

MAT 328 Project: Malique Russell

Structure and Shape of the Data

The dataset is structured in a rectangular (tabular) format, consisting of rows and columns. It includes a mix of quantitative and qualitative data.

Quantitative Data:

- Extremely Low Income Units: Units with rents at 0–30% of the area median income
- Very Low Income Units: Rents at 31–50% of the area median income
- Low Income Units: Rents at 51–80% of the area median income
- Moderate Income Units: Rents at 81–120% of the area median income
- Middle Income Units: Rents at 121–165% of the area median income
- Other Income Units: Units reserved for building superintendents
- Counted Rental Units: Units counted under the Housing New York plan where assistance was provided to landlords
- Counted Homeownership Units: Units counted under the Housing New York plan where assistance was provided directly to homeowners
- All Counted Units: Total affordable units counted under the Housing New York plan
- Total Units: Sum of all units in the dataset
- Senior Units: Units specifically designated for senior households

Qualitative Data:

- Project ID: Unique identifier for each project
- Project Name: Name assigned by the Housing Preservation and Development (HPD).
- Program Group: Type of housing initiative
- Project Start Date: Date of project loan or agreement closure
- Project Completion Date: Date of the last building completion in a project
- Extended Affordability Only: Indicates whether the project qualifies for extended affordability
- Prevailing Wage Status: Specifies if the project adheres to prevailing wage requirements (e.g., Davis-Bacon Act)
- Planned Tax Benefit: Expected tax incentives associated with the project

Granularity of the Data:

The dataset has a low level of granularity, as each row represents aggregated unit data rather than individual housing units. A more granular dataset would provide detailed information at the unit level rather than summaries by category

Scope and Completeness of the Data

The dataset is well-suited for analyzing affordable housing trends in New York City. However, its scope is too broad for hyper-localized questions (e.g., borough-specific trends) and too narrow for state-wide analysis

Temporality of the Data

The dataset spans eight years, covering January 1, 2014, to December 31, 2021. It is managed by the Department of Housing Preservation and Development (HPD) and was last updated on March 3, 2025

Faithfulness of the Data

The dataset appears highly reliable, as it is compiled by a reputable city agency with direct oversight and access to housing records, ensuring accuracy and completeness

```
In [3]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
import statsmodels.formula.api as smf
from collections import Counter
```

```
In [4]: affordable_housing = pd.read_csv("Affordable_Housing_Production_by_Project.csv")
affordable_housing["Project Completion Date"] = pd.to_datetime(affordable_housing["
affordable_housing = affordable_housing.sort_values(by='Project Completion Date', a
affordable_housing.head()
```

Out[4]:

	Project ID	Project Name	Program Group	Project Start Date	Project Completion Date	Extended Affordability Only	Prevailing Wage Status
549	55759	CONFIDENTIAL	CONFIDENTIAL	01/03/2014	2014-01-03	No	No Prevailing Wage
523	55647	CONFIDENTIAL	CONFIDENTIAL	01/07/2014	2014-01-07	No	No Prevailing Wage
555	55773	CONFIDENTIAL	CONFIDENTIAL	01/10/2014	2014-01-10	No	No Prevailing Wage
641	57341	CONFIDENTIAL	CONFIDENTIAL	01/10/2014	2014-01-10	No	No Prevailing Wage
533	55697	CONFIDENTIAL	CONFIDENTIAL	01/14/2014	2014-01-14	No	No Prevailing Wage

```
In [5]: # Dropping Incomplete projects
complete_projects = affordable_housing.dropna(subset=['Project Completion Date'])
complete_projects.reset_index(drop=True, inplace=True)
complete_projects.head()
```

