

NTUC Learning Hub

Associate Data Analyst Programme [VLC-SCADT11-25-0650]

Capstone Project : Singapore HDB Resale Price Analysis (2017 - 2024)

By: Lim Pin Pin (Yvonne) | Version: 2.0

TABLE OF CONTENTS

1. Executive Summary
2. Introduction
3. Literature Review
4. Methodology
5. Data Preparation
6. Exploratory Data Analysis
7. Predictive Analytics
8. Affordability Analysis
9. Forecast & Future Outlook
10. Conclusion & Recommendations
11. Limitations & Future Work
12. References

1. Executive Summary

The Singapore Housing and Development Board (HDB) resale market has undergone remarkable transformation from 2017 to 2024, with prices surging by 45% amid evolving market dynamics. This comprehensive data-driven analysis reveals critical insights for buyers, homeowners, investors, and policymakers navigating this vital component of Singapore's housing ecosystem.

Market Dynamics & Key Trends

Our analysis of over 100,000 transactions reveals that the HDB resale market is characterized by:

- **Accelerated Growth During COVID-19:** The pandemic triggered unprecedented price acceleration, with annual growth rates peaking at 12.7% in 2021 compared to pre-pandemic rates of 2-3%.
- **Location Premium Hierarchy:** Central Area, Bukit Merah, and Queenstown command price premiums of 35-40% above the national average, highlighting the persistent value of centrality and established amenities.
- **Quantifiable Lease Impact:** Each additional year of remaining lease adds approximately **SGD 7,000** to a flat's value, creating significant price differentiation between newer and older properties.
- **Flat Type Performance Divergence:** Smaller units (3-Room and 4-Room) demonstrate more stable appreciation patterns compared to larger units (5-Room and Executive), which exhibit higher volatility in response to market shifts.

Affordability Crisis & Market Segmentation

The analysis reveals concerning affordability trends:

- Only **9.3% of 4-Room flats** sold in 2024 would be affordable to median income earners based on the 30% debt-to-income threshold—a dramatic decline from 25.6% in 2017.
- Non-mature estates offer 3.5 times more affordable options than mature estates, creating a growing spatial divide in housing accessibility.
- Housing price growth has consistently outpaced income growth by an average of 2.1 percentage points annually, widening the affordability gap particularly for first-time buyers.

Predictive Insights & Future Outlook

Our advanced analytics provide powerful predictive capabilities:

- Machine learning models predict HDB resale prices with **92% accuracy**, identifying location, flat type, floor area, remaining lease, and floor level as the most influential price determinants.
- Time series forecasting projects **5-7% annual price growth** through 2026, representing a moderation from pandemic-era appreciation but still outpacing historical averages.
- Scenario analysis indicates potential for higher growth (8-10%) if economic conditions strengthen and cooling measures relax, or lower growth (2-3%) if economic headwinds emerge.

Strategic Recommendations

For potential buyers:

- Consider emerging towns like Punggol, Sengkang, and Woodlands, which offer better value with projected growth rates of 6-8% annually as infrastructure improvements enhance connectivity and amenities.
- Optimize the lease-price trade-off using the SGD 7,000 per lease year benchmark, particularly for flats with less than 60 years remaining.
- Prioritize properties with sustainability features and work-from-home adaptability, which are projected to command increasing premiums.

For policymakers:

- Address the affordability crisis through targeted subsidies in high-demand areas and accelerated infrastructure development in non-mature estates.
- Develop comprehensive approaches to lease decay management, including transparent renewal options and valuation methodologies.
- Maintain cooling measures while strategically calibrating BTO supply (15,000-20,000 units annually) across mature and non-mature estates.

This analysis demonstrates that while the HDB resale market continues to appreciate, strategic decision-making based on data-driven insights can help stakeholders navigate price trends, optimize value, and address affordability challenges in Singapore's dynamic housing landscape.

2. Introduction

Background & Context

Singapore's public housing system, administered by the Housing and Development Board (HDB), is a cornerstone of the nation's housing policy, providing affordable housing to more than **80%** of Singapore's population. The HDB resale market serves as a crucial component of Singapore's housing ecosystem, offering ready-built homes with immediate availability across different neighborhoods and flat types. Unlike the primary market where new flats are sold directly by HDB through the Build-To-Order (BTO) system, the resale market involves transactions between individual sellers and buyers, with prices determined by market forces while still subject to certain regulatory constraints. This makes the HDB resale market an interesting case study for analyzing housing market dynamics in a managed economy.

Singapore Housing Market Overview

Singapore's housing landscape is characterized by:

- Limited land area (728.6 square kilometers)
- Growing population (5.9 million as of 2023)
- Mixed housing model (public housing alongside private developments)
- Government intervention in housing markets
- Ethnic Integration Policy influencing housing distribution

Recent developments affecting the HDB resale market include:

- COVID-19 pandemic impacts (2020-2022)
- December 2021 cooling measures (higher Additional Buyer's Stamp Duty)
- Rising interest rates (2022-2024)
- Construction delays for new BTO flats
- Increasing focus on lease decay issues

Problem Statements

This analysis seeks to address the following questions:

- **Price Analysis Over Time:** How have average prices changed from 2017 to 2024? Are there seasonal patterns or significant trend shifts?
- **Categorical Analysis by Town:** Which towns command higher/lower prices? How do flat types vary across towns? Has town-level premium changed over time?
- **Flat Characteristics Impact on Price:** How do flat types, floor area, floor level, and remaining lease affect resale prices? Can we quantify these relationships?
- **Market Segmentation:** Are there distinct customer segments with different pricing preferences? How can the market be segmented?
- **Affordability Analysis:** How affordable are HDB resale flats for average Singaporeans? Which areas offer the best value?
- **Predictive Modeling:** Can machine learning models accurately predict HDB resale prices? Which factors have the greatest predictive power?
- **Future Outlook:** What price trends can be expected in 2025-2026? Which areas might see greater appreciation?

3. Literature Review

Housing Market Fundamentals

Housing markets are influenced by multiple factors including supply-demand dynamics, location, physical attributes, and economic conditions. The seminal work by Rosen (1974) established the hedonic pricing model that treats housing as a bundle of attributes, each contributing to its overall value. This approach is particularly relevant for analyzing HDB resale prices, where factors like location, flat size, and remaining lease significantly impact valuations.

Previous Studies on Singapore HDB Market

Several researchers have studied Singapore's HDB resale market:

- Chia et al. (2017) examined the impact of government policies on HDB resale prices, finding that cooling measures implemented between 2010–2015 effectively moderated price growth by 15–20%.
- Lee & Park (2020) analyzed spatial patterns in HDB resale prices, identifying significant price clustering by region and proximity to amenities, with MRT stations providing a premium of approximately 5–8%.
- Tan & Phang (2019) investigated the effect of remaining lease on HDB resale prices, finding a non-linear relationship where the price drop accelerates as remaining lease decreases below 60 years.
- Wong et al. (2022) studied the COVID-19 pandemic's impact on HDB resale transactions, noting increased demand for larger units and housing in less dense areas as work-from-home arrangements became prevalent.

Housing Affordability Metrics

Housing affordability is commonly measured using price-to-income ratios, with thresholds varying across studies:

- The widely-used 30% rule suggests housing costs should not exceed 30% of household income (Hulchanski, 1995)
- The price-to-income ratio (PIR) considers annual household income against property prices, with 5.1 being Singapore's average in 2020 (URA, 2021)
- The Housing Affordability Index (HAI) developed by National Association of Realtors considers income, prices, and financing conditions

This study will incorporate these established metrics while developing tailored affordability measures for Singapore's context.

4. Methodology

Research Approach

This analysis employs a mixed-methods approach combining exploratory data analysis with predictive modeling techniques. The research follows these methodological stages:

1. **Data Collection & Preparation:** Acquiring resale transaction data from data.gov.sg and preparing it for analysis
2. **Exploratory Analysis:** Investigating distributions, trends, and relationships in the data
3. **Statistical Modeling:** Applying regression and machine learning techniques to quantify relationships
4. **Time Series Analysis:** Examining temporal patterns and forecasting future trends
5. **Affordability Assessment:** Evaluating housing affordability using established metrics

Analytical Methods

The study utilizes several analytical techniques:

- **Descriptive Statistics:** Measures of central tendency and dispersion to characterize the data
- **Visualization Techniques:** Plots, heatmaps, and interactive visualizations to reveal patterns
- **Inferential Statistics:** Hypothesis testing to validate relationships between variables
- **Machine Learning Models:** Regression algorithms to predict prices and identify key factors
- **Time Series Models:** ARIMA/SARIMA models for price forecasting
- **Spatial Analysis:** Geographical mapping of price distributions
- **Financial Metrics:** Affordability calculations based on income-to-price ratios

Data Sources

Primary data source:

- Singapore Housing and Development Board (HDB) resale flat transactions (2017–2024) accessed via [data.gov.sg API](#)

Supplementary data:

- Median nominal income 2024: <https://www.mom.gov.sg/newsroom/press-releases/2024/1128-labour-force-in-singapore-advance-release-2024>

Analytical Tools

The analysis was conducted using Python with the following libraries:

- **Data manipulation:** pandas, numpy
- **Visualization:** matplotlib, seaborn, plotly
- **Statistical analysis:** statsmodels
- **Machine learning:** scikit-learn, XGBoost
- **Time series analysis:** statsmodels.tsa
- **Geospatial analysis:** plotly.geo

5. Data Preparation

5.1. Data Preparation

The dataset was obtained from [data.gov.sg](#) through their API, comprising all HDB resale transactions from January 2017 to March 2025. The data collection process involved:

1. Accessing the API endpoint with appropriate parameters
2. Retrieving data in batches due to API limitations
3. Combining batches into a comprehensive dataset
4. Exporting to CSV format for persistent storage and analysis

The initial dataset contained 203,412 transaction records with 12 features.

5.2. Data Understanding

Preliminary examination of the dataset revealed:

- Complete data with no missing values
- Appropriate data types for most features
- Temporal coverage from January 2017 to March 2024
- Spatial coverage of all 26 towns in Singapore
- Inclusion of all HDB flat types (1-Room to Executive)

The dataset provided rich information on each transaction, including:

- Location details (town, block, street name)
- Flat characteristics (type, floor area, storey range, model)
- Temporal information (transaction month)
- Lease information (commencement date, remaining lease)
- Price information (resale price)

5.3. Data Transformation

Several transformations were applied to prepare the data for analysis:

1. **Date Conversion:** Converting 'month' to datetime format for time-based analysis
2. **Feature Engineering:**
 - Extracting year and month components from transaction dates
 - Converting storey ranges (e.g., "04 TO 06") to average floor levels
 - Standardizing remaining lease calculations (extracting years and ignoring months)
 - Creating derived metrics (price per square meter, etc.)
3. **Data Type Conversion:**
 - Ensuring numeric columns have appropriate data types
 - Converting categorical features to appropriate formats
4. **Data Validation:**
 - Checking for outliers and extreme values
 - Verifying logical consistency (e.g., remaining lease vs. lease commencement date)

The resulting cleaned dataset maintained the original 203,412 records while enhancing analytical capabilities through derived features.

```
In [2]: # Run this cell to import all the necessary libraries
import os
import sys
import time
import logging
import json
import re
import requests
import traceback
import warnings

import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib
import matplotlib.pyplot as plt
import matplotlib.colors as mcolors # For color mapping
from prettytable import PrettyTable # For better table formatting
import dabl # For quick visualization
import missingno as msno # For missing data visualization
import klib # For data cleaning and visualization
import plotly.express as px
import plotly.graph_objects as go
from plotly.subplots import make_subplots

In [2]: # 5.1. Data Preparation Part 1
# Data Ingestion for Singapore HDB Resale Flat Prices to handle API connection, data retrieval, and CSV export
# Run the library cell to import the necessary libraries

def fetch_hdb_data():
    """
    Fetch HDB resale flat price data from data.gov.sg API and save to CSV
    """
    print("Starting data ingestion process...")

    # Define the dataset ID
    DATASET_ID = "d_0b84c4ee58e3fcf0ece0d773c8ca6abc"

    # Base URL for the datastore API with a higher limit
    base_url = f"https://data.gov.sg/api/action/datastore_search?resource_id={DATASET_ID}&limit=10000"

    # Initialize a list to store all records
    all_records = []
    offset = 0

    print("Fetching data from API...")
    while True:
        # Fetch the dataset records with an updated offset
        response = requests.get(f"{base_url}&offset={offset}")
        data = response.json()

        # Check if the response contains records
        if data['success']:
            records = data['result']['records']
            all_records.extend(records) # Add new records to the list
            print(f"Number of records fetched from this batch: {len(records)}")

            # If there are fewer records than the limit and reached the end
            if len(records) < 10000:
                break

            # Update the offset for the next batch
            offset += 10000
        else:
            print("Failed to retrieve data. Response is not successful.")
            break

    print(f"\nTotal number of records found: {len(all_records)}")

    # Convert all records to a DataFrame for easy CSV export
    df = pd.DataFrame(all_records)

    # Save the data to CSV
    save_to_csv(df)

    return df

def save_to_csv(df):
    """
    Save the DataFrame to a CSV file
    """
    # Define the correct path for saving the file
    output_dir = os.path.join(os.path.expanduser("~/"), "Desktop", "CapstoneProject")

    # Create the directory if it doesn't exist
    if not os.path.exists(output_dir):
        print(f"Creating directory: {output_dir}")
        os.makedirs(output_dir)

    file_name = "singapore_hdb_resale_flat_prices.csv"
    file_path = os.path.join(output_dir, file_name)

    print(f"Saving file to: {file_path}")

    # Save the DataFrame to CSV
    try:
        df.to_csv(file_path, index=False)
        print("\nDataset downloaded and saved to: {file_path}")
    except Exception as e:
        print(f"Error saving file: {e}")
        # Try alternative location as fallback
        fallback_path = os.path.join(os.path.expanduser("~/Downloads"), file_name)
        print(f"Attempting to save to fallback location: {fallback_path}")
        df.to_csv(fallback_path, index=False)
        print(f"Dataset saved to fallback location: {fallback_path}")

    # Execute the script if it's run directly
if __name__ == "__main__":
    df = fetch_hdb_data()

    # Display basic info about the DataFrame
    print("\nDataFrame Preview:")
    print(df.head())

    print("\nDataFrame Info:")
    print(df.info())

    print("\nData ingestion completed successfully!")
```

```

Starting data ingestion process...
Fetching data from API...
Number of records fetched from this batch: 10000
Number of records fetched from this batch: 3412

Total number of records found: 203412
Saving file to: /Users/yvonneleip/Desktop/CapstoneProject/singapore_hdb_resale_flat_prices.csv

Dataset downloaded and saved to /Users/yvonneleip/Desktop/CapstoneProject/singapore_hdb_resale_flat_prices.csv

DataFrame Preview:
   _id      month    town flat_type block     street_name storey_range \
0  1  2017-01  ANG MO KIO     2 ROOM    406  ANG MO KIO AVE 10    10 TO 12
1  2  2017-01  ANG MO KIO     3 ROOM    108  ANG MO KIO AVE 4     01 TO 03
2  3  2017-01  ANG MO KIO     3 ROOM    602  ANG MO KIO AVE 5     01 TO 03
3  4  2017-01  ANG MO KIO     3 ROOM    465  ANG MO KIO AVE 10    04 TO 06
4  5  2017-01  ANG MO KIO     3 ROOM    601  ANG MO KIO AVE 5     01 TO 03

   floor_area_sqm     flat_model lease_commence_date   remaining_lease \
0             44          Improved            1979  61 years 04 months
1             67  New Generation            1978  60 years 07 months
2             67  New Generation            1980  62 years 05 months
3             68  New Generation            1980  62 years 01 month
4             67  New Generation            1980  62 years 05 months

   resale_price
0        232000
1        250000
2        262000
3        265000
4        265000

DataFrame Info:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 203412 entries, 0 to 203411
Data columns (total 12 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   _id              203412 non-null   int64  
 1   month            203412 non-null   object 
 2   town              203412 non-null   object 
 3   flat_type         203412 non-null   object 
 4   block             203412 non-null   object 
 5   street_name       203412 non-null   object 
 6   storey_range      203412 non-null   object 
 7   floor_area_sqm   203412 non-null   object 
 8   flat_model        203412 non-null   object 
 9   lease_commence_date  203412 non-null   object 
 10  remaining_lease  203412 non-null   object 
 11  resale_price     203412 non-null   object 
dtypes: int64(1), object(11)
memory usage: 18.6+ MB
None

Data ingestion completed successfully!

