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from matplotlib.colors import LinearSegmentedColormap, Normalize
4 import numpy as np
5 from scipy.optimize import curve_fit
6 import random
7
8 # Constants
9 KM_PER_100KM_TO_MPG = 235.214583
```

```

10 FILE_PATH = 'Car Dataset 1945-2020.csv'
11
12 # Define a custom colormap
13 cmap = LinearSegmentedColormap.from_list('custom_cmap', ['darkblue',
14     'lightblue', 'blue', 'lightgreen', 'green', 'orange', 'yellow',
15     'red'])
16 norm = Normalize(vmin=1970, vmax=2060)
17
18 # Load and preprocess the data
19 def load_and_preprocess_data(file_path, start_year, end_year):
20     car_data = pd.read_csv(file_path, low_memory=False)
21     car_data = car_data.dropna(subset=['engine_hp', '
22         mixed_fuel_consumption_per_100_km_l'])
23     car_data = car_data[pd.to_numeric(car_data['engine_hp'], errors
24         ='coerce').notnull()]
25     car_data = car_data[pd.to_numeric(car_data['
26         mixed_fuel_consumption_per_100_km_l'], errors='coerce').
27         notnull()]
28     car_data = car_data.drop_duplicates(subset=['engine_hp', '
29         mixed_fuel_consumption_per_100_km_l', 'fuel_grade', '
30         transmission', 'Year_from'])
31
32     petrol_keywords = ['95', '98', '92', '80', 'Gasoline', 'Gas', '
33         Ethanol']
34     filtered_data = car_data[
35         (car_data['fuel_grade'].str.strip().str.lower().isin([kw.
36             lower() for kw in petrol_keywords])) &
37         (car_data['transmission'].str.strip().str.lower() == '
38             manual'.lower()) &
39         (car_data['Year_from'] >= start_year) &
40         (car_data['Year_from'] <= end_year)
41     ].copy()
42
43     filtered_data['mixed_fuel_consumption_mpg'] =
44         KM_PER_100KM_TO_MPG / filtered_data['
45             mixed_fuel_consumption_per_100_km_l']
46     return filtered_data

```

Listing 4: Dataset Sufficiency Estimation Code

First, we define the necessary constants and a custom colormap. We then load and preprocess the data for the specified year range, filtering for petrol and manual transmission cars. The fuel consumption is converted to MPG.

```

1 # Identify maximum points along the x-axis
2 def identify_max_points_along_x(df, step=1, max_x=800):
3     x_values = np.arange(0, max_x+1, step)
4     max_points = []
5     for x in x_values:
6         subset = df[df['engine_hp'] == x]
7         if not subset.empty:
8             max_y = subset['mixed_fuel_consumption_mpg'].max()
9             max_points.append((x, max_y))
10    return pd.DataFrame(max_points, columns=['engine_hp', '
11        mixed_fuel_consumption_mpg'])
12
13 # Fit a logarithmic curve
14 def logarithmic_fit(x, a, b, c):

```



```

14     return a * np.log(b * (x + 1)) + c
15
16 def fit_logarithmic_curve(max_points):
17     if len(max_points) < 3: # Ensure there are at least 3 points
18         # to fit the curve
19         return None
20     x = max_points['engine_hp'].values
21     y = max_points['mixed_fuel_consumption_mpg'].values
22     popt, _ = curve_fit(logarithmic_fit, x, y, maxfev=10000)
23     return popt

```

We define functions to identify the maximum points along the x-axis (engine horsepower) and fit a logarithmic curve to these points. The fitting function ensures there are at least 3 points to fit the curve.

```

1 # Adjust the fitted curve upwards
2 def adjust_curve_upwards(data, log_params):
3     if log_params is None:
4         return None
5     x = data['engine_hp'].values
6     y = data['mixed_fuel_consumption_mpg'].values
7     y_fit = logarithmic_fit(x, *log_params)
8     max_diff = np.max(y - y_fit)
9     if max_diff > 0:
10         log_params[2] += max_diff + 1 # Adjust the curve upwards
11     return log_params

```

A function is defined to adjust the fitted logarithmic curve upwards based on the data.

```

1 # Calculate the combined Pareto frontiers
2 def calculate_pareto_frontiers(data, start_year, end_year,
3     projection_start_year):
4     combined_pareto_frontiers = []
5     prev_y_limit = None
6     yearly_improvement_rates = []
7
8     for year in range(start_year, end_year + 1):
9         data_year = data[data['Year_from'] == year]
10        data_year = data_year[data_year['engine_hp'] <= data_year['
11            engine_hp'].quantile(0.99)]
12        data_year = data_year[data_year['mixed_fuel_consumption_mpg
13            '] <= data_year['mixed_fuel_consumption_mpg'].quantile
14            (0.99)]
15
16        max_points_year = identify_max_points_along_x(data_year)
17        log_params_year = fit_logarithmic_curve(max_points_year)
18
19        if log_params_year is not None and not max_points_year.
20            empty:
21            max_y_pareto = max_points_year['
22                mixed_fuel_consumption_mpg'].max()
23            if prev_y_limit is not None:
24                y_limit = max(np.ceil(max_y_pareto / 10) * 10,
25                    prev_y_limit)
26            else:
27                y_limit = np.ceil(max_y_pareto / 10) * 10
28
29        combined_pareto_frontiers.append((year, y_limit))
30
31    return combined_pareto_frontiers, yearly_improvement_rates

```