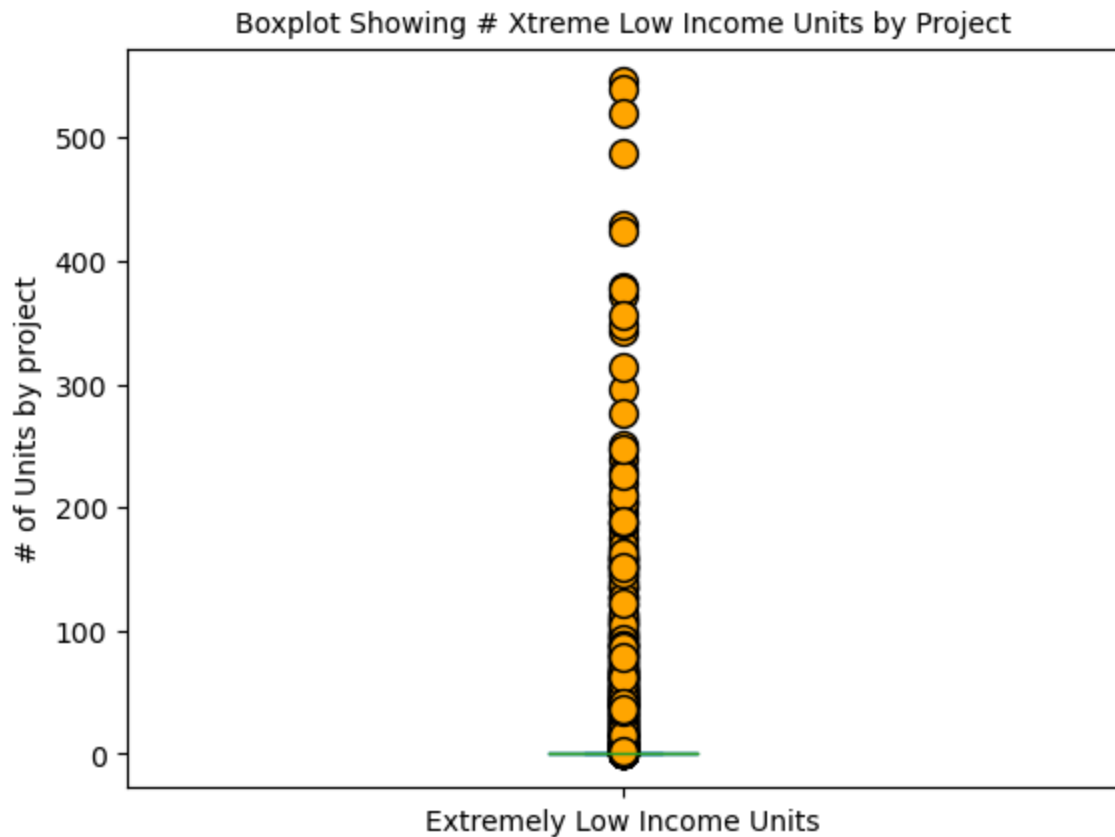
Out[5]:

	Project ID	Project Name	Program Group	Project Start Date	Project Completion Date	Extended Affordability Only	Prevailing Wage Status
0	55759	CONFIDENTIAL	CONFIDENTIAL	01/03/2014	2014-01-03	No	Non Prevailing Wage
1	55647	CONFIDENTIAL	CONFIDENTIAL	01/07/2014	2014-01-07	No	Non Prevailing Wage
2	55773	CONFIDENTIAL	CONFIDENTIAL	01/10/2014	2014-01-10	No	Non Prevailing Wage
3	57341	CONFIDENTIAL	CONFIDENTIAL	01/10/2014	2014-01-10	No	Non Prevailing Wage
4	55697	CONFIDENTIAL	CONFIDENTIAL	01/14/2014	2014-01-14	No	Non Prevailing Wage

```
In [6]: xtreme = complete_projects["Extremely Low Income Units"]
very = complete_projects["Very Low Income Units"]
low = complete_projects["Low Income Units"]
moderate = complete_projects["Moderate Income Units"]
middle = complete_projects["Middle Income Units"]
other = complete_projects["Other Income Units"]
owned = complete_projects["Counted Homeownership Units"]
total = complete_projects["All Counted Units"]
```

```
In [7]: xtreme.plot(kind = "box", flierprops=dict(marker='o', markersize=10, markerfacecolor='r'),
plt.ylabel("# of Units by project")
plt.title('Boxplot Showing # Xtreme Low Income Units by Project', fontsize = 10)
```

Out[7]: Text(0.5, 1.0, 'Boxplot Showing # Xtreme Low Income Units by Project')



This graph shows the distribution of completed extremely low-income units by project built in New York City from anuary 1, 2014 to December 30, 2025

Some notable deductions:

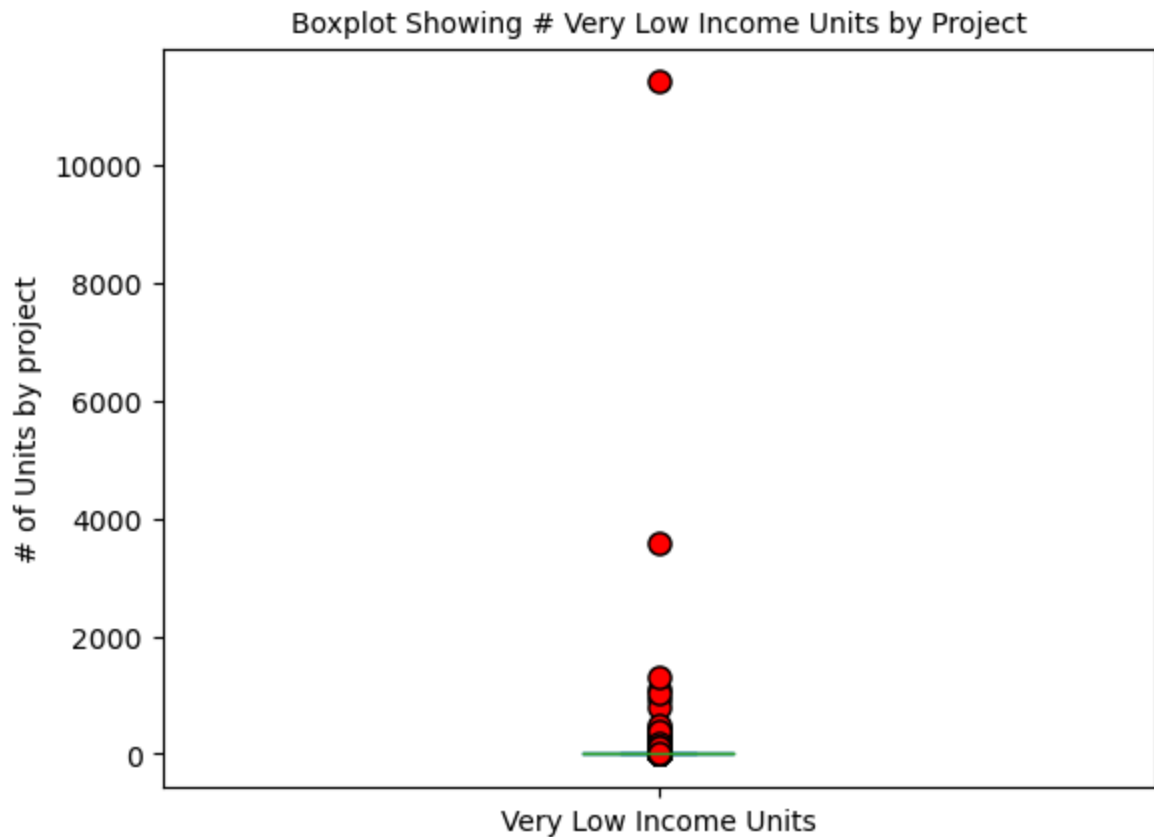
- The majority of projects had under 300 extremely low-income unit
- Most projects had less than 250 of these units

```
In [9]: # Percent of Completed Extremely Low Income Units
Xtreme =xtreme.sum()/total.sum()
xtreme_per = round(Xtreme*100,2)
xtreme_per
```

Out[9]: 16.84

```
In [10]: very.plot(kind = "box", flierprops=dict(marker='o', markersize = 8, markerfacecolor=
plt.ylabel("# of Units by project")
plt.title('Boxplot Showing # Very Low Income Units by Project', fontsize = 10)
```

Out[10]: Text(0.5, 1.0, 'Boxplot Showing # Very Low Income Units by Project')



This graph shows the distribution of completed very low-income units by project built in New York City from January 1, 2014 to December 30, 2025

Some notable deductions:

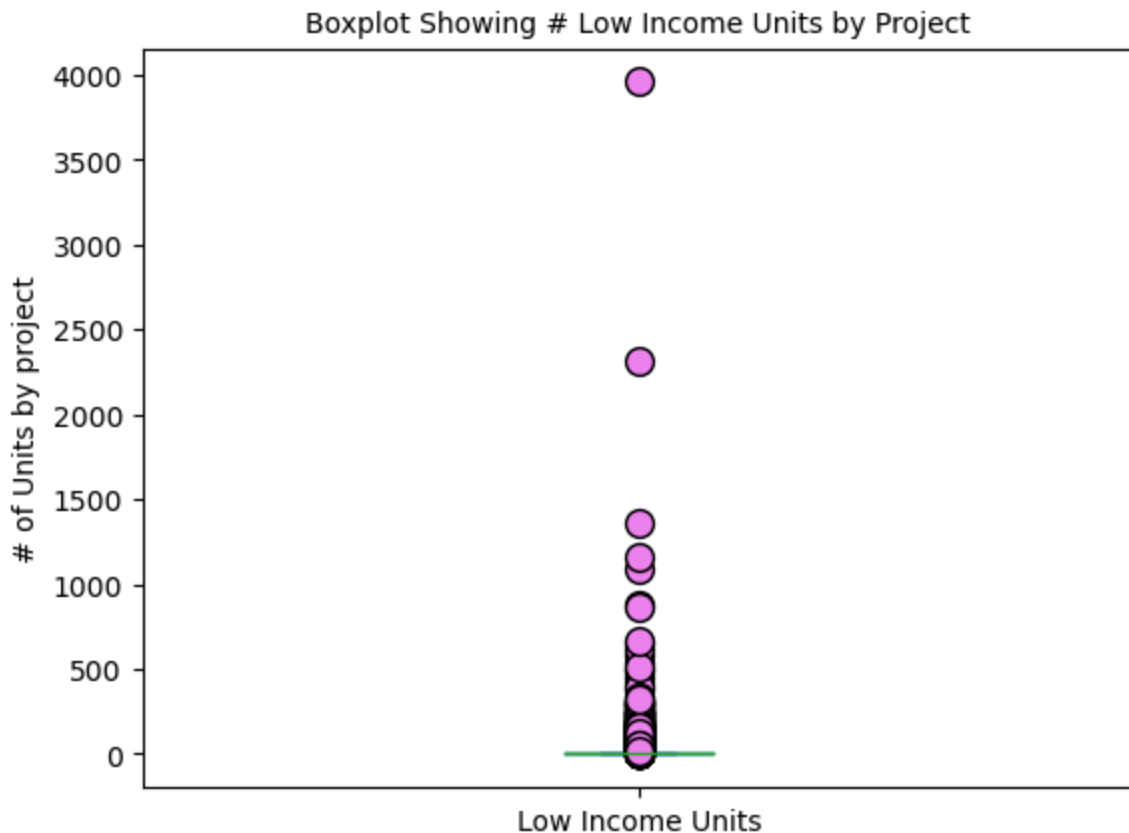
- The majority of projects had under 2000 very low-income units
- One project had 10,000+ of these units

```
In [12]: # Percent of Completed Very Low Income Units
Very = very.sum()/total.sum()
round(Very*100,2)
very_per = round(Very*100,2)
very_per
```