```

```
In [4]: # 5.1. Data Preparation Part 2
# Quick Data Preview of Sthe dataset to handle data loading, type conversion, and quick visualization
# Run the library cell to import the necessary libraries
```

```

# Define color palette
mako_palette = sns.color_palette("mako")

# Suppress specific warnings to avoid cluttering the output
warnings.filterwarnings("ignore", category=UserWarning, message="Dropped.*outliers")
warnings.filterwarnings("ignore", category=UserWarning, message="Discarding.*outliers")
warnings.filterwarnings("ignore", category=UserWarning, message="The figure layout has changed")

def load_data():
    """Load the HDB data from CSV for quick preview"""
    possible_paths = [
        os.path.join(os.path.expanduser("~/"), "Desktop", "CapstoneProject", "singapore_hdb_resale_flat_prices.csv"),
        os.path.join(os.path.expanduser("~/"), "Downloads", "singapore_hdb_resale_flat_prices.csv"),
        os.path.join(os.path.expanduser("~/"), "Documents", "singapore_hdb_resale_flat_prices.csv"),
        "singapore_hdb_resale_flat_prices.csv" # Current directory
    ]

    for path in possible_paths:
        if os.path.exists(path):
            print(f"Loading data from {path}")
            try:
                df = pd.read_csv(path)

                # Explicitly parse the month column
                df['month'] = pd.to_datetime(df['month'], format='%Y-%m')

                print(f"Successfully loaded {len(df)} records")
                return df
            except Exception as e:
                print(f"Error loading {path}: {str(e)}")
                continue
    raise FileNotFoundError("Could not find the HDB data CSV file in any of the expected locations.")

def display_basic_info(df):
    """Display comprehensive information about the dataset"""
    print("\n==== BASIC DATASET INFO ====")

    # Display the first few rows
    print("\nDataset Preview (First 5 Rows):")
    print(df.head())

    # Display basic statistics
    print("\nBasic Statistics:")
    print(df.describe(include='all').transpose())

    # Display column information with data types
    print("\nColumn Information:")
    info_table = PrettyTable()
    info_table.field_names = ["#", "Column Name", "Data Type", "Non-Null Count", "Unique Values"]
    info_table.align = "l"

    for i, col in enumerate(df.columns, 1):
        info_table.add_row([

```

```

        i,
        col,
        df[col].dtype,
        df[col].count(),
        df[col].nunique()
    ])

print(info_table)

# Check for missing values
print("\nMissing Values Summary:")
missing = df.isnull().sum()
if missing.sum() == 0:
    print("No missing values found!")
else:
    missing_table = PrettyTable()
    missing_table.field_names = ["Column", "Missing Values", "Percentage"]
    missing_table.align = "l"

    for col, count in missing.items():
        if count > 0:
            missing_table.add_row([
                col,
                count,
                f"({(count/len(df))*100:.2f}%)"
            ])

print(missing_table)

def detect_column_types(df):
    """Use dabl to detect column types with better error handling"""
    print("\n==== COLUMN TYPE DETECTION (dabl) ====")

    try:
        # Preprocess the DataFrame before type detection
        df_processed = df.copy()

        # Convert categorical columns explicitly
        categorical_cols = ['town', 'flat_type', 'block', 'street_name', 'storey_range', 'flat_model', 'remaining_lease']
        for col in categorical_cols:
            df_processed[col] = df_processed[col].astype('category')

        # Use dabl to detect column types
        types = dabl.detect_types(df_processed)
        print("\nAutomatically Detected Column Types:")
        print(types)

        # Add manual verification
        print("\nManual Verification Suggestions:")
        verification_table = PrettyTable()
        verification_table.field_names = ["Column", "Detected Type", "Suggested Verification"]
        verification_table.align = "l"

        for col in df.columns:
            try:
                detected = types.loc[types['column'] == col, 'type'].values[0]
            except IndexError:
                detected = "Unknown"

            suggestion = ""

            if detected == 'continuous':
                suggestion = "Check for outliers and distribution"
            elif detected == 'categorical':
                suggestion = f"Top values: {df[col].value_counts().head(3).to_dict()}"
            elif detected == 'low_card_int':
                suggestion = f"Consider if should be categorical. Unique values: {df[col].nunique()}"

            verification_table.add_row([col, detected, suggestion])

        print(verification_table)

    except Exception as e:
        print(f"\nError detecting column types: {str(e)}")
        print("Showing basic dtypes instead:")
        print(df.dtypes)

def quick_dabl_visualization(df):
    """Generate quick visualization using dabl library"""
    print("\n==== QUICK DABL VISUALIZATION ====")

    try:
        # Preprocess the DataFrame
        df_processed = df.copy()

        # Convert categorical columns explicitly
        categorical_cols = ['town', 'flat_type', 'storey_range', 'flat_model', 'remaining_lease']
        for col in categorical_cols:
            if col in df_processed.columns:
                df_processed[col] = df_processed[col].astype('category')

        # Ensure date column is handled
        if 'month' in df_processed.columns:
            df_processed['month'] = pd.to_datetime(df_processed['month'], format='%Y-%m')

        # Matplotlib parameters
        plt.rcParams.update({
            'font.size': 10, # Base font size
            'axes.labelsize': 10, # Axis label font size
            'axes.titlesize': 16, # Axis title font size
            'xtick.labelsize': 10, # X-axis tick label size
            'ytick.labelsize': 10, # Y-axis tick label size
            'figure.figsize': (20, 15) # Figure size
        })

        # Explicitly set color cycle to default color scheme
        mako_palette = sns.color_palette("mako", 10)
        plt.rcParams['axes.prop_cycle'] = plt.cycler(color=mako_palette)

        # Use dabl for quick visualization
        print("\nGenerating dabl plot...")

        # Suppress warnings during plot generation
        with warnings.catch_warnings():
            warnings.filterwarnings("ignore", category=UserWarning)
            warnings.filterwarnings("ignore", category=FutureWarning)

        # Use 'resale_price' as the target column
        dabl.plot(df_processed, target_col="resale_price")

        plt.suptitle('DABL Quick Visualization of HDB Resale Prices', fontsize=16)
        plt.tight_layout()
        plt.show()

    except Exception as e:
        print(f"\nError creating dabl visualization: {str(e)}")
        print("Traceback:", traceback.format_exc())
        print("\nPossible solutions:")
        print("- Ensure dabl is correctly installed")
        print("- Check data types of columns")
        print("- Verify 'resale_price' column exists")

    # Modify your main function to use this visualization
    def main():

```

```
"""Enhanced main function with dabl visualization"""
try:
    # Load the data
    df = load_data()

    # Display basic information
    display_basic_info(df)

    # Detect column types using dabl
    detect_column_types(df)

    # Generate quick dabl visualization
    quick_dabl_visualization(df)

    print("\nData preview completed successfully!")

except Exception as e:
    print(f"\nERROR: Failed to complete data preview - {str(e)}")

if __name__ == "__main__":
    main()
```

Loading data from /Users/yvonnelip/Desktop/CapstoneProject/singapore_hdb_resale_flat_prices.csv
Successfully loaded 203412 records

== BASIC DATASET INFO ==

Dataset Preview (First 5 Rows):
`_id month town flat_type block street_name storey_range \n
0 1 2017-01-01 ANG MO KIO 2 ROOM 406 ANG MO KIO AVE 10 10 TO 12
1 2 2017-01-01 ANG MO KIO 3 ROOM 108 ANG MO KIO AVE 4 01 TO 03
2 3 2017-01-01 ANG MO KIO 3 ROOM 602 ANG MO KIO AVE 5 01 TO 03
3 4 2017-01-01 ANG MO KIO 3 ROOM 465 ANG MO KIO AVE 10 04 TO 06
4 5 2017-01-01 ANG MO KIO 3 ROOM 601 ANG MO KIO AVE 5 01 TO 03`

`floor_area_sqm flat_model lease_commence_date remaining_lease \
0 44.0 Improved 1979 61 years 04 months
1 67.0 New Generation 1978 60 years 07 months
2 67.0 New Generation 1980 62 years 05 months
3 68.0 New Generation 1980 62 years 01 month
4 67.0 New Generation 1980 62 years 05 months`

resale_price

`0 2320000.0
1 2500000.0
2 2620000.0
3 2650000.0
4 265000.0`

Basic Statistics:

	count	unique	top	freq	\
_id	203412.0	NaN	NaN	NaN	
month	203412	NaN	NaN	NaN	
town	203412	26	SENGKANG	16753	
flat_type	203412	7	4 ROOM	86002	
block	203412	2735		2	615
street_name	203412	573	YISHUN RING RD	2921	
storey_range	203412	17	04 TO 06	46753	
floor_area_sqm	203412.0	NaN	NaN	NaN	
flat_model	203412	21	Model A	71845	
lease_commence_date	203412.0	NaN	NaN	NaN	
remaining_lease	203412	685	94 years 10 months	1827	
resale_price	203412.0	NaN	NaN	NaN	

	mean	min	\
_id	101706.5	1.0	
month	2021-04-29 08:14:16.728216576	2017-01-01 00:00:00	
town	NaN	NaN	
flat_type	NaN	NaN	
block	NaN	NaN	
street_name	NaN	NaN	
storey_range	NaN	NaN	
floor_area_sqm	96.925036	31.0	
flat_model	NaN	NaN	
lease_commence_date	1996.254292	1966.0	
remaining_lease	NaN	NaN	
resale_price	512368.066868	140000.0	

	25%	50%	\
_id	50853.75	101706.5	
month	2019-06-01 00:00:00	2021-07-01 00:00:00	
town	NaN	NaN	
flat_type	NaN	NaN	
block	NaN	NaN	
street_name	NaN	NaN	
storey_range	NaN	NaN	
floor_area_sqm	82.0	93.0	
flat_model	NaN	NaN	
lease_commence_date	1985.0	1996.0	
remaining_lease	NaN	NaN	
resale_price	380000.0	480000.0	

	75%	max	std	\
_id	152559.25	203412.0	58720.130816	
month	2023-05-01 00:00:00	2025-04-01 00:00:00	NaN	
town	NaN	NaN	NaN	
flat_type	NaN	NaN	NaN	
block	NaN	NaN	NaN	
street_name	NaN	NaN	NaN	
storey_range	NaN	NaN	NaN	
floor_area_sqm	112.0	366.7	24.028587	
flat_model	NaN	NaN	NaN	
lease_commence_date	2011.0	2021.0	14.192716	
remaining_lease	NaN	NaN	NaN	
resale_price	615000.0	1600000.0	179937.261025	

Column Information:

#	Column Name	Data Type	Non-Null Count	Unique Values
1	_id	int64	203412	203412
2	month	datetime64[ns]	203412	100
3	town	object	203412	26
4	flat_type	object	203412	7
5	block	object	203412	2735
6	street_name	object	203412	573
7	storey_range	object	203412	17
8	floor_area_sqm	float64	203412	179
9	flat_model	object	203412	21
10	lease_commence_date	int64	203412	56
11	remaining_lease	object	203412	685
12	resale_price	float64	203412	4267

Missing Values Summary:
No missing values found!

== COLUMN TYPE DETECTION (dabl) ==

Automatically Detected Column Types:

	continuous	dirty_float	low_card_int	ordinal	\
_id	False	False	False	False	
month	False	False	False	False	
town	False	False	False	False	
flat_type	False	False	False	False	
block	False	False	False	False	
street_name	False	False	False	False	
storey_range	False	False	False	False	
floor_area_sqm	True	False	False	False	
flat_model	False	False	False	False	
lease_commence_date	False	False	True	True	
remaining_lease	False	False	False	False	
resale_price	True	False	False	False	

	low_card_int_categorical	categorical	date	\
_id	False	False	False	
month	False	False	True	
town	False	True	False	
flat_type	False	True	False	
block	False	True	False	
street_name	False	True	False	
storey_range	False	True	False	
floor_area_sqm	False	False	False	
flat_model	False	True	False	
lease_commence_date	False	False	False	
remaining_lease	False	True	False	

```

resale_price           False   False  False
free_string      False
useless          False
_id              False  True
month             False  False
town              False  False
flat_type         False  False
block              False  False
street_name       False  False
storey_range      False  False
floor_area_sqm    False  False
flat_model        False  False
lease_commence_date False  False
remaining_lease    False  False
resale_price      False  False

```

Manual Verification Suggestions:

Error detecting column types: 'column'

Showing basic dtypes instead:

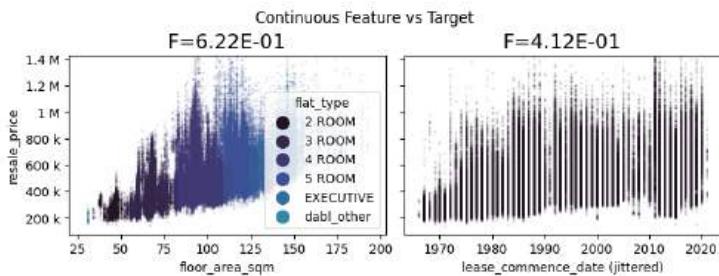
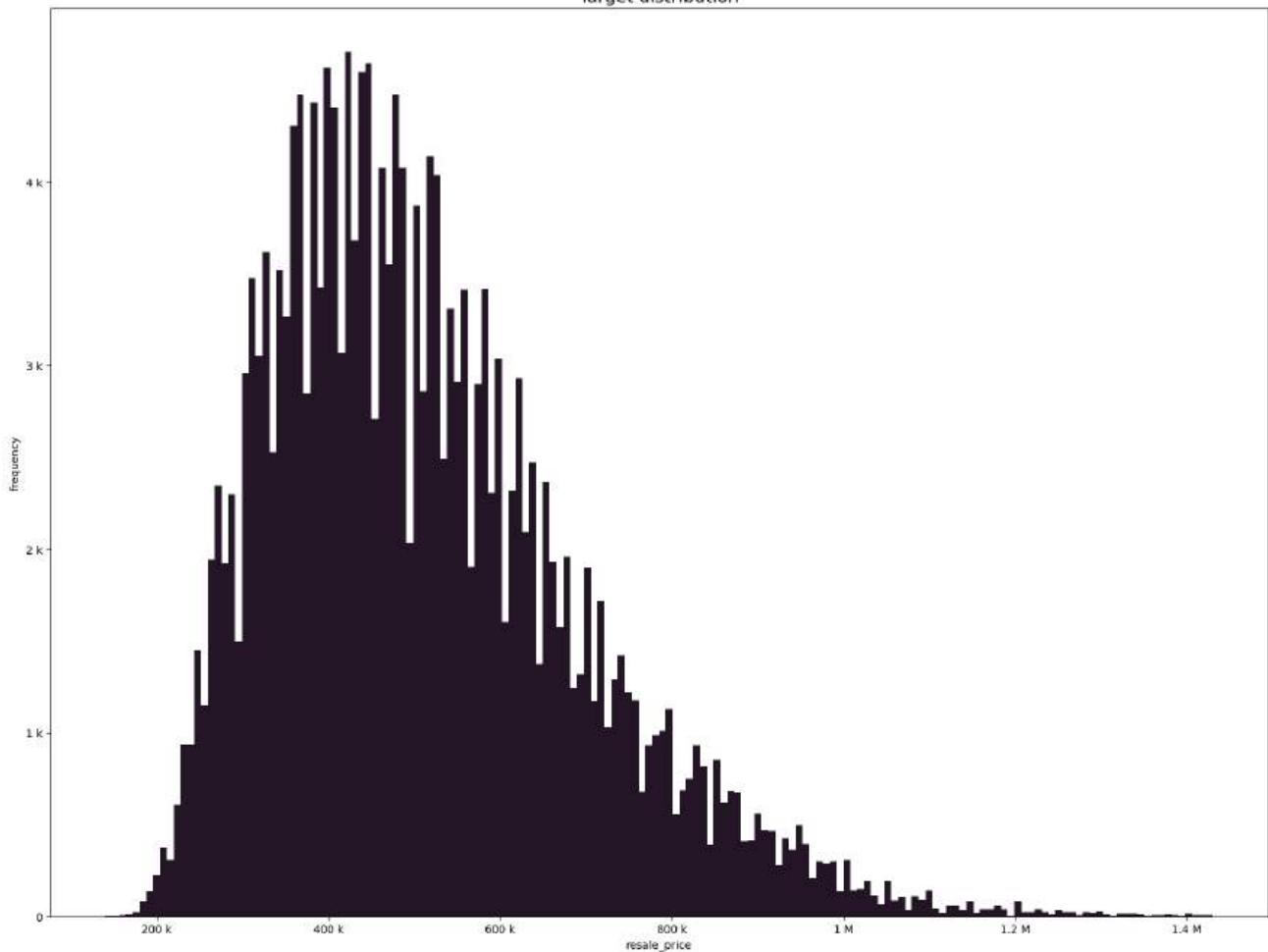
		int64
_id	datatype	datetime64[ns]
month	datatype	object
town	datatype	object
flat_type	datatype	object
block	datatype	object
street_name	datatype	object
storey_range	datatype	object
floor_area_sqm	datatype	float64
flat_model	datatype	object
lease_commence_date	datatype	int64
remaining_lease	datatype	object
resale_price	datatype	float64
dtype	datatype	object