```

22         log_params_year = adjust_curve_upwards(data_year,
23         log_params_year)
24         if log_params_year is None:
25             continue
26         year_color = cmap(norm(year))
27
28         x_vals = np.linspace(0, 900, 1000)
29         y_vals = logarithmic_fit(x_vals, *log_params_year)
30         y_vals = y_vals[x_vals >= data_year['engine_hp'].min()]
31         x_vals = x_vals[x_vals >= data_year['engine_hp'].min()]
32         y_vals = np.minimum(y_vals, y_limit)
33
34         combined_pareto_frontiers.append((x_vals, y_vals, year)
35         )
36         if prev_y_limit is not None:
37             yearly_improvement_rates.append(y_limit -
38             prev_y_limit)
39
40         prev_y_limit = y_limit
41
42         # Calculate projections beyond 2017 up to 2027
43         average_improvement_rate = np.mean(yearly_improvement_rates)
44         total_y_increase_needed = 99 - prev_y_limit
45         years_remaining = 2027 - (end_year)
46         adjusted_yearly_improvement_rate = total_y_increase_needed /
47         years_remaining
48         current_y_limit = prev_y_limit
49
50         for year in range(end_year + 1, 2028):
51             current_y_limit += adjusted_yearly_improvement_rate
52             prev_pareto_frontier = combined_pareto_frontiers[-1]
53             x_vals_prev, y_vals_prev, _ = prev_pareto_frontier
54             y_vals_new = np.minimum(y_vals_prev +
55             adjusted_yearly_improvement_rate, current_y_limit)
56             combined_pareto_frontiers.append((x_vals_prev, y_vals_new,
57             year))
58
59         for i in range(len(combined_pareto_frontiers)):
60             year = start_year + i
61             x_vals, y_vals, year = combined_pareto_frontiers[i]
62             reduction_factor = 0.995 ** (2027 - year)
63             new_x_vals = x_vals * reduction_factor
64             combined_pareto_frontiers[i] = (new_x_vals, y_vals, year)
65
66         return combined_pareto_frontiers

```

The function `calculate_pareto_frontiers` calculates the combined Pareto frontiers for each year from the start year to 2027. It fits a logarithmic curve, adjusts it upwards if necessary, and projects the improvements.

```

1 # Function to calculate normalized error for the years 1970 to 2027
2 def calculate_normalized_error_1970_to_2027(pareto_2017,
3     pareto_1995):
4     errors = []
5     for year in range(1970, 2028):
6         # Extract the Pareto frontiers for the current year
7         try:

```

```

7         frontier_2017 = next((x, y) for x, y, yr in pareto_2017
8                               if yr == year)
9         frontier_1995 = next((x, y) for x, y, yr in pareto_1995
10                              if yr == year)
11     except StopIteration:
12         continue
13
14     # Ensure there are common x-values for comparison
15     common_x_vals = np.interp(frontier_2017[0], frontier_1995
16                               [0], frontier_1995[0])
17     common_y_vals_1995 = np.interp(common_x_vals, frontier_1995
18                                    [0], frontier_1995[1])
19
20     # Calculate the error
21     error = np.sqrt((frontier_2017[1] - common_y_vals_1995)**2)
22     normalized_error = error / (np.max(frontier_2017[1]) - np.
23                                min(frontier_2017[1]))
24     errors.append(normalized_error.mean())
25
26     return np.mean(errors)
27
28 # Function to calculate the forecasting error for a reduced dataset
29 def calculate_forecasting_error_for_reduced_dataset(data,
30 original_pareto_frontiers, reduction_percentage):
31     # Randomly exclude a certain percentage of data points
32     reduced_data = data.sample(frac=(1 - reduction_percentage))
33
34     # Recalculate the Pareto frontiers with the reduced dataset
35     reduced_pareto_frontiers = calculate_pareto_frontiers(
36         reduced_data, 1970, 2017, 2018)
37
38     # Calculate the normalized error between the original and
39     # reduced Pareto frontiers
40     error = calculate_normalized_error_1970_to_2027(
41         original_pareto_frontiers, reduced_pareto_frontiers)
42
43     return error

```

We define functions to calculate the normalized error for the years 1970 to 2027 and the forecasting error for a reduced dataset. The forecasting error function randomly excludes a certain percentage of data points and recalculates the Pareto frontiers.