Out[12]: 24.88

```
In [13]: low.plot(kind = "box", flierprops=dict(marker='o', markersize=10, markerfacecolor =
plt.ylabel("# of Units by project")
plt.title('Boxplot Showing # Low Income Units by Project', fontsize = 10)
```

Out[13]: Text(0.5, 1.0, 'Boxplot Showing # Low Income Units by Project')



This graph shows the distribution of completed low-income units by project built in New York City from January 1, 2014 to December 30, 2025

Some notable deductions:

- The majority of projects had under 1000 low-income unit
- Most projects had less than 600 of these units

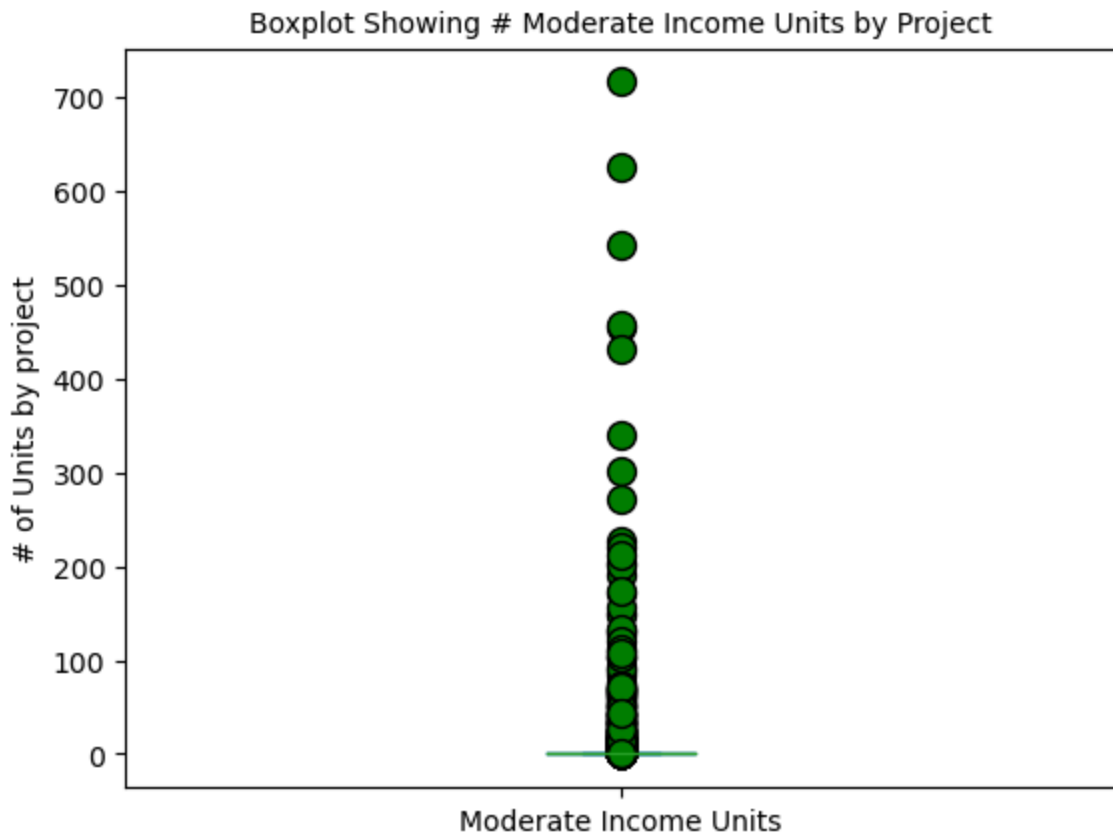
In [15]: *# Percent of Completed Low Income Units*

```
Low = low.sum()/total.sum()
round(Low*100,2)
low_per = round(Low*100,2)
low_per
```