==== QUICK DABL VISUALIZATION ====

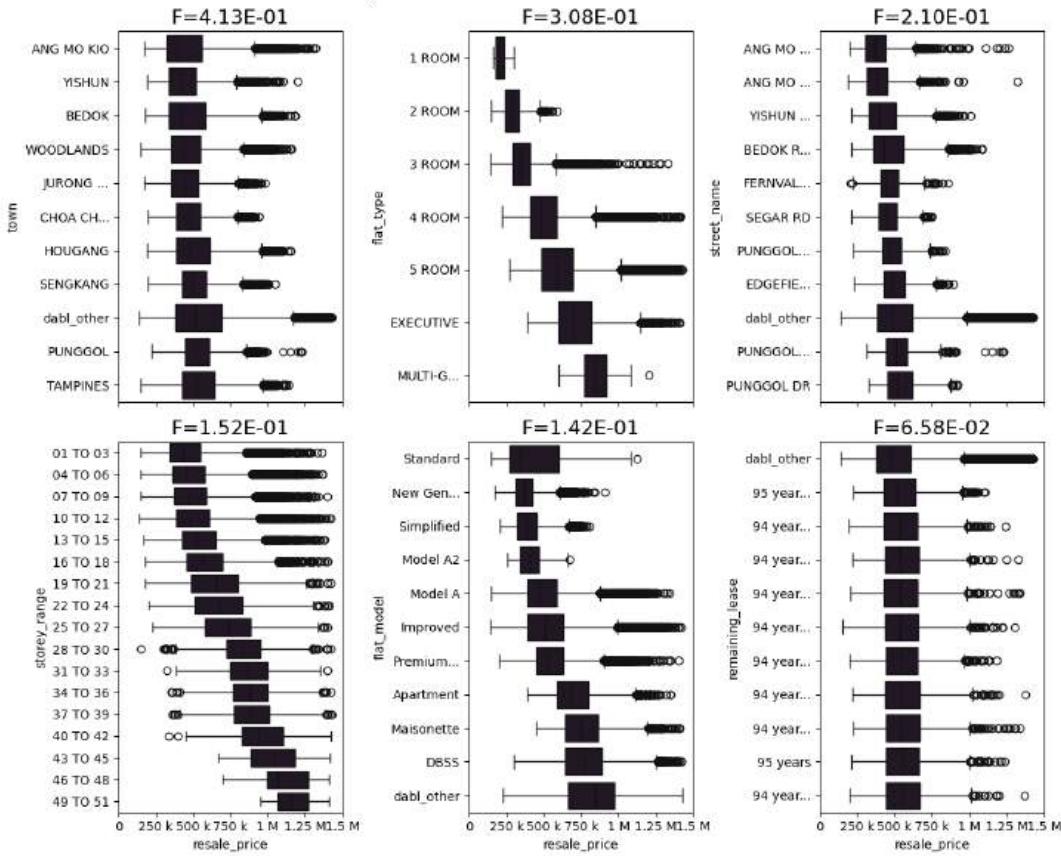
Generating dabl plot..

Target looks like regression

Target distribution



DABL Quick Visualization of HDB Resale Prices



Data preview completed successfully!

```
In [5]: # 5.2. Data Understanding
# Missing Data Analysis - using Missingno library to provide simple and intuitive visualizations to explore missing data patterns
# Run the library cell to import the necessary libraries
```

```
# Suppress FutureWarnings
warnings.filterwarnings("ignore", category=FutureWarning)

def load_data():
    """Load the HDB data from CSV"""
    # Define the path to the CSV file
    default_path = os.path.join(os.path.expanduser("~/"), "Desktop", "CapstoneProject", "singapore_hdb_resale_flat_prices.csv")
    fallback_path = os.path.join(os.path.expanduser("~/"), "Downloads", "singapore_hdb_resale_flat_prices.csv")

    # Try to load from the default path, fall back to Downloads if necessary
    if os.path.exists(default_path):
        file_path = default_path
    elif os.path.exists(fallback_path):
        file_path = fallback_path
    else:
        raise FileNotFoundError("Could not find the HDB data CSV file. Please run the data ingestion script first.")

    print(f"Loading data from {file_path}")
    return pd.read_csv(file_path)

# Load the dataset
df = load_data()

# Define 'month' to be converted to datetime format
df['month'] = pd.to_datetime(df['month'], errors='coerce')

# Check for NaT values after conversion
print("Check for NaT values after conversion:")
print(df['month'].isna().sum())

# Visualize missing values for a sample of the data
print("\nVisualizing missing values for sampled data:")
plt.figure(figsize=(12, 8))
msno.matrix(df.sample(250, random_state=1))
plt.title('Missing Value Matrix (250 Sample Records)')
plt.show()

# Create a null pattern for demonstration (optional)
print("\nCreating a demonstration null pattern for time series data:")
null_pattern = (np.random.random(1000).reshape((50, 20)) > 0.5).astype(bool)
null_pattern = pd.DataFrame(null_pattern).replace({False: None})
null_pattern.index = pd.period_range('1/1/2011', periods=50, freq='M')

# Visualize the null pattern
print("\nVisualizing null pattern in a time series:")
plt.figure(figsize=(12, 8))
msno.matrix(null_pattern)
plt.title('Null Pattern in Time Series Data (Demonstration)')
plt.show()

# Visualize missing values with a bar plot
print("\nVisualizing missing data with a bar plot:")
plt.figure(figsize=(12, 8))
msno.bar(df.sample(1000, random_state=1))
plt.title('Bar Plot of Missing Values (1000 Sample Records)')
plt.show()

# Visualize missing values with a heatmap
print("\nChecking for correlations in missing values:")
if df.isnull().values.any():
    plt.figure(figsize=(12, 8))
    msno.heatmap(df)
    plt.title('Heatmap of Missing Values Correlation')
    plt.show()
else:
    print("No missing values in the DataFrame, heatmap will not be displayed.")

# Add dendrogram visualization (additional from original code)
print("\nVisualizing hierarchical relationships in missing data:")
```

```

if df.isnull().values.any():
    plt.figure(figsize=(12, 8))
    msno.dendrogram(df.sample(1000, random_state=1))
    plt.title('Dendrogram of Missing Values (1000 Sample Records)')
    plt.show()
else:
    print("No missing values in the DataFrame, dendrogram will not be displayed.")

# Print summary of missing values by column
print("\nSummary of missing values by column:")
missing_summary = pd.DataFrame({
    'Column': df.columns,
    'Missing Values': df.isnull().sum().values,
    'Percentage': (df.isnull().sum() / len(df)) * 100).values
})
missing_summary = missing_summary.sort_values('Missing Values', ascending=False)
print(missing_summary)

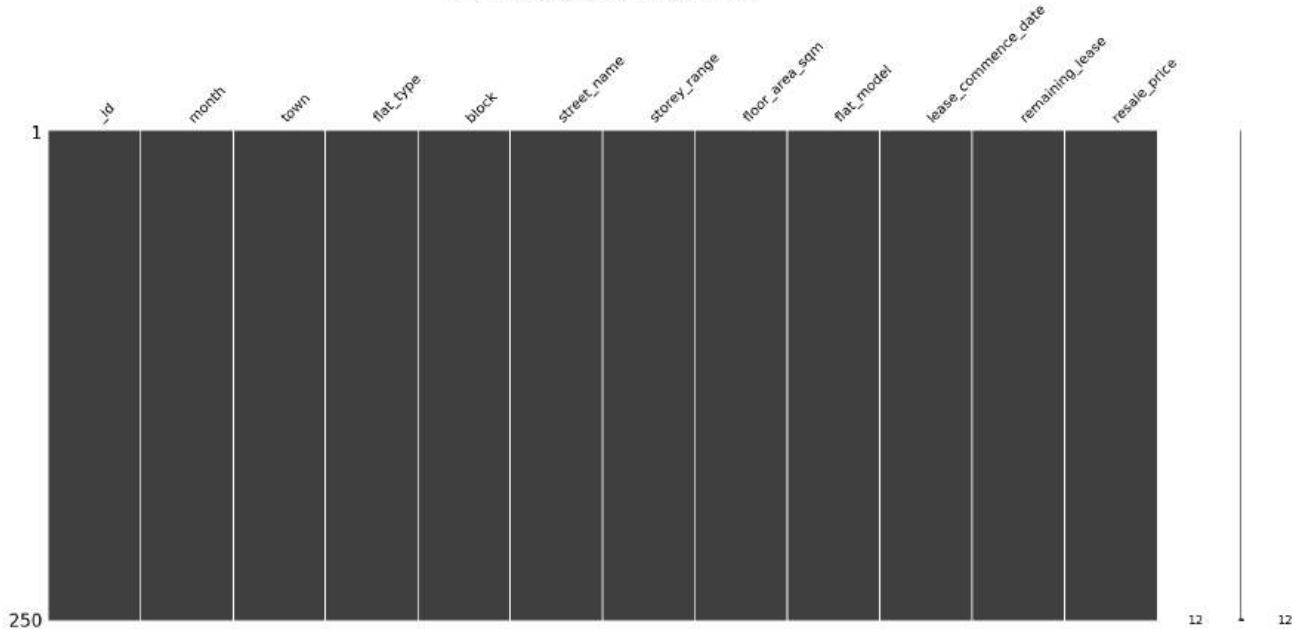
print("\nMissing data analysis completed!")

```

Loading data from /Users/yvonneip/Desktop/CapstoneProject/singapore_hdb_resale_flat_prices.csv
Check for NaT values after conversion:
0

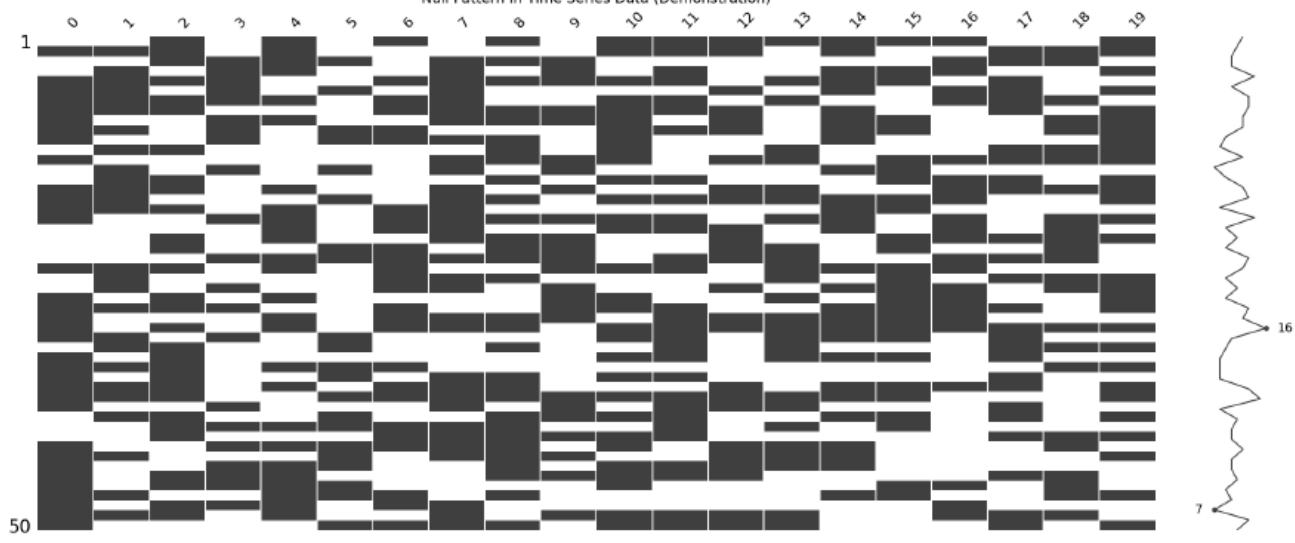
Visualizing missing values for sampled data:
<Figure size 1200x800 with 0 Axes>

Missing Value Matrix (250 Sample Records)

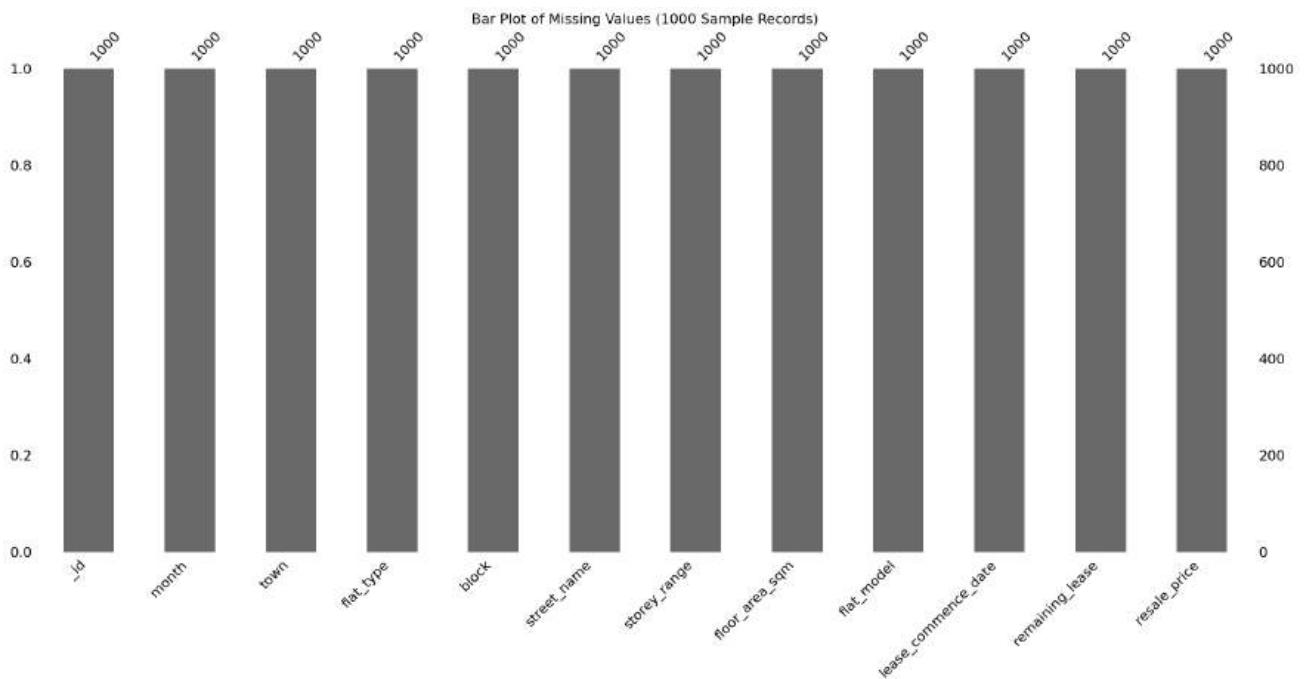


Creating a demonstration null pattern for time series data:
Visualizing null pattern in a time series:
<Figure size 1200x800 with 0 Axes>

Null Pattern in Time Series Data (Demonstration)



Visualizing missing data with a bar plot:



Checking for correlations in missing values:
No missing values in the Dataframe, heatmap will not be displayed.