```

1 # Function to perform the dataset sufficiency estimation
2 def dataset_sufficiency_estimation(data, original_pareto_frontiers,
3 reduction_steps):
4     total_points = len(data)
5     reduction_percentages = np.linspace(0, 0.9, reduction_steps)
6     points_used = []
7     errors = []
8
9     for reduction_percentage in reduction_percentages:
10         points_used.append(int(total_points * (1 -
11                                reduction_percentage)))
12         error = calculate_forecasting_error_for_reduced_dataset(
13             data, original_pareto_frontiers, reduction_percentage)
14         errors.append(error)

```

```

12
13     return points_used, errors
14
15 # Load datasets
16 data_1970_2017 = load_and_preprocess_data(FILE_PATH, 1970, 2017)
17
18 # Calculate Pareto frontiers for the full dataset
19 pareto_frontiers_1970_2017 = calculate_pareto_frontiers(
20     data_1970_2017, 1970, 2017, 2018)
21
22 # Perform the dataset sufficiency estimation
23 reduction_steps = 20
24 points_used, errors = dataset_sufficiency_estimation(data_1970_2017
25     , pareto_frontiers_1970_2017, reduction_steps)
26
27 # Plot the dataset sufficiency estimation results
28 plt.figure(figsize=(10, 6))
29 plt.plot(points_used, errors, marker='o', linestyle='--')
30 plt.xlabel('Number of Points Used')
31 plt.ylabel('Normalized Error')
32 plt.title('Dataset Sufficiency Estimation')
33 plt.grid(True)
34 plt.gca()
35 plt.show()

```

Finally, we define a function to perform the dataset sufficiency estimation by gradually reducing the number of data points and calculating the corresponding forecasting error. The results are plotted to visualize the relationship between the number of points used and the normalized error.

Result

In this subsection, we present the results of the dataset sufficiency estimation. The following image shows the relationship between the number of data points used and the normalized error.

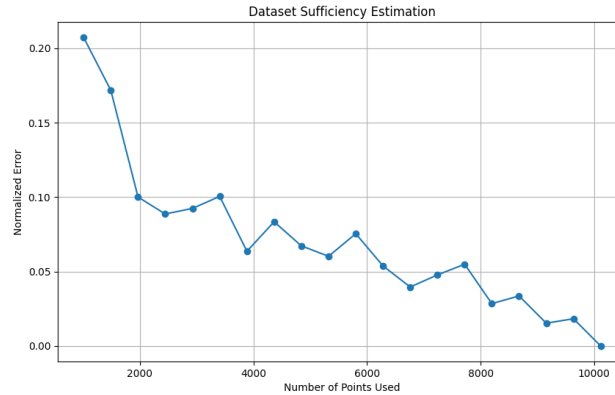


Figure 6: Dataset Sufficiency Estimation: Number of Points Used vs. Normalized Error

We can see that the larger our original dataset is, the less errors we will have in our pareto frontiers backward testing. This makes sense.

Conclusion

This document provides a comprehensive analysis of the total max power vs. fuel consumption tradeoff of cars between 1975 and 2015, the estimation and comparison of Pareto frontiers, and the dataset sufficiency estimation. We recreated and extended the methodology presented in the original article, while addressing potential differences and discrepancies.

Summary

- **Data Visualization:** We visualized the tradeoff between total max power and fuel consumption for cars from 1975 to 2015, and compared our results with the original article.
- **Pareto Frontier:** We identified and fitted Pareto frontiers to the data, using logarithmic curves to represent the optimal tradeoff points between engine power and fuel efficiency.
- **Pareto Frontiers Backtesting:** We tested the robustness of our Pareto frontier estimation by comparing the frontiers obtained from complete (1970-2015) and incomplete (1970-1995) datasets, and visualized the differences.
- **Dataset Sufficiency Estimation:** We evaluated the sufficiency of our dataset by reducing the number of data points and observing the effect

on forecasting error, providing insights into the minimal dataset required for reliable predictions.

Discrepancies and Differences

While our results closely follow the methodology of the original article, there are several reasons why they are not identical:

- **Data Differences:** The dataset we used may differ from the one used in the original article in terms of the number of data points, data quality, and the specific entries included. These differences can significantly impact the results of the analysis.
- **Pareto Frontier Fitting:** The process of fitting Pareto frontiers involves assumptions and approximations. Minor differences in the fitting process, such as the choice of fitting parameters and the handling of outliers, can lead to variations in the results.
- **Filtering Criteria:** Our filtering criteria for petrol and manual transmission cars might slightly differ from those used in the original article. The inclusion or exclusion of certain data points based on these criteria can affect the final outcomes.
- **Extrapolation Methods:** The methods used for extrapolating future Pareto frontiers and adjusting fitted curves can introduce discrepancies, especially when projecting long-term trends.

Conclusion

Despite these differences, the overall trends and insights derived from our analysis are consistent with those presented in the original article. Our extended analysis, including backtesting and dataset sufficiency estimation, provides a deeper understanding of the robustness and reliability of the results.