Out[15]: 36.67

In [16]: `moderate.plot(kind = "box", flierprops=dict(marker='o', markersize=10, markerfaceco`
`plt.ylabel("# of Units by project")`
`plt.title('Boxplot Showing # Moderate Income Units by Project', fontsize = 10)`

Out[16]: Text(0.5, 1.0, 'Boxplot Showing # Moderate Income Units by Project')



This graph shows the distribution of completed moderate income units by project built in New York City from January 1, 2014 to December 30, 2025

Some notable deductions:

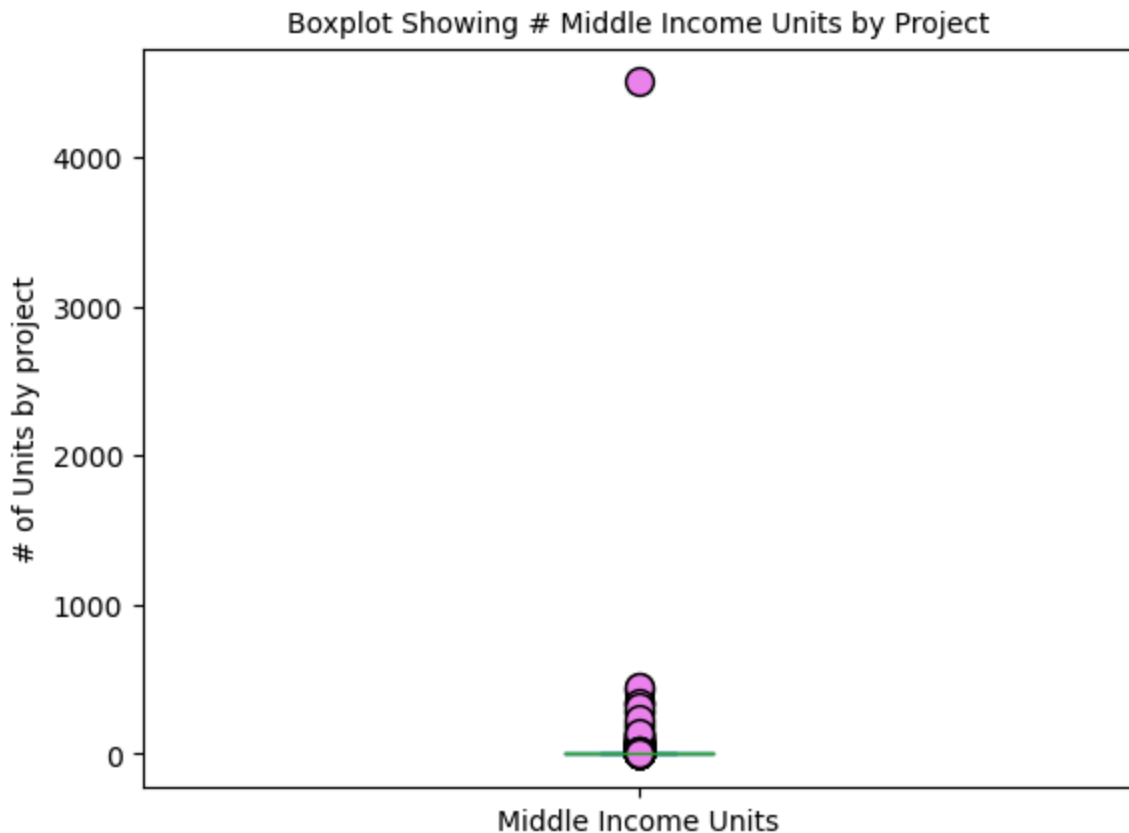
- The majority of projects had under 250 moderate income unit
- Most projects had less than 100 of these units

```
In [18]: # Percent of Completed Moderate Income Units
Moderate = moderate.sum()/total.sum()
round(Moderate*100,2)
moderate_per = round(Moderate*100,2)
moderate_per
```

Out[18]: 6.85

```
In [19]: middle.plot(kind = "box", flierprops=dict(marker='o', markersize=10, markerfacecolor=
plt.ylabel("# of Units by project")
plt.title('Boxplot Showing # Middle Income Units by Project', fontsize = 10)
```

Out[19]: Text(0.5, 1.0, 'Boxplot Showing # Middle Income Units by Project')



This graph shows the distribution of completed middle income units by project built in New York City from January 1, 2014 to December 30, 2025

Some notable deductions:

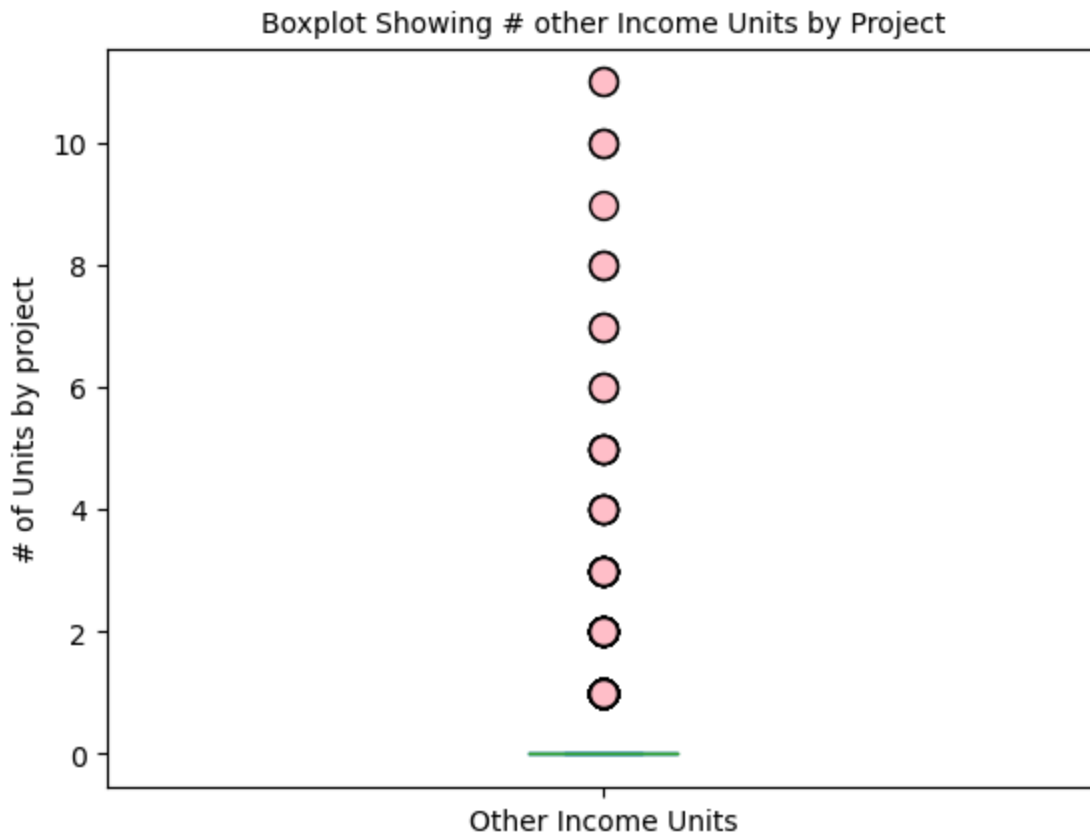
- The majority of projects had less than 150 middle low-income unit

```
In [21]: # Percent of Completed Middle Income Units
Middle = middle.sum()/total.sum()
round(Middle*100,2)
middle_per = round(Middle*100,2)
middle_per
```

Out[21]: 14.28

```
In [22]: other.plot(kind = "box", flierprops=dict(marker='o', markersize=10, markerfacecolor
plt.ylabel("# of Units by project")
plt.title('Boxplot Showing # other Income Units by Project', fontsize = 10)
```

Out[22]: Text(0.5, 1.0, 'Boxplot Showing # other Income Units by Project')



This graph shows the distribution of completed other income units by project built in New York City from January 1, 2014 to December 30, 2025

Some notable deductions:

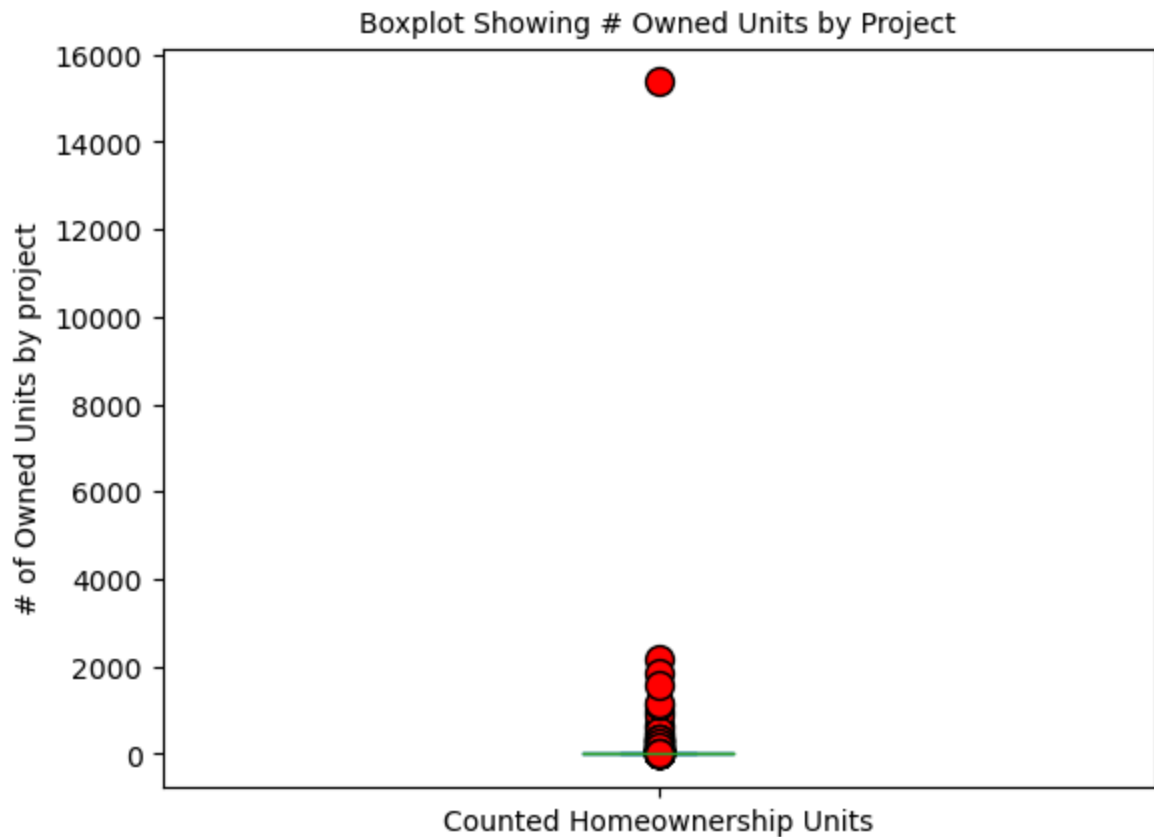
- These units seem to be fairly distributed

```
In [24]: # Percentage of Completed Other Income Units
Other = other.sum()/total.sum()
round(Other*100,2)
other_per = round(Other*100,2)
other_per
```

Out[24]: 0.47

```
In [25]: owned.plot(kind = "box" , flierprops=dict(marker='o', markersize=10, markerfacecolor='r'))
plt.ylabel("# of Owned Units by project")
plt.title('Boxplot Showing # Owned Units by Project', fontsize = 10)
```

Out[25]: Text(0.5, 1.0, 'Boxplot Showing # Owned Units by Project')



This graph shows the distribution of completed owned income units by project built in New York City from January 1, 2014 to December 30, 2025

Some notable deductions:

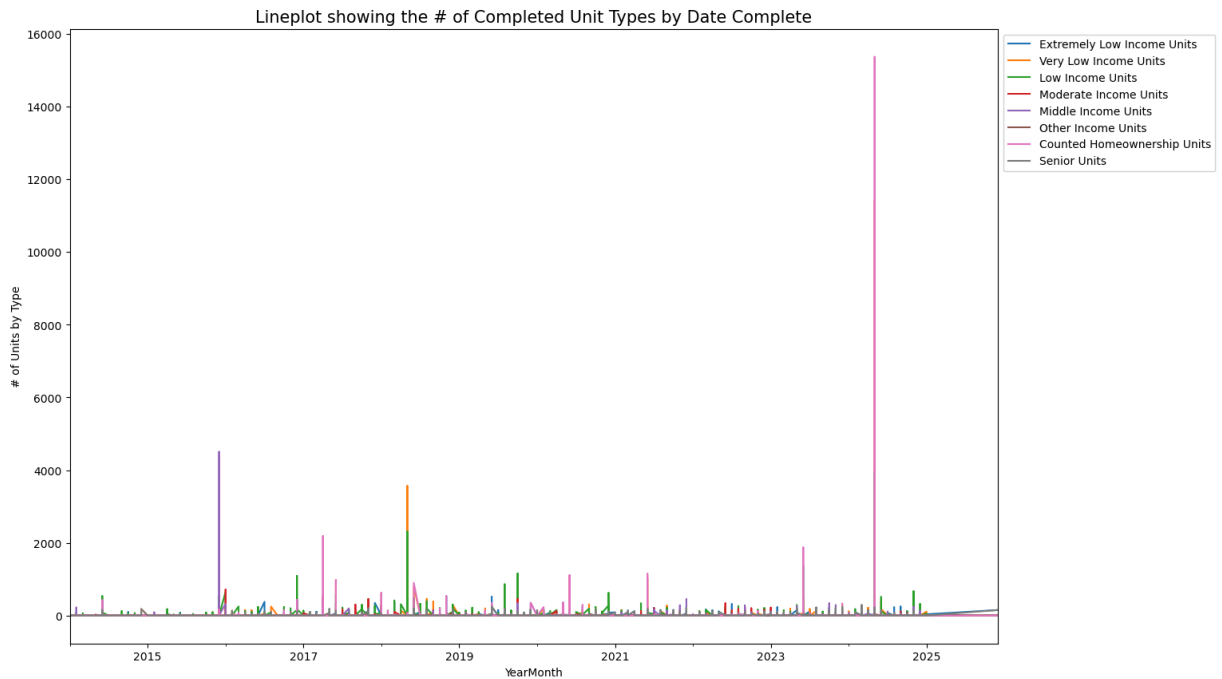
- The majority of projects had less than 2000 owned unit
- One project had 15,000+ units owned, which is an obvious outlier