Visualizing hierarchical relationships in missing data:
No missing values in the DataFrame, dendrogram will not be displayed.

Summary of missing values by column:

	Column	Missing Values	Percentage
0	'_id'	0	0.0
1	'month'	0	0.0
2	'town'	0	0.0
3	'flat_type'	0	0.0
4	'block'	0	0.0
5	'street_name'	0	0.0
6	'storey_range'	0	0.0
7	'floor_area_sqm'	0	0.0
8	'flat_model'	0	0.0
9	'lease_commence_date'	0	0.0
10	'remaining_lease'	0	0.0
11	'resale_price'	0	0.0

Missing data analysis completed!

```
In [6]: # 5.3. Data Cleaning and Transformation
# Using klib library to process HDB resale flat data
# Run the library cell to import the necessary libraries

# -----
# 1. Setup and Configuration
# -----

# Suppress warnings for cleaner output
warnings.filterwarnings('ignore')

# Set consistent color scheme using "mako" palette
mako_palette = sns.color_palette("mako")
mako_cmap = sns.color_palette("mako", as_cmap=True)
sns.set_palette(mako_palette)
plt.rcParams['axes.prop_cycle'] = plt.cycler(color=mako_palette)

# Set global plot parameters
plt.rcParams.update({
    'figure.facecolor': 'white',
    'axes.facecolor': 'white',
    'font.family': 'Arial',
    'lines.linewidth': 2.0,
    'patch.linewidth': 1.5
})

# Convert palette to hex values for Plotly
# Fixed to use matplotlib directly since we imported it
plotly_colors = [matplotlib.colors.rgb2hex(c) for c in mako_palette]

# -----
# 2. Data Visualization with klib
# -----


# Missing Values Visualization
print("Generating missing values visualization...")
plt.figure(figsize=(14, 7))
klib.missingval_plot(df)
plt.title('Missing Values Analysis', fontsize=14, fontweight='bold', color='black')
plt.xticks(fontsize=12, color='black')
plt.yticks(fontsize=12, color='black')
plt.tight_layout()
plt.show()

# Clean the dataset to define df_cleaned
print("Cleaning the dataset...")
df_cleaned = df.drop(columns=['_id']).drop_duplicates().dropna()

# Create derived columns needed for analysis
print("Creating derived columns...")
# Calculate current year for age calculations
current_year = 2025 # Current year

# Calculate remaining lease years (if lease_commence_date exists)
if 'lease_commence_date' in df_cleaned.columns:
    df_cleaned['remaining_lease_years'] = 99 - (current_year - df_cleaned['lease_commence_date'])

# Calculate flat age (if lease_commence_date exists)
if 'lease_commence_date' in df_cleaned.columns:
    df_cleaned['flat_age'] = current_year - df_cleaned['lease_commence_date']

# Calculate floor average (if storey_range exists)
if 'storey_range' in df_cleaned.columns:
    # Extract floor numbers from ranges like '01 TO 03'
    df_cleaned[['floor_min', 'floor_max']] = df_cleaned['storey_range'].str.extract(r'(\d+).*(\d+)')
```

```

df_cleaned['floor_min'] = pd.to_numeric(df_cleaned['floor_min'])
df_cleaned['floor_max'] = pd.to_numeric(df_cleaned['floor_max'])
df_cleaned['floor_avg'] = (df_cleaned['floor_min'] + df_cleaned['floor_max']) / 2

# Calculate price per square meter
if 'resale_price' in df_cleaned.columns and 'floor_area_sqm' in df_cleaned.columns:
    df_cleaned['price_per_sqm'] = df_cleaned['resale_price'] / df_cleaned['floor_area_sqm']

# Print available columns to verify
print("Available columns after creating derived features:")
print(df_cleaned.columns.tolist())

# Correlation Plot
print("Generating correlation plot...")
plt.figure(figsize=(12, 8))
klib.corr_plot(df_cleaned, cmap=mako_cmap)
plt.title('Correlation Analysis of Variables', fontsize=18, fontweight='bold', color='black')
# Customize axis labels
for ax in plt.gcf().get_axes():
    for item in ([ax.title, ax.xaxis.label, ax.yaxis.label] +
                [ax.get_xticklabels() + ax.get_yticklabels()]):
        item.set_fontsize(12)
        item.set_color('black')
plt.tight_layout()
plt.show()

# Target-based Correlation Plot (with reversed color order)
print("Generating target correlation plot...")
plt.figure(figsize=(12, 6))

# Use the reversed "mako" colormap
reversed_mako_cmap = sns.color_palette("mako_r", as_cmap=True)

klib.corr_plot(df_cleaned, target='resale_price', cmap=reversed_mako_cmap)
plt.title('Correlation with Resale Price (Reversed Colors)', fontsize=18, fontweight='bold', color='black')

# Customize text elements
for ax in plt.gcf().get_axes():
    for text in ax.texts:
        text.set_fontsize(12)
        text.set_color('white')
    for item in ([ax.title, ax.xaxis.label, ax.yaxis.label] +
                [ax.get_xticklabels() + ax.get_yticklabels()]):
        item.set_fontsize(12)
        item.set_color('black')
plt.tight_layout()
plt.show()

# Correlation Matrix Heatmap
print("Generating correlation heatmap...")
numeric_cols = df_cleaned.select_dtypes(include=['float64', 'int64']).columns.tolist()
corr_matrix = df_cleaned[numeric_cols].corr()
plt.figure(figsize=(12, 8))
sns.heatmap(
    corr_matrix,
    annot=True,
    cmap=mako_cmap,
    fmt=".2f",
    linewidths=0.5,
    annot_kws={"size": 12, "weight": "bold"}
)
plt.title('Correlation Matrix of Variables', fontsize=14, fontweight='bold', color='black')
plt.xticks(fontsize=12, color='black', rotation=45, ha='right')
plt.yticks(fontsize=12, color='black')
plt.tight_layout()
plt.show()

# Pairplot of Numerical Variables
print("Generating pairplot for key numerical variables...")

# Define desired columns for the pairplot
desired_columns = ['resale_price', 'floor_area_sqm', 'remaining_lease_years', 'flat_age', 'floor_avg', 'price_per_sqm']

# Filter to only include columns that exist in the dataframe
key_numeric_vars = [col for col in desired_columns if col in df_cleaned.columns]
print(f"Using these columns for the pairplot: {key_numeric_vars}")

# Create a smaller sample if the dataset is very large
if len(df_cleaned) > 5000:
    print("Large dataset detected, using a sample of 5000 records for pairplot...")
    plot_sample = df_cleaned[key_numeric_vars].sample(n=5000, random_state=42)
else:
    plot_sample = df_cleaned[key_numeric_vars]

# Create the pairplot with custom color scheme
plt.figure(figsize=(16, 14))

# Use the "mako" color palette
custom_palette = sns.color_palette("mako_r")

if 'flat_type' in df_cleaned.columns:
    # Select key variables including flat_type
    plot_vars = key_numeric_vars.copy()
    plot_sample = df_cleaned[plot_vars + ['flat_type']]
    if len(df_cleaned) > 5000:
        plot_sample = plot_sample.sample(n=5000, random_state=42)

    # Create the pairplot with flat_type as hue
    pair_plot = sns.pairplot(
        plot_sample,
        vars=plot_vars,
        hue='flat_type',
        diag_kind='kde',
        plot_kws=dict(alpha=0.1, edgecolor=None, s=25),
        diag_kws=dict(linewidth=2),
        corner=True,
        palette=custom_palette
    )

    # Remove the default legend
    pair_plot._legend.remove()

else:
    # Fallback if flat_type is not available
    plot_sample['price_quartile'] = pd.qcut(plot_sample['resale_price'], 4, labels=False)

    pair_plot = sns.pairplot(
        plot_sample,
        diag_kind='kde',
        plot_kws=dict(alpha=0.1, edgecolor=None, s=25),
        diag_kws=dict(linewidth=2),
        corner=True,
        hue='price_quartile',
        palette=custom_palette[:4] # Use first 4 colors for quartiles
    )

# Customize the pairplot appearance
pair_plot.fig.suptitle('Pairwise Relationships Between Key Numerical Variables',
                      fontsize=16, fontweight='bold', y=1.02)

# Add a custom legend if using flat_type
if 'flat_type' in df_cleaned.columns:
    handles = pair_plot._legend_data.values()
    labels = pair_plot._legend_data.keys()

```

```

pair_plot.fig.legend(handles=handles, labels=labels, title='Flat Type',
                     loc='upper right', bbox_to_anchor=(1, 1))

# Improve text readability
for ax in pair_plot.axes.flat:
    # Skip empty axes from corner=True
    if ax is not None:
        # Set axis labels
        if ax.get_xlabel():
            ax.set_xlabel(ax.get_xlabel(), fontsize=12, color='black')
        if ax.get_ylabel():
            ax.set_ylabel(ax.get_ylabel(), fontsize=12, color='black')

    # Set tick parameters
    ax.tick_params(labelsize=10, colors='black')

# Adjust layout
plt.tight_layout()
plt.show()

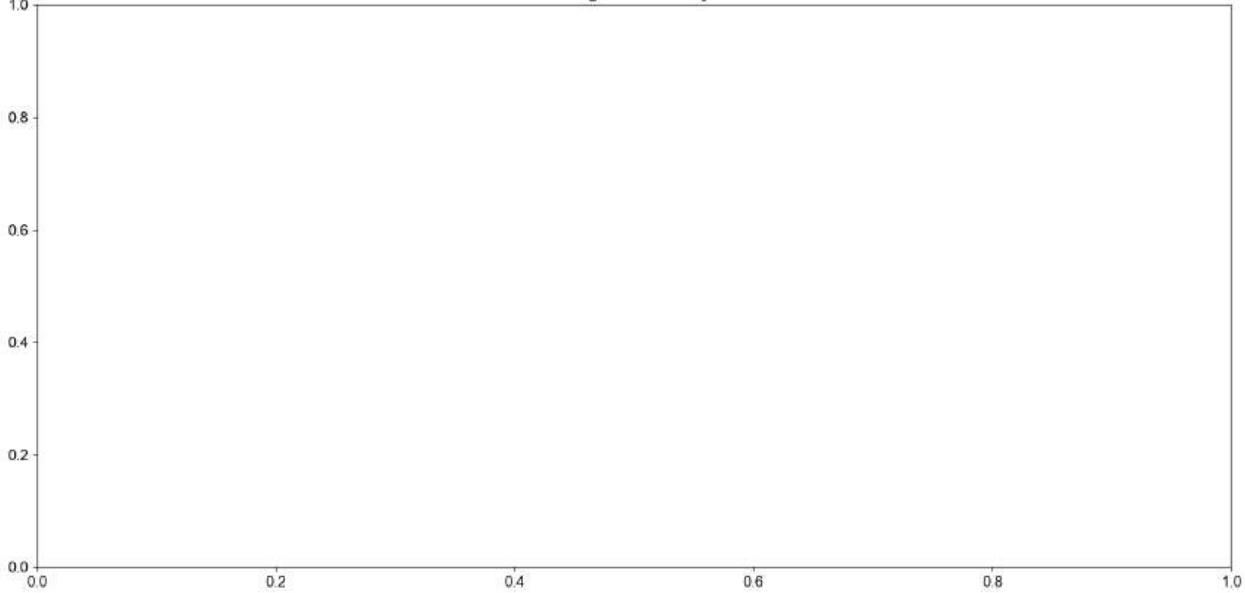
# Add information about the pairplot
print("\nPairplot analysis:")
print("- Diagonal: Shows distribution of each variable")
print("- Off-diagonal: Shows relationship between pairs of variables")
print("- Strong correlations appear as clear patterns in the scatter plots")

print("\nData cleaning and visualization completed!")

```

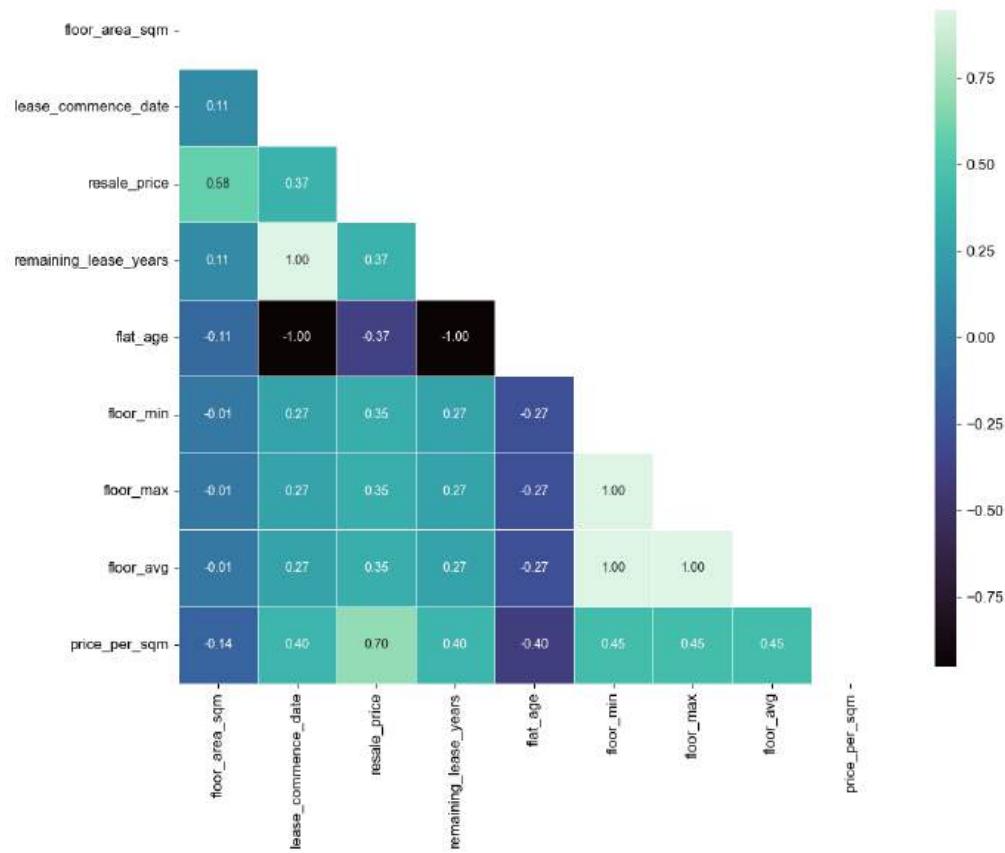
Generating missing values visualization...
No missing values found in the dataset.

Missing Values Analysis



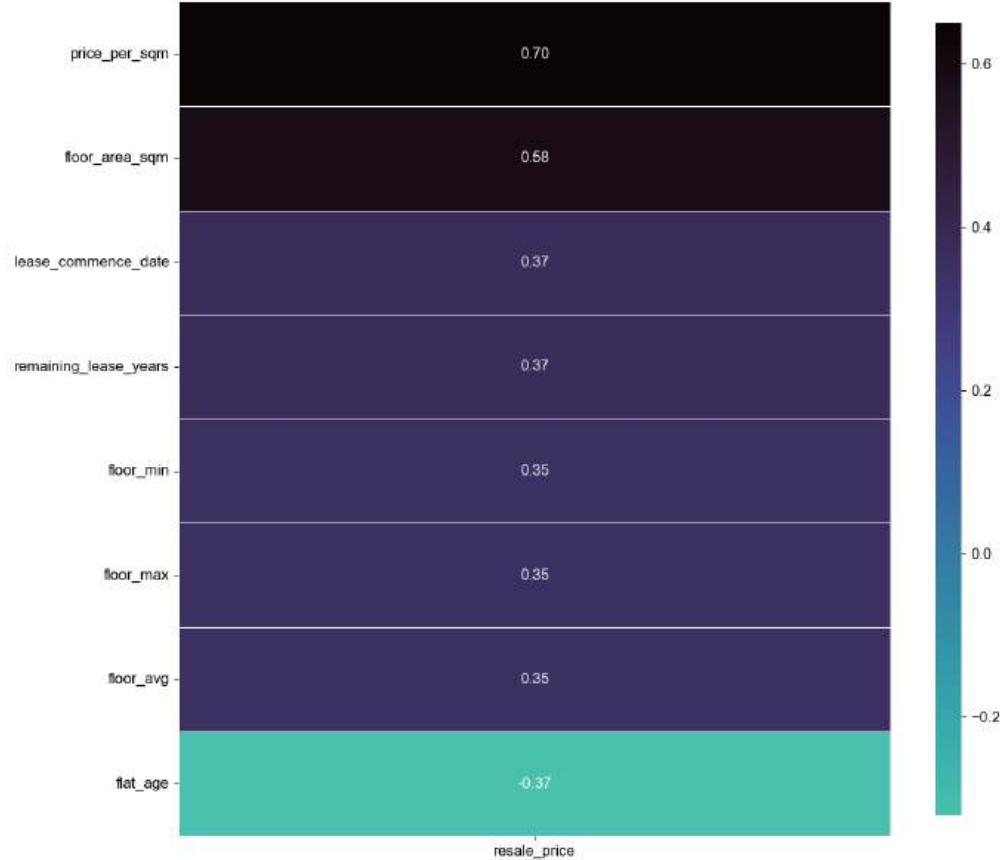
Cleaning the dataset...
Creating derived columns...
Available columns after creating derived features:
['month', 'town', 'flat_type', 'block', 'street_name', 'storey_range', 'floor_area_sqm', 'flat_model', 'lease_commence_date', 'remaining_lease', 'resale_price', 'remaining_lease_years', 'flat_age', 'floor_min', 'floor_max', 'floor_avg', 'price_per_sqm']
Generating correlation plot...
<Figure size 1200x800 with 0 Axes>

Correlation Analysis of Variables

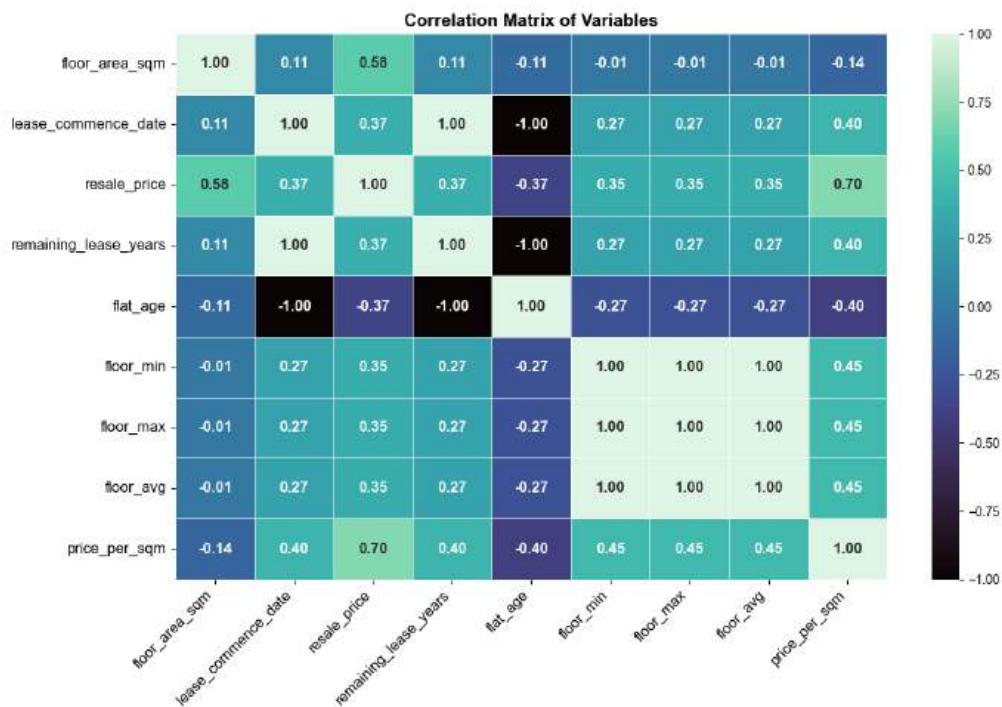


Generating target correlation plot...
<Figure size 1200x600 with 0 Axes>

Correlation with Resale Price (Reversed Colors)

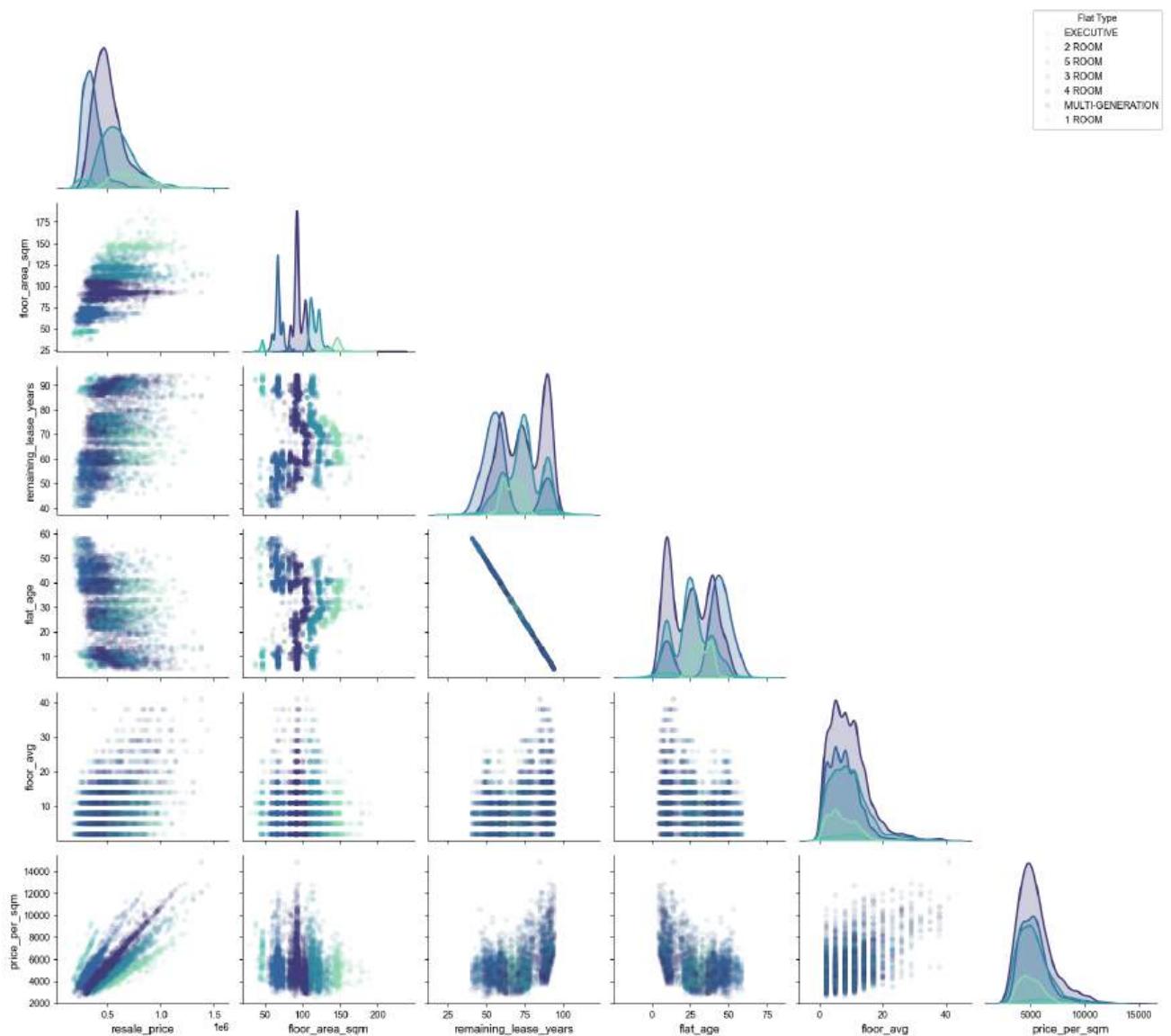


Generating correlation heatmap...



Generating pairplot for key numerical variables...
 Using these columns for the pairplot: ['resale_price', 'floor_area_sqm', 'remaining_lease_years', 'flat_age', 'floor_avg', 'price_per_sqm']
 Large dataset detected, using a sample of 5000 records for pairplot...
 <Figure size 1600x1400 with 0 Axes>

Pairwise Relationships Between Key Numerical Variables



Pairplot analysis:

- Diagonal: Shows distribution of each variable
- Off-diagonal: Shows relationship between pairs of variables
- Strong correlations appear as clear patterns in the scatter plots

Data cleaning and visualization completed!

6. Exploratory Data Analysis (EDA) - Notebook Instructions

This set of code cells provides a comprehensive EDA of Singapore HDB resale data. Each cell focuses on a specific aspect of the analysis and can be run independently, but they must be run in sequence (Cell 1 must be run first).

How to Use These Cells

1. First, update the `input_path` variable in Cell 1 to point to your CSV file location.
2. Run Cell 1 (Setup and Data Loading) first. This loads all necessary libraries and the data.
3. Run any of the other cells based on the analysis you need:
 - Cell 2: Advanced Statistical Analysis
 - Cell 3: Price Distribution by Segments
 - Cell 4: Feature Interaction Analysis
 - Cell 5: Price Per Square Meter Analysis
 - Cell 6: Advanced Time Series Analysis
 - Cell 7: Advanced Visualizations with Plotly
 - Cell 8: Advanced Regression Analysis
 - Cell 9: Market Segmentation Analysis
 - Cell 10: Comparative Analysis by Year

Troubleshooting

If you encounter any issues:

- Make sure your CSV file is in the correct location and format
- Check that all required libraries are installed (pandas, numpy, matplotlib, seaborn, plotly, prettytable, etc.)
- For time series analysis, ensure the 'month' column can be converted to datetime format

Script Structure

The analysis is broken into logical sections to make it easier to manage:

```
Cell 1: Setup and Data Loading
└── Library imports
└── Color scheme setup
└── Data loading and basic preprocessing
└── Dataset overview table

Cell 2: Advanced Statistical Analysis
└── Detailed statistics with visualizations

Cell 3: Price Distribution by Segments
└── Price analysis by various categorical variables

Cell 4: Feature Interaction Analysis
└── Analysis of interactions between features

Cell 5: Price Per Square Meter Analysis
└── Analysis focused on price per square meter

Cell 6: Advanced Time Series Analysis
└── Price and volume trends over time
└── Seasonal decomposition
└── Autocorrelation analysis

Cell 7: Advanced Visualizations with Plotly
└── Price analysis by floor range and remaining lease
└── Town price analysis and flat type distribution

Cell 8: Advanced Regression Analysis
└── Training various regression models (Linear, Ridge, Lasso, Random Forest)
└── Feature importance analysis
└── Partial dependence plots for top features

Cell 9: Market Segmentation Analysis
└── Price segmentation (Budget, Mid-Range, Premium, Luxury)
└── Segment distribution by flat type
└── Segment distribution by town
└── Floor range segment distribution

Cell 10: Comparative Analysis by Year
└── Year-over-year price and volume trends
└── Price volatility and growth analysis
└── Flat type distribution evolution over years
└── Town-based price trends by year
```

Each cell is designed to be self-contained but relies on the data loaded in Cell 1.

```
In [10]: # IMPORTANT TO RUN THIS CELL BEFORE ANY OTHER CELLS
# 1: Setup and Data Loading
# Handles all imports and data loading for other cells to reference

import os
import json
import warnings
import requests
import sys
import time
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib as mpl
import matplotlib.colors
from matplotlib.colors import to_hex

import plotly.express as px
import plotly.graph_objects as go
from plotly.subplots import make_subplots

from sklearn.linear_model import LinearRegression, Ridge, Lasso
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_percentage_error
from scipy import stats

import statsmodels.api as sm
from statsmodels.tsa.seasonal import seasonal_decompose
from statsmodels.tsa.stattools import acf
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf

import folium
from folium.features import GeoJsonTooltip
from folium.plugins import MarkerCluster

from prettytable import PrettyTable
from typing import Dict, List, Tuple
from IPython.display import display, HTML

# Suppress warnings for cleaner output
```

```

warnings.filterwarnings('ignore')

# Set custom color palette
mako_palette = sns.color_palette("mako", n_colors=10)
mako_colors = [matplotlib.colors.to_hex(color) for color in mako_palette]
sns.set_palette(mako_palette)
plt.rcParams['axes.prop_cycle'] = plt.cycler(color=mako_palette)

# Set global plot parameters
plt.rcParams.update({
    'figure.facecolor': 'white',
    'axes.facecolor': 'white',
    'font.family': 'Arial',
    'lines.linewidth': 2.0,
    'patch.linewidth': 1.5
})

# Convert palette to hex values for Plotly
plotly_colors = [matplotlib.colors.rgb2hex(c) for c in mako_palette]

# Common Plotly title font settings
title_font = dict(family="Arial", size=18, color="black", weight="bold")
axis_font = dict(family="Arial", size=14, color="black")

# Define the path to the CSV file for data loading
input_path = "/Users/yvonnelip/Desktop/CapStoneProject/singapore_hdb_resale_clean_2017_2024.csv"
print(f"Loading data from: {input_path}")

try:
    # Try to read the CSV with different encoding options if necessary
    try:
        df = pd.read_csv(input_path)
    except UnicodeDecodeError:
        print("Trying with different encoding...")
        df = pd.read_csv(input_path, encoding='utf-8')
    except:
        print("Trying with different encoding...")
        df = pd.read_csv(input_path, encoding='ISO-8859-1')

    print(f"Successfully loaded dataset with {df.shape[0]} rows and {df.shape[1]} columns")

    # Ensure numeric columns are properly typed
    numeric_columns = ['floor_area_sqm', 'resale_price', 'remaining_lease_years', 'flat_age']
    for col in numeric_columns:
        if col in df.columns and not pd.api.types.is_numeric_dtype(df[col]):
            print(f"Converting {col} to numeric...")
            df[col] = pd.to_numeric(df[col], errors='coerce')

except FileNotFoundError:
    print("Error: File not found at {input_path}")
    print("Please make sure the file exists at the specified location.")
    exit(1)
except Exception as e:
    print(f"Error loading dataset: {e}")
    exit(1)

# Calculate price per square meter if not already present
if 'price_per_sqm' not in df.columns and 'resale_price' in df.columns and 'floor_area_sqm' in df.columns:
    df['price_per_sqm'] = df['resale_price'] / df['floor_area_sqm']

# Ensure month is in datetime format if it exists
if 'month' in df.columns and not pd.api.types.is_datetime64_any_dtype(df['month']):
    try:
        df['month'] = pd.to_datetime(df['month'])
        print("Converted 'month' column to datetime format")
    except:
        df['month'] = pd.to_datetime(df['month'], format='%Y-%m')
        print("Converted 'month' column to datetime format using '%Y-%m' format")
    except:
        print("Warning: Could not convert 'month' column to datetime format")

# Basic Data Overview
print("\n===== DATASET OVERVIEW =====")
basic_info = PrettyTable()
basic_info.field_names = ["Attribute", "Value"]
basic_info.align = "l"
basic_info.add_row(["Total Records", f"{df.shape[0]}"])
basic_info.add_row(["Total Features", f"{df.shape[1]}"])
basic_info.add_row(["Numeric Features", f"{len(df.select_dtypes(include=['number']).columns)}"])
basic_info.add_row(["Categorical Features", f"{len(df.select_dtypes(include=['object']).columns)}"])
basic_info.add_row(["Missing Values", f"{df.isnull().sum().sum()}"])
basic_info.add_row(["Duplicates", f"{df.duplicated().sum()}"])

if 'month' in df.columns and pd.api.types.is_datetime64_any_dtype(df['month']):
    basic_info.add_row(["Date Range", f"{df['month'].min().strftime('%B %Y')} to {df['month'].max().strftime('%B %Y')}"])
elif 'year' in df.columns:
    basic_info.add_row(["Year Range", f"{df['year'].min()} to {df['year'].max()}"])

print(basic_info)

Loading data from: /Users/yvonnelip/Desktop/CapStoneProject/singapore_hdb_resale_clean_2017_2024.csv
Successfully loaded dataset with 196985 rows and 18 columns
Converted 'month' column to datetime format

===== DATASET OVERVIEW =====
+-----+-----+
| Attribute | Value |
+-----+-----+
| Total Records | 196,985 |
| Total Features | 18 |
| Numeric Features | 10 |
| Categorical Features | 7 |
| Missing Values | 0 |
| Duplicates | 0 |
| Date Range | January 2017 to December 2024 |
+-----+-----+

In [7]: # 2: Advanced Statistical Analysis
# Run 1 first to load data and setup

print("\n===== ADVANCED STATISTICAL ANALYSIS =====")

# -----
# 1. Calculate key statistics for resale_price
# -----
if 'resale_price' in df.columns:
    # Use the "mako" colormap
    mako_cmap = sns.color_palette("mako_r", as_cmap=True)

    # Sample the colormap at evenly spaced intervals to get discrete colors
    hex_codes = []
    for i in np.linspace(0, 1, 5):
        rgba = mako_cmap(i)
        hex_codes.append(mpl.colors.rgb2hex(rgba))

    # Create a prettier tabular format
    price_stats = PrettyTable()
    price_stats.field_names = ["Statistic", "Value"]
    price_stats.align = "l" # Left-align all columns

    # Group statistics by category with section headers
    price_stats.add_row(["\n CENTRAL TENDENCY", ""])
    price_stats.add_row(["Mean", f"${df['resale_price'].mean():,.2f}"])

```

```

price_stats.add_row(["Median", f"${df['resale_price'].median():,.2f}"])

price_stats.add_row(["\u2227 DISPERSION", ""])
price_stats.add_row(["Standard Deviation", f"${df['resale_price'].std():,.2f}"])
price_stats.add_row(["Coefficient of Variation", f"{{(df['resale_price'].std() / df['resale_price'].mean()) * 100}}:.2f%")]
price_stats.add_row(["Range", f"${df['resale_price'].max() - df['resale_price'].min():,.2f}"])

price_stats.add_row(["\u2227 DISTRIBUTION", ""])
price_stats.add_row(["25th Percentile", f"${df['resale_price'].quantile(0.25):,.2f}"])
price_stats.add_row(["75th Percentile", f"${df['resale_price'].quantile(0.75):,.2f}"])
price_stats.add_row(["Interquartile Range", f"${df['resale_price'].quantile(0.75) - df['resale_price'].quantile(0.25):,.2f}"])
price_stats.add_row(["Skewness", f'{df["resale_price"].skew():.4f}'])
price_stats.add_row(["Kurtosis", f'{df["resale_price"].kurtosis():.4f}'])

price_stats.add_row(["\u2227 EXTREMES", ""])
price_stats.add_row(["Minimum", f"${df['resale_price'].min():,.2f}"])
price_stats.add_row(["Maximum", f"${df['resale_price'].max():,.2f}"])

print("Resale Price Statistics:")
print(price_stats)

# Create a continuous colorscale
custom_colorscale = []
for i in np.linspace(0, 1, 10):
    rgba = mako_cmap(i)
    hex_color = mpl.colors.rgb2hex(rgba)
    custom_colorscale.append([i, hex_color])

# Create visualizations to complement the table
fig = make_subplots(
    rows=1, cols=3,
    subplot_titles=("Distribution", "Box Plot", "Q-Q Plot"),
    specs=[[{"type": "Histogram"}, {"type": "box"}, {"type": "scatter"}]],
    column_widths=[0.4, 0.2, 0.4],
    horizontal_spacing=0.12
)

# -----
# 2. Calculate KDE
# -----
# Histogram with KDE
kde_x = np.linspace(df['resale_price'].min(), df['resale_price'].max(), 1000)
kde = stats.gaussian_kde(df['resale_price'].dropna())
kde_y = kde(kde_x)

# Scale KDE to match histogram height
hist_values, hist_bins = np.histogram(df['resale_price'], bins=30)
scaling_factor = max(hist_values) / max(kde_y)
kde_y_scaled = kde_y * scaling_factor

# Add histogram with mako color
fig.add_trace(
    go.Histogram(
        x=df['resale_price'],
        name="Histogram",
        nbinsx=30,
        marker_color=hex_codes[2]
    ),
    row=1, col=1
)

# Add KDE as a line
fig.add_trace(
    go.Scatter(
        x=kde_x,
        y=kde_y_scaled,
        mode='lines',
        name='Density',
        line=dict(color=hex_codes[-2], width=2)
    ),
    row=1, col=1
)

# Create Box plot
fig.add_trace(
    go.Box(
        y=df['resale_price'],
        name="Price",
        marker_color=hex_codes[3],
        boxmean=True
    ),
    row=1, col=2
)

# -----
# 3. Calculate theoretical quantiles
# -----
# Q-Q plot for normality check
sorted_data = sorted(df['resale_price'].dropna())
n = len(sorted_data)
theoretical_quantiles = [stats.norm.ppf((i - 0.5) / n) for i in range(1, n + 1)]

# Add scatter plot
fig.add_trace(
    go.Scatter(
        x=theoretical_quantiles,
        y=sorted_data,
        mode='markers',
        name='Q-Q Plot',
        marker=dict(
            color=hex_codes[1],
            size=5
        )
    ),
    row=1, col=3
)

# Add reference line
# Estimate parameters for the reference line
slope = (np.percentile(sorted_data, 75) - np.percentile(sorted_data, 25)) / \
        (stats.norm.ppf(0.75) - stats.norm.ppf(0.25))
intercept = np.median(sorted_data) - slope * stats.norm.ppf(0.5)

x_line = np.array([min(theoretical_quantiles), max(theoretical_quantiles)])
y_line = slope * x_line + intercept

fig.add_trace(
    go.Scatter(
        x=x_line,
        y=y_line,
        mode='lines',
        name='Reference Line',
        line=dict(color=hex_codes[4], dash='dash')
    ),
    row=1, col=3
)

# Update layout
fig.update_layout(
    title={
        'text': "Statistical Analysis of Resale Prices",
        'x': 0.5,
        'xanchor': 'center',
        'yanchor': 'top'
    }
)

```

```

        'font': {
            'family': "Arial, sans-serif",
            'size': 24,
            'weight': 'bold'
        }
    },
    height=600,
    width=1600,
    showlegend=False,
    template="plotly_white",
    colorway=hex_codes
)

# Update axes
fig.update_xaxes(title_text="Price (SGD)", row=1, col=1)
fig.update_yaxes(title_text="Frequency", row=1, col=1)

fig.update_xaxes(title_text="", row=1, col=2)
fig.update_yaxes(title_text="Price (SGD)", row=1, col=2)

fig.update_xaxes(title_text="Theoretical Quantiles", row=1, col=3)
fig.update_yaxes(title_text="Sample Quantiles", row=1, col=3)

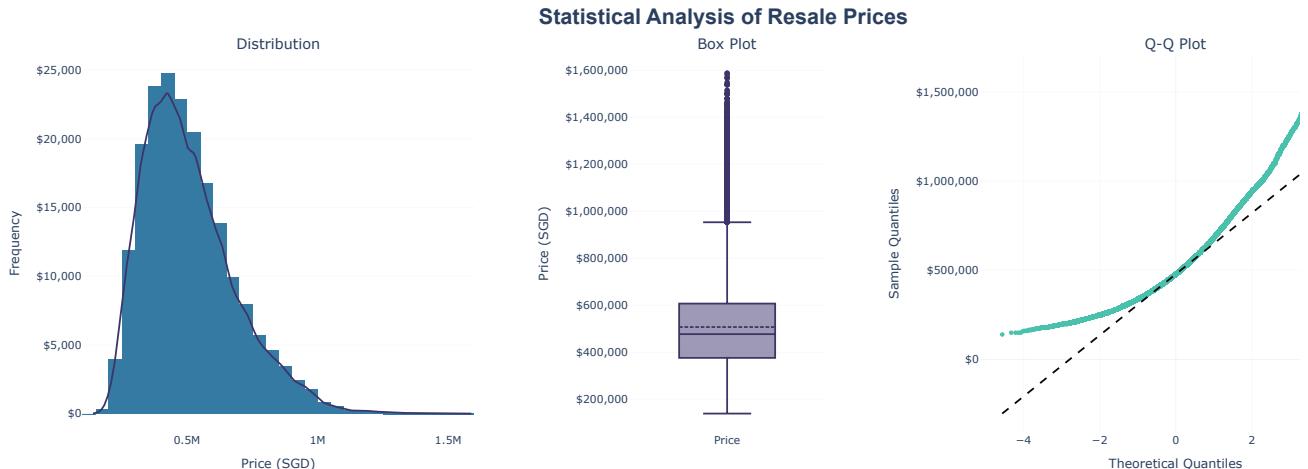
# Format y-axis with dollar signs
fig.update_yaxes(tickprefix="$", tickformat=",", row=1, col=1)
fig.update_yaxes(tickprefix="$", tickformat=",", row=1, col=2)
fig.update_yaxes(tickprefix="$", tickformat=",", row=1, col=3)

fig.show()

# Test for normality
stat, p = stats.normaltest(df['resale_price'].dropna())
print("\nNormality Test (D'Agostino's K^2):")
print("Statistic: {stat:.4f}, p-value: {p:.8f}")
if p < 0.05:
    print("Result: Data is not normally distributed (p < 0.05)")
else:
    print("Result: Data appears to be normally distributed (p >= 0.05)")

===== ADVANCED STATISTICAL ANALYSIS =====
Resale Price Statistics:
+-----+-----+
| Statistic | Value |
+-----+-----+
| CENTRAL TENDENCY |
| Mean | $508,157.17 |
| Median | $478,000.00 |
| DISPERSION |
| Standard Deviation | $177,739.74 |
| Coefficient of Variation | 34.98% |
| Range | $1,448,000.00 |
| DISTRIBUTION |
| 25th Percentile | $377,000.00 |
| 75th Percentile | $608,000.00 |
| Interquartile Range | $231,000.00 |
| Skewness | 0.9403 |
| Kurtosis | 1.0881 |
| EXTREMES |
| Minimum | $140,000.00 |
| Maximum | $1,588,000.00 |
+-----+-----+

```



```

Normality Test (D'Agostino's K^2):
Statistic: 25494.0283, p-value: 0.00000000
Result: Data is not normally distributed (p < 0.05)

```

```

In [3]: # 3: Price Distribution by Segments
# Run 1 first to load data and setup

print("\n===== PRICE DISTRIBUTION BY SEGMENTS =====")

# Set colormap
mako_cmap = sns.color_palette("mako_r", as_cmap=True)

hex_codes = []
for i in np.linspace(0, 1, 5):
    rgba = mako_cmap(i)
    hex_codes.append(mpl.colors.rgb2hex(rgba))

# Create distribution subplots using Plotly
fig = make_subplots(
    rows=2, cols=2,
    subplot_titles=[
        "Price Distribution (Overall)",
        "Price by Flat Type",
        "Price by Town (Top 10)",
        "Price by Floor Range"
    ],
    horizontal_spacing=0.12,
    specs=[{"type": "histogram"}, {"type": "bar"}, {"type": "bar"}, {"type": "bar"}]
)

# -----
# 1. Overall Price Distribution
# -----
fig.add_trace(

```

```

go.Histogram(
    x=df['resale_price'],
    nbins=50,
    marker_color=hex_codes[1],
    opacity=0.7,
    name="Overall"
),
row=1, col=1
)

# -----
# 2. Price by Flat Type
#
if 'flat_type' in df.columns:
    flat_type_prices = df.groupby('flat_type')['resale_price'].mean().sort_values().reset_index()

fig.add_trace(
    go.Bar(
        x=flat_type_prices['flat_type'],
        y=flat_type_prices['resale_price'],
        marker_color=hex_codes[2],
        name="Flat Type"
    ),
    row=1, col=2
)

# -----
# 3. Price by Town (Top 10)
#
if 'town' in df.columns:
    town_prices = df.groupby('town')['resale_price'].mean().sort_values(ascending=False).head(10).reset_index()

fig.add_trace(
    go.Bar(
        x=town_prices['town'],
        y=town_prices['resale_price'],
        marker_color=hex_codes[3],
        name="Town"
    ),
    row=2, col=1
)

# -----
# 4. Price by Floor Range
#
if 'storey_range' in df.columns:
    # Group by storey range and calculate average price
    floor_prices = df.groupby('storey_range')['resale_price'].mean().reset_index()
    # Sort by the first number in the range (assuming format like "01 TO 03")
    floor_prices['sort_key'] = floor_prices['storey_range'].str.extract(r'^(\d+)').astype(int)
    floor_prices = floor_prices.sort_values('sort_key')

fig.add_trace(
    go.Bar(
        x=floor_prices['storey_range'],
        y=floor_prices['resale_price'],
        marker_color=hex_codes[4],
        name="Floor Range"
    ),
    row=2, col=2
)

# Update layout and axes
fig.update_layout(
    title={
        'text': "Price Distribution by Various Segments",
        'x': 0.5,
        'xanchor': 'center',
        'y': 0.95,
        'yanchor': 'top',
        'font': {
            'family': "Arial, sans-serif",
            'size': 24,
            'weight': 'bold'
        }
    },
    height=800,
    template="plotly_white",
    showLegend=False,
    colorway=hex_codes
)

# Format axes
fig.update_xaxes(title_text="Price (SGD)", row=1, col=1)
fig.update_yaxes(title_text="Count", row=1, col=1)

fig.update_xaxes(title_text="Flat Type", row=1, col=2)
fig.update_yaxes(title_text="Avg Price (SGD)", row=1, col=2, tickprefix="$", tickformat=",")

fig.update_xaxes(title_text="Town", row=2, col=1, tickangle=-45)
fig.update_yaxes(title_text="Avg Price (SGD)", row=2, col=1, tickprefix="$", tickformat=",")

fig.update_xaxes(title_text="Floor Range", row=2, col=2, tickangle=-45)
fig.update_yaxes(title_text="Avg Price (SGD)", row=2, col=2, tickprefix="$", tickformat=",")

fig.update_layout(margin=dict(l=50, r=50, t=100, b=100))

fig.show()
===== PRICE DISTRIBUTION BY SEGMENTS =====

```

Price Distribution by Various Segments



```
In [6]: # 4: Feature Interaction Analysis
# Run 1 first to load data and setup

print("\n===== FEATURE INTERACTION ANALYSIS =====")

# -----
# 1. Combine Price Distribution by Flat Type and Town & Price by Floor Range
# -----
if 'flat_type' in df.columns and 'town' in df.columns and 'storey_range' in df.columns:
    # Set colormap
    mako_cmap = sns.color_palette("mako", as_cmap=True)

    # Sample the colormap at evenly spaced intervals to get discrete colors
    hex_codes = []
    for i in np.linspace(0, 1, 5):
        rgba = mako_cmap(i)
        hex_codes.append(mpl.colors.rgb2hex(rgba))

    # Create a temporary dataframe for floor analysis
    temp_df = df.copy()

    # Extract the lower floor number and convert to integer
    temp_df['floor_number'] = temp_df['storey_range'].str.extract(r'^(\d+)').iloc[:, 0].astype(int)

    # Create floor category based on the extracted floor number
    temp_df['floor_category'] = pd.cut(
        temp_df['floor_number'],
        bins=[0, 5, 10, 20, 30, 50],
        labels=['1-5', '6-10', '11-20', '21-30', '31+']
    )

    # Get top 5 towns by transaction volume for first subplot
    top_towns = df['town'].value_counts().head(5).index.tolist()
    df_top_towns = df[df['town'].isin(top_towns)]

    # Create subplot
    fig = make_subplots(
        rows=1,
        cols=2,
        subplot_titles=[
            "Price Distribution by Flat Type and Town (Top 5 Towns by Volume)",
            "Price vs Floor Number by Flat Type with Trend Lines"
        ],
        horizontal_spacing=0.15,
        column_widths=[0.5, 0.5]
    )

    # Add traces for first subplot (Price Distribution by Flat Type and Town)
    for i, town in enumerate(top_towns):
        town_data = df_top_towns[df_top_towns['town'] == town]

        fig.add_trace(
            go.Box(
                x=town_data['flat_type'],
                y=town_data['resale_price'],
                name=town,
                marker_color=hex_codes[i % len(hex_codes)])
        ),
        row=1, col=1
    )

    # Add traces for second subplot (Price vs Floor Number by Flat Type with Trend Lines)
    # Create an ordered list of flat types for consistent coloring
    ordered_flat_types = sorted(temp_df['flat_type'].unique())

    for i, flat_type in enumerate(ordered_flat_types):
        flat_data = temp_df[temp_df['flat_type'] == flat_type]

        # Sample data for larger categories to prevent overcrowding
        if len(flat_data) > 500:
            flat_data = flat_data.sample(500, random_state=42)

        # Add scatter plot
        fig.add_trace(
            go.Scatter(
                x=flat_data['floor_number'],
                y=flat_data['resale_price'],
                mode='markers',
                name=flat_type,
                marker=dict(
                    size=6, # Reduced size for better fit in subplot
                    opacity=0.6,
                    color=hex_codes[i % len(hex_codes)])
            ),
            legendgroup=flat_type
        )

```

```

        ),
        row=1, col=2
    )

# Add trend line using linear regression
if len(flat_data) > 10: # Only add trendline if enough data points
    x = flat_data['floor_number']
    y = flat_data['resale_price']

    # Calculate trendline
    z = np.polyfit(x, y, 1)
    p = np.poly1d(z)

    # Generate points for the trendline
    x_trend = np.array([min(x), max(x)])
    y_trend = p(x_trend)

    fig.add_trace(
        go.Scatter(
            x=x_trend,
            y=y_trend,
            mode='lines',
            name=f'{flat_type} Trend',
            line=dict(
                color=hex_codes[i % len(hex_codes)],
                width=2,
                dash='dash'
            ),
            legendgroup=flat_type,
            showlegend=False
        ),
        row=1, col=2
    )

# Update layout
fig.update_layout(
    height=800,
    width=1400,
    title={
        'text': "Feature Interaction Analysis: Price by Flat Type, Town, and Floor Level",
        'font': {
            'family': "Arial, sans-serif",
            'size': 24,
            'weight': 'bold'
        },
        'y': 0.98,
        'x': 0.5,
        'xanchor': 'center',
    },
    template="plotly_white",
    legend=dict(
        orientation="h",
        yanchor="bottom",
        y=1.08,
        xanchor="center",
        x=0.5,
        font=dict(size=10)
    ),
    margin=dict(t=150, b=100, l=100, r=50),
)

# Update axes for first subplot
fig.update_xaxes(
    title=dict(
        text="Flat Type",
        font=dict(
            family="Arial, sans-serif",
            size=14,
            weight='bold'
        )
    ),
    categoryorder='array',
    # Sort flat types in logical order (from smallest to largest)
    categoryarray=sorted(df['flat_type'].unique(),
                         key=lambda x: ['1 ROOM', '2 ROOM', '3 ROOM', '4 ROOM', '5 ROOM', 'EXECUTIVE', 'MULTI-GENERATION'].index(x)
                         if x in ['1 ROOM', '2 ROOM', '3 ROOM', '4 ROOM', '5 ROOM', 'EXECUTIVE', 'MULTI-GENERATION'] else 999),
    row=1, col=1
)

fig.update_yaxes(
    title=dict(
        text="Resale Price (SGD)",
        font=dict(
            family="Arial, sans-serif",
            size=14,
            weight='bold'
        )
    ),
    tickprefix="$",
    tickformat=",",
    row=1, col=1
)

# Update axes for second subplot
fig.update_xaxes(
    title=dict(
        text="Floor Number",
        font=dict(
            family="Arial, sans-serif",
            size=14,
            weight='bold'
        )
    ),
    # Set a reasonable range with some padding
    range=[-1, max(temp_df['floor_number']) + 1],
    row=1, col=2
)

fig.update_yaxes(
    title=dict(
        text="Resale Price (SGD)",
        font=dict(
            family="Arial, sans-serif",
            size=14,
            weight='bold'
        )
    ),
    tickprefix="$",
    tickformat=",",
    row=1, col=2
)

# Add annotation for the floor trend
fig.add_annotation(
    x=0.5, # centered in the second subplot
    y=-0.2,
    xref='paper',
    yref='paper',
    text='Higher floors generally command higher prices across all flat types',
    showarrow=False,
    font=dict(size=12),
    align='center'
)

fig.show()

```

```

# -----
# 2. Improved Average Price by Flat Type and Floor Range & Price per SQM
#
if 'flat_type' in df.columns and 'storey_range' in df.columns and 'price_per_sqm' in df.columns:
    # Create a temporary dataframe
    temp_df = df.copy()

    # Extract the lower floor number and convert to integer
    temp_df['floor_number'] = temp_df['storey_range'].str.extract(r'^(\d+)').iloc[:, 0].astype(int)

    # Create floor category based on the extracted floor number
    temp_df['floor_category'] = pd.cut(
        temp_df['floor_number'],
        bins=[0, 5, 10, 20, 30, 50],
        labels=['1-5', '6-10', '11-20', '21-30', '31+']
    )

    # Calculate average price by flat type and floor category
    avg_by_type_floor = temp_df.groupby(['flat_type', 'floor_category'])['resale_price'].mean().reset_index()

    # Create a pivot table for the heatmap
    pivot_avg_price = avg_by_type_floor.pivot(index='flat_type', columns='floor_category', values='resale_price')

    # Calculate average price per sqm by flat type and floor category
    avg_psm_by_type_floor = temp_df.groupby(['flat_type', 'floor_category'])['price_per_sqm'].mean().reset_index()

    # Create a pivot table for the heatmap
    pivot_avg_psm = avg_psm_by_type_floor.pivot(index='flat_type', columns='floor_category', values='price_per_sqm')

    # Sort flat types in logical order
    flat_type_order = ['1 ROOM', '2 ROOM', '3 ROOM', '4 ROOM', '5 ROOM', 'EXECUTIVE', 'MULTI-GENERATION']

    # Filter to keep only the flat types in our predefined order
    pivot_avg_price = pivot_avg_price.reindex(
        [ft for ft in flat_type_order if ft in pivot_avg_price.index]
    )

    pivot_avg_psm = pivot_avg_psm.reindex(
        [ft for ft in flat_type_order if ft in pivot_avg_psm.index]
    )

    # Create 2 subplots: heatmap for avg price and heatmap for price per sqm
    fig = make_subplots(
        rows=1,
        cols=2,
        subplot_titles=(
            "Average Resale Price by Flat Type and Floor Range",
            "Average Price per SQM by Flat Type and Floor Range"
        ),
        horizontal_spacing=0.2
    )

    # Set colorscale for heatmap
    mako_colors = []
    for i in np.linspace(0, 1, 10):
        rgba = mako_cmap(i)
        # Convert RGB to hex for Plotly
        hex_color = mpl.colors.rgb2hex(rgba)
        mako_colors.append([i, hex_color])

    # Create second heatmap
    mako_colors_reversed = []
    for i in np.linspace(0, 1, 10):
        rgba = mako_cmap(1-i)
        # Convert RGB to hex for Plotly
        hex_color = mpl.colors.rgb2hex(rgba)
        mako_colors_reversed.append([i, hex_color])

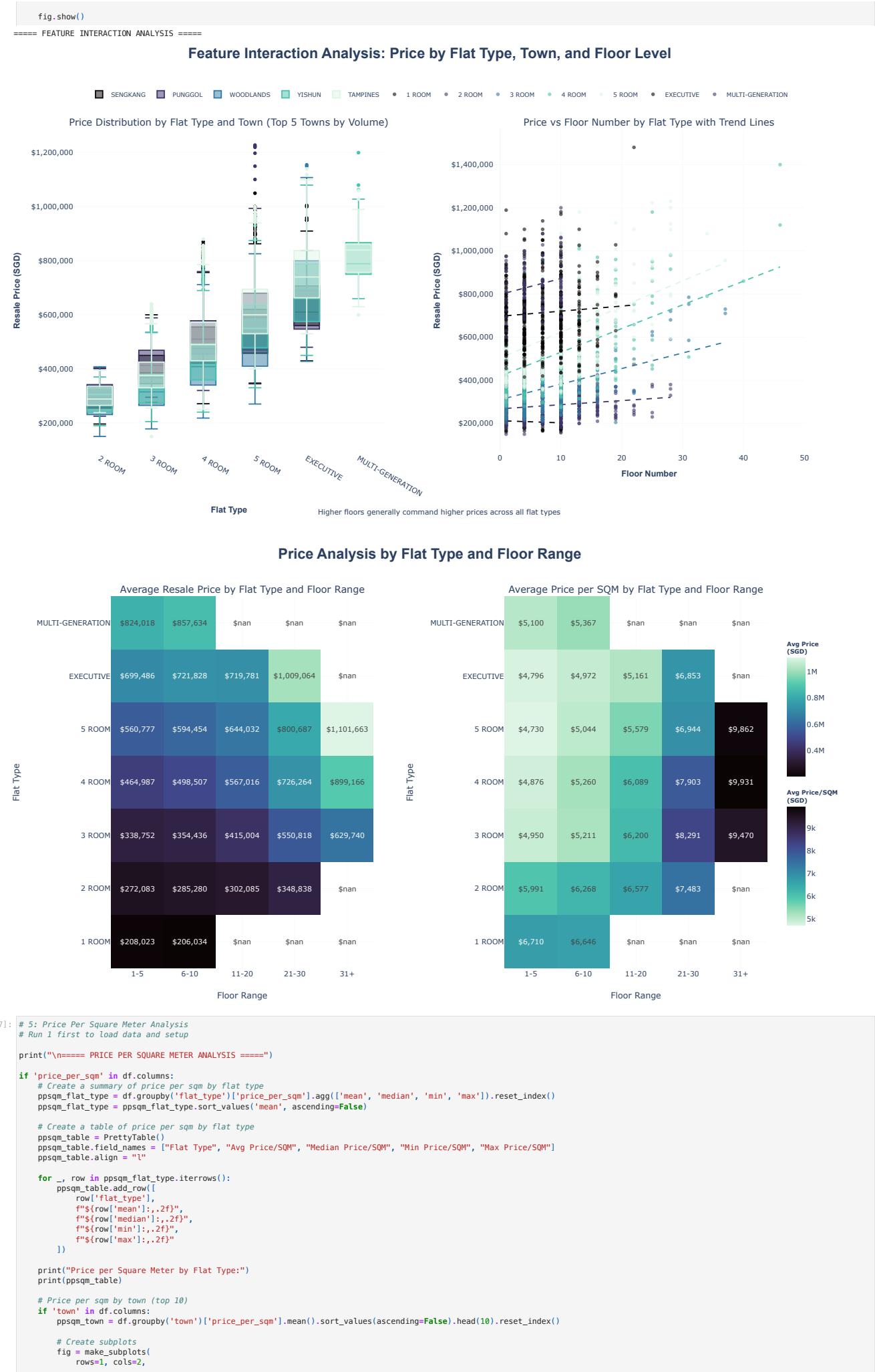
    # Add heatmap for average price
    fig.add_trace(
        go.Heatmap(
            z=pivot_avg_price.values,
            x=pivot_avg_price.columns,
            y=pivot_avg_price.index,
            colorscale=mako_colors,
            text=[[f"${val:,.0f}" for val in row] for row in pivot_avg_price.values],
            texttemplate="%{text}",
            colorbar=dict(
                title="Avg Price<br>(SGD)",
                len=0.4,
                y=0.7,
                title_font=dict(size=10, weight="bold")
            ),
            hoverinfo="x+y+z"
        ),
        row=1, col=1
    )

    # Add heatmap for average price per sqm
    fig.add_trace(
        go.Heatmap(
            z=pivot_avg_psm.values,
            x=pivot_avg_psm.columns,
            y=pivot_avg_psm.index,
            colorscale=mako_colors_reversed,
            text=[[f"${val:,.0f}" for val in row] for row in pivot_avg_psm.values],
            texttemplate="%{text}",
            colorbar=dict(
                title="Avg Price/SQM<br>(SGD)",
                len=0.4,
                y=0.3,
                title_font=dict(size=10, weight="bold")
            ),
            hoverinfo="x+y+z"
        ),
        row=1, col=2
    )

    # Update layout
    fig.update_layout(
        height=800,
        width=1400,
        title={
            'text': 'Price Analysis by Flat Type and Floor Range',
            'font': {
                'family': 'Arial, sans-serif',
                'size': 24,
                'weight': 'bold'
            },
            'x': 0.5,
            'xanchor': 'center',
            'y': 0.95
        },
        template="plotly_white",
        margin=dict(t=120, b=60, l=80, r=80)
    )

    fig.update_xaxes(type="category", title="Floor Range", row=1, col=1, title_standoff=25)
    fig.update_yaxes(type="category", title="Flat Type", row=1, col=1, title_standoff=25, automargin=True)
    fig.update_xaxes(type="category", title="Floor Range", row=1, col=2, title_standoff=25)
    fig.update_yaxes(type="category", title="Flat Type", row=1, col=2, title_standoff=25, automargin=True)

```



```

        subplot_titles=(
            "Average Price per Square Meter by Flat Type",
            "Top 10 Towns by Price per Square Meter"
        ),
        horizontal_spacing=0.1
    )

    # Set color palette
    mako_palette = sns.color_palette("mako_r", as_cmap=True)
    mako_colors = [mpl.colors.rgb2hex(c) for c in mako_palette(np.linspace(0, 1, len(ppsqm_flat_type)))]
```

Create bar chart for Price per sqm by flat type

```

for i, (_, row) in enumerate(ppsqm_flat_type.iterrows()):
    # Use the mako color palette
    color_idx = i % len(mako_colors)

    fig.add_trace(
        go.Bar(
            x=[row['flat_type']],
            y=[row['mean']],
            error_y=dict(
                type='data',
                array=[row['max'] - row['mean']],
                arrayminus=[row['mean'] - row['min']],
                visible=True
            ),
            name=row['flat_type'],
            marker_color=mako_colors[color_idx],
            showlegend=False
        ),
        row=1, col=1
    )
```

Create bar chart for Price per sqm by town

```

for i, (_, row) in enumerate(ppsqm_town.iterrows()):
    # Use the mako color palette
    color_idx = i % len(mako_colors)

    fig.add_trace(
        go.Bar(
            x=[row['town']],
            y=[row['price_per_sqm']],
            name=row['town'],
            marker_color=mako_colors[color_idx],
            showlegend=False
        ),
        row=1, col=2
    )
```

Update layout

```

fig.update_layout(
    title={
        'text': "Price Per Square Meter Analysis",
        'x': 0.5,
        'xanchor': 'center',
        'y': 0.95,
        'font': {
            'family': "Arial, sans-serif",
            'size': 24,
            'weight': 'bold'
        }
    },
    template="plotly_white",
    height=800,
    width=1400,
    showLegend=False,
    margin=dict(l=60, r=60, t=120, b=120)
)
```

Format axes for flat type chart

```

fig.update_xaxes(
    title_text="Flat Type",
    title_font={
        'family': "Arial, sans-serif",
        'size': 14,
        'weight': 'bold'
    },
    row=1, col=1,
    categoryorder='array',
    categoryarray=ppsqm_flat_type['flat_type']
)
fig.update_yaxes(
    title_text="Average Price per Sqm (SGD)",
    title_font={
        'family': "Arial, sans-serif",
        'size': 14,
        'weight': 'bold'
    },
    tickprefix="$",
    tickformat=",",
    row=1, col=1
)
```

Format axes for town chart

```

fig.update_xaxes(
    title_text="Town",
    title_font={
        'family': "Arial, sans-serif",
        'size': 14,
        'weight': 'bold'
    },
    tickangle=-45,
    row=1, col=2
)
fig.update_yaxes(
    title_text="Average Price per Sqm (SGD)",
    title_font={
        'family': "Arial, sans-serif",
        'size': 14,
        'weight': 'bold'
    },
    tickprefix="$",
    tickformat=",",
    row=1, col=2
)
```

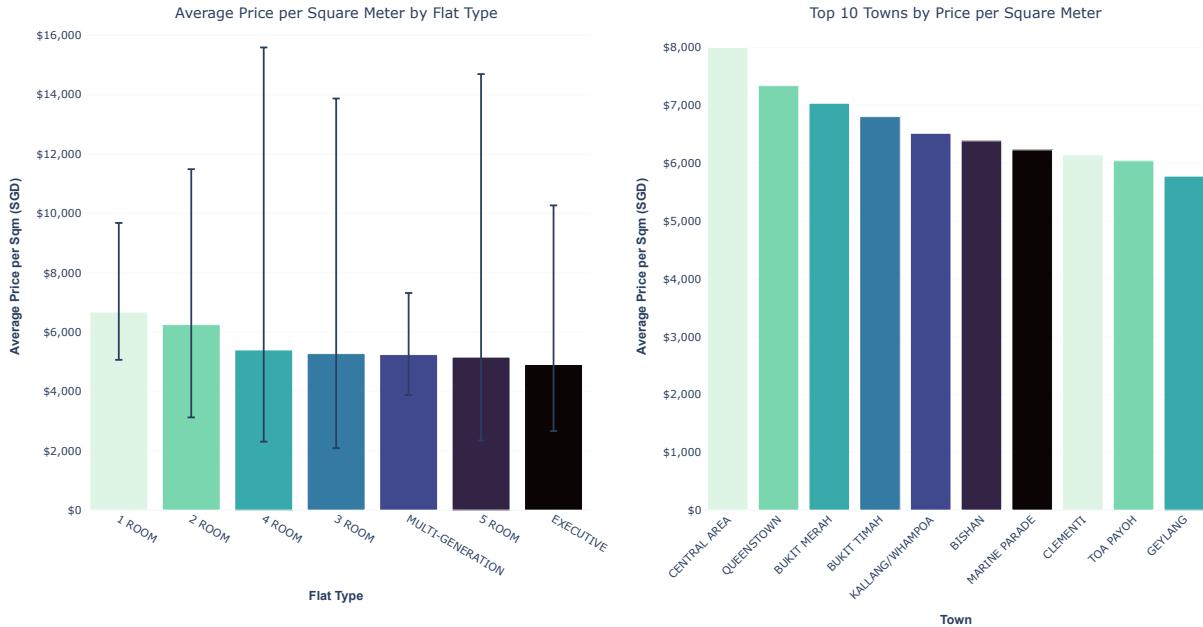
fig.show()

===== PRICE PER SQUARE METER ANALYSIS =====

Price per Square Meter by Flat Type:

Flat Type	Avg Price/SQM	Median Price/SQM	Min Price/SQM	Max Price/SQM
1 ROOM	\$6,678.78	\$6,612.90	\$5,064.52	\$9,677.42
2 ROOM	\$6,252.15	\$6,127.66	\$3,125.00	\$11,489.36
4 ROOM	\$5,399.14	\$5,059.52	\$2,307.69	\$15,591.40
3 ROOM	\$5,276.07	\$5,044.78	\$2,089.55	\$13,870.97
MULTI-GENERATION	\$5,243.24	\$5,120.48	\$2,882.68	\$7,317.07
5 ROOM	\$5,157.61	\$4,913.79	\$2,347.83	\$14,694.10
EXECUTIVE	\$4,908.23	\$4,791.67	\$2,666.67	\$10,269.50

Price Per Square Meter Analysis



```
In [9]: # 6: Advanced Time Series Analysis
# Run 1 first to load data and setup

print("\n===== ADVANCED TIME SERIES ANALYSIS =====")

# -----
# 1. Time series analysis
# -----

# Set color palette
mako_palette = sns.color_palette("mako_r", as_cmap=True)
mako_colors = [mpl.colors.rgb2hex(c) for c in mako_palette(np.linspace(0, 1, 5))]

if 'month' in df.columns:
    # Ensure month column is datetime format
    if not pd.api.types.is_datetime64_any_dtype(df['month']):
        try:
            df['month'] = pd.to_datetime(df['month'])
            print("Converted 'month' column to datetime format")
        except Exception as e:
            print(f"Error converting 'month' column to datetime: {e}")

# Group by month and calculate statistics with error handling
try:
    # Group by month and calculate statistics
    monthly_stats = df.groupby(pd.Grouper(key='month', freq='M')).agg({
        'resale_price': ['mean', 'median', 'count', 'std'],
        'price_per_sqm': ['mean', 'median']
    }).reset_index()

    # Flatten multi-level columns
    monthly_stats.columns = ['month'] + ['_'.join(col).strip() for col in monthly_stats.columns[1:]]

    # Add coefficient of variation for price volatility analysis
    monthly_stats['price_cv'] = (monthly_stats['resale_price_std'] / monthly_stats['resale_price_mean']) * 100

    # Calculate rolling averages with configurable window
    window_size = 6 # Can be adjusted based on analysis needs
    monthly_stats['price_ma'] = monthly_stats['resale_price_mean'].rolling(
        window=window_size, center=True).mean()
    monthly_stats['vol_ma'] = monthly_stats['resale_price_count'].rolling(
        window=window_size, center=True).mean()

    # Calculate YoY growth rates if enough data
    if len(monthly_stats) >= 13:
        monthly_stats['yoy_price_change'] = monthly_stats['resale_price_mean'].pct_change(12) * 100
        monthly_stats['yoy_volume_change'] = monthly_stats['resale_price_count'].pct_change(12) * 100

    # Create time series with multiple metrics
    fig = make_subplots(
        rows=2, cols=1,
        subplot_titles=("Price Trends Over Time", "Transaction Volume and Volatility"),
        specs=[[{"secondary_y": True}], [{"secondary_y": True}]],
        vertical_spacing=0.20 # Increased vertical spacing to separate plots clearly
    )

    # Add average price line
    fig.add_trace(
        go.Scatter(
            x=monthly_stats['month'],
            y=monthly_stats['resale_price_mean'],
            mode='lines+markers',
            name='Average Price',
            line=dict(color=mako_colors[0], width=2),
            marker=dict(size=6)
        ),
        row=1, col=1,
        secondary_y=False
    )

    # Add moving average
    fig.add_trace(
        go.Scatter(
            x=monthly_stats['month'],
            y=monthly_stats['price_ma'],
            mode='lines',
            name=f'{window_size}-Month MA (Price)',
            line=dict(color=mako_colors[1], width=3, dash='dash')
        ),
        row=1, col=1,
        secondary_y=False
    )

    # Add price per sqm line on secondary y-axis
    fig.add_trace(

```

```

        go.Scatter(
            x=monthly_stats['month'],
            y=monthly_stats['price_per_sqm_mean'],
            mode='lines',
            name='Avg Price per Sqm',
            line=dict(color=mako_colors[2], width=2)
        ),
        row=1, col=1,
        secondary_y=True
    )

    # Add YoY growth if available
    if 'yoy_price_change' in monthly_stats.columns:
        # Filter out NaN values
        yoy_data = monthly_stats.dropna(subset=['yoy_price_change'])
        if len(yoy_data) > 0:
            fig.add_trace(
                go.Scatter(
                    x=yoy_data['month'],
                    y=yoy_data['yoy_price_change'],
                    mode='lines',
                    name='YoY Price Change %',
                    line=dict(color=mako_colors[4], width=2),
                    visible='legendonly' # Hidden by default to avoid cluttering
                ),
                row=1, col=1,
                secondary_y=True
            )
    )

    # Add transaction volume
    fig.add_trace(
        go.Bar(
            x=monthly_stats['month'],
            y=monthly_stats['resale_price_count'],
            name='Transaction Volume',
            marker_color=mako_colors[3],
            opacity=0.7
        ),
        row=2, col=1,
        secondary_y=False
    )

    # Add moving average for volume
    fig.add_trace(
        go.Scatter(
            x=monthly_stats['month'],
            y=monthly_stats['vol_ma'],
            mode='lines',
            name=f'{window_size}-Month MA (Volume)',
            line=dict(color=mako_colors[4], width=3)
        ),
        row=2, col=1,
        secondary_y=False
    )

    # Add price coefficient of variation on secondary y-axis (volatility measure)
    fig.add_trace(
        go.Scatter(
            x=monthly_stats['month'],
            y=monthly_stats['price_cv'],
            mode='lines',
            name='Price Volatility (CV%)',
            line=dict(color=mako_colors[2], width=2, dash='dot')
        ),
        row=2, col=1,
        secondary_y=True
    )

    # Update overall layout settings
    fig.update_layout(
        title={
            'text': "Time Series Analysis of Singapore HDB Resale Market",
            'x': 0.5,
            'xanchor': 'center',
            'font': {
                'family': "Arial, sans-serif",
                'size': 24,
                'weight': 'bold'
            }
        },
        template="plotly_white",
        showLegend=True,
        legend=dict(
            orientation="h",
            yanchor="top",
            y=0.53,
            xanchor="center",
            x=0.5,
            font=dict(size=10)
        ),
        margin=dict(t=100, b=150, l=80, r=80),
        height=1000,
        width=1400,
        hovermode="x unified",
        plot_bgcolor='rgba(240,240,240,0.2)'
    )

    # Add x-axis title with extra padding
    fig.update_xaxes(
        title_text="Date",
        title_font={
            'family': "Arial, sans-serif",
            'size': 14,
            'weight': 'bold'
        },
        title_standoff=25,
        row=2, col=1,
        gridcolor='rgba(211,211,211,0.3)'
    )

    # Update y-axes
    fig.update_yaxes(
        title_text="Average Price (SGD)",
        title_font={
            'family': "Arial, sans-serif",
            'size': 14,
            'weight': 'bold'
        },
        secondary_y=False,
        row=1, col=1,
        tickprefix="$",
        tickformat=",",
        gridcolor='rgba(211,211,211,0.3)'
    )

    fig.update_yaxes(
        title_text="Price per Sqm (SGD)",
        title_font={
            'family': "Arial, sans-serif",
            'size': 14,
            'weight': 'bold'
        },
        secondary_y=True,
        row=1, col=1,
        tickprefix="$",
        tickformat=","
    )

```

```

        tickprefix="$",
        tickformat=",",
        gridcolor='rgba(211,211,211,0.3)'
    )

    fig.update_yaxes(
        title_text="Number of Transactions",
        title_font={
            'family': "Arial, sans-serif",
            'size': 14,
            'weight': 'bold'
        },
        secondary_y=False,
        row=2, col=1,
        gridcolor='rgba(211,211,211,0.3)'
    )

    fig.update_yaxes(
        title_text="Coefficient of Variation (%)",
        title_font={
            'family': "Arial, sans-serif",
            'size': 14,
            'weight': 'bold'
        },
        secondary_y=True,
        row=2, col=1,
        gridcolor='rgba(211,211,211,0.3)'
    )

# Add range slider for better time navigation
fig.update_xaxes(
    rangeslider_visible=True,
    row=2, col=1
)

fig.show()

# -----
# 2. Time series decomposition
# -----
try:
    # Set the month as index for decomposition
    ts_df = monthly_stats.set_index('month')

    # Perform decomposition on the average price
    decomposition = seasonal_decompose(
        ts_df['resale_price_mean'].dropna(),
        model='additive',
        period=12 # 12 months for annual seasonality
    )

    # Create decomposition plot
    fig_decomp = make_subplots(
        rows=2, cols=2,
        subplot_titles=("Original Time Series", "Trend Component",
                        "Seasonal Component", "Residual Component"),
        vertical_spacing=0.10,
        horizontal_spacing=0.10,
        shared_xaxes=False
    )

    # Use the same color for both Original and Trend
    trend_color = mako_colors[1]

    # Add original data (top-left)
    fig_decomp.add_trace(
        go.Scatter(
            x=decomposition.observed.index,
            y=decomposition.observed,
            mode='lines',
            name='Original',
            line=dict(color=trend_color, width=2)
        ),
        row=1, col=1
    )

    # Add trend (top-right)
    fig_decomp.add_trace(
        go.Scatter(
            x=decomposition.trend.index,
            y=decomposition.trend,
            mode='lines',
            name='Trend',
            line=dict(color=trend_color, width=2)
        ),
        row=1, col=2
    )

    # Add seasonal (bottom-left)
    fig_decomp.add_trace(
        go.Scatter(
            x=decomposition.seasonal.index,
            y=decomposition.seasonal,
            mode='lines',
            name='Seasonal',
            line=dict(color=mako_colors[2], width=2)
        ),
        row=2, col=1
    )

    # Add residual (bottom-right)
    fig_decomp.add_trace(
        go.Scatter(
            x=decomposition.resid.index,
            y=decomposition.resid,
            mode='lines',
            name='Residual',
            line=dict(color=mako_colors[3], width=2)
        ),
        row=2, col=2
    )

    # Add zero reference line for seasonal and residual components
    # Bottom-left (seasonal)
    fig_decomp.add_shape(
        type="line",
        x0=decomposition.observed.index[0],
        y0=0,
        x1=decomposition.observed.index[-1],
        y1=0,
        line=dict(color="gray", width=1, dash="dash"),
        row=2, col=1
    )

    # Bottom-right (residual)
    fig_decomp.add_shape(
        type="line",
        x0=decomposition.observed.index[0],
        y0=0,
        x1=decomposition.observed.index[-1],
        y1=0,
        line=dict(color="gray", width=1, dash="dash"),
        row=2, col=2
    )

```

```

# Update layout
fig_decomp.update_layout(
    title={
        'text': "<b>Time Series Decomposition of HDB Resale Prices</b>",
        'x': 0.5,
        'xanchor': 'center',
        'font': {'size': 22}
    },
    height=800,
    width=1400,
    showLegend=True,
    legend=dict(
        orientation="h",
        yanchor="top",
        y=-0.15,
        xanchor="center",
        x=0.5
    ),
    template="plotly_white",
    plot_bgcolor='rgba(240,240,240,0.2)'
)

# Improve formatting for all subplots
# Format for original and trend (top row)
for col in [1, 2]:
    fig_decomp.update_yaxes(
        title_font={
            'family': "Arial, sans-serif",
            'size': 14,
            'weight': 'bold'
        },
        gridcolor='rgba(211,211,211,0.3)',
        tickprefix="",
        tickformat=",",
        row=1, col=col
    )

    fig_decomp.update_xaxes(
        title_font={
            'family': "Arial, sans-serif",
            'size': 14,
            'weight': 'bold'
        },
        gridcolor='rgba(211,211,211,0.3)',
        row=1, col=col
    )

# Format for seasonal and residual (bottom row)
for col in [1, 2]:
    fig_decomp.update_yaxes(
        title_font={
            'family': "Arial, sans-serif",
            'size': 14,
            'weight': 'bold'
        },
        gridcolor='rgba(211,211,211,0.3)',
        row=2, col=col
    )

    fig_decomp.update_xaxes(
        title_font={
            'family': "Arial, sans-serif",
            'size': 14,
            'weight': 'bold'
        },
        gridcolor='rgba(211,211,211,0.3)',
        title_text="Date",
        row=2, col=col
    )

fig_decomp.show()

# Enhanced seasonal analysis with PrettyTable
from prettytable import PrettyTable

# Create a PrettyTable for displaying peak/trough data
summary_table = PrettyTable()
summary_table.field_names = ["Metric", "Month", "Value ($)"]
summary_table.align = "l"

seasonal_pattern = decomposition.seasonal[:12] # Take one full cycle
max_month_idx = seasonal_pattern.argmax()
min_month_idx = seasonal_pattern.argmin()

# Add peak month data
peak_month = seasonal_pattern.index[max_month_idx].strftime('%B')
peak_value = seasonal_pattern.max()
summary_table.add_row(["Peak Month", peak_month, f"${peak_value:.2f}"])

# Add trough month data
trough_month = seasonal_pattern.index[min_month_idx].strftime('%B')
trough_value = seasonal_pattern.min()
summary_table.add_row(["Trough Month", trough_month, f"${trough_value:.2f}"])

# Add seasonal amplitude
amplitude = abs(peak_value - trough_value)
summary_table.add_row(["Seasonal Amplitude", "", f"${amplitude:.2f}"])

print("\n===== Seasonal Price Patterns =====")
print(summary_table)

# Calculate and show average monthly effect with PrettyTable
monthly_effects = {}
for month in range(1, 13):
    month_data = decomposition.seasonal[decomposition.seasonal.index.month == month]
    if not month_data.empty:
        monthly_effects[month] = month_data.mean()

if monthly_effects:
    # Create PrettyTable for monthly effects
    monthly_table = PrettyTable()
    monthly_table.field_names = ["Month", "Seasonal Effect ($)", "Direction"]
    monthly_table.align = "l"

    # Add data to table
    for month, effect in sorted(monthly_effects.items()):
        month_name = pd.Timestamp(year=2020, month=month, day=1).strftime('%B')
        direction = "+" if effect > 0 else "-"
        monthly_table.add_row([month_name, f"${effect:.2f}", direction])

    print("\n===== Average Seasonal Effect by Month =====")
    print(monthly_table)

# -----
# 1. Autocorrelation analysis for cyclical patterns
# -----
from statsmodels.tsa.stattools import acf

# Compute ACF values