In [27]: *# Percentage of Completed Owned Units*

```
Owned = owned.sum()/total.sum()
round(Owned*100,2)
owned_per = round(Owned*100,2)
owned_per
```

Out[27]: 18.93

```
In [28]: projects = complete_projects.drop(["Project ID", "Total Units", "All Counted Units"]
projects['YearMonth'] = projects['Project Completion Date'].dt.to_period('M')
projects.drop(["Project Completion Date"], axis = 1).plot(x="YearMonth", figsize =
plt.ylabel("# of Units by Type")
_=plt.title('Lineplot showing the # of Completed Unit Types by Date Complete', font
```



This line graph shows the number of completed units by type built in New York City from January 1, 2014 to December 30, 2025

Some key points shown in the graph are:

- Most projects had under 2000 completed units regardless of type
- At least 4 projects had over 2000 units completed which makes them outliers in the data
- The most units completed for a single project type fall under the home-owner category, completed after 2024

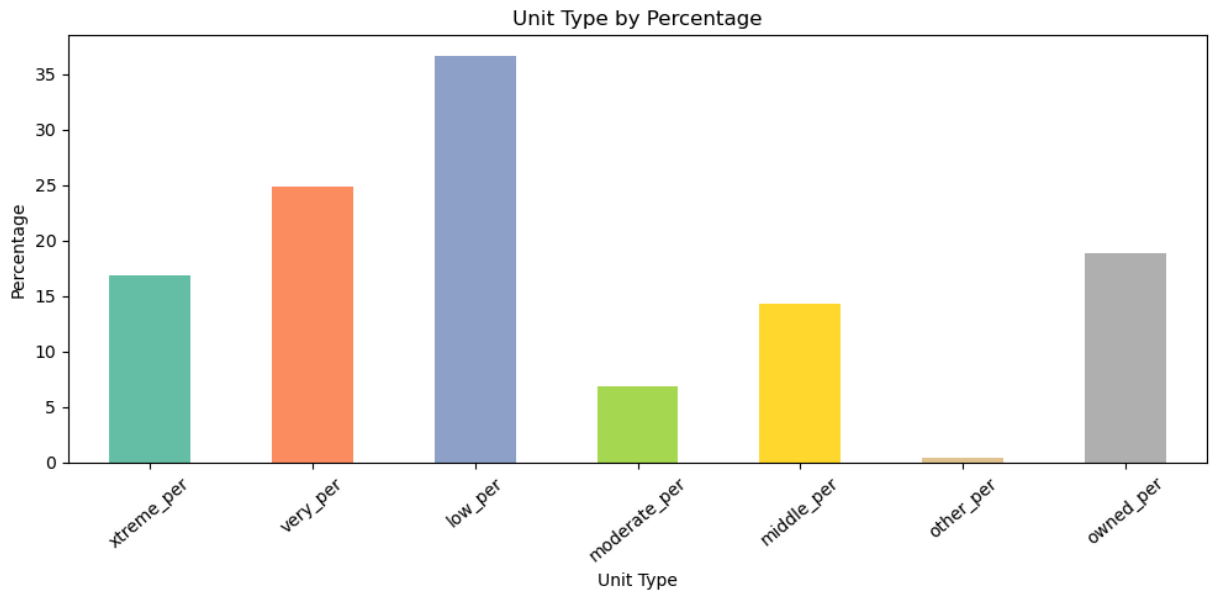
In [116...

```
data = {
    'Unit Type': ["xtreme_per", "very_per", "low_per", "moderate_per", "middle_per",
                  "other_per", "counted_per", "senior_per"],
    'Percentages': [16.84, 24.88, 36.67, 6.85, 14.28, 0.47, 18.93]}

df = pd.DataFrame(data)
colors = plt.cm.Set2(np.linspace(0, 1, len(data['Unit Type'])))

df.plot(kind='bar', x='Unit Type', y='Percentages', color=colors, legend=False, fig=fig)

plt.title("Unit Type by Percentage")
plt.ylabel("Percentage")
plt.xticks(rotation=40)
plt.tight_layout()
plt.show()
```



This bar graph shows the percentage of completed units by type built in New York City from January 1, 2014 to December 30, 2025

Some key points shown on the chart are:

- Low income units were the most built in New York City during the 10 year period
- Other income units were the least built in the period
- Units falling under the xtreme, very and low income categories account for the bulk of units built 78.39%
- Only 19% of units completed are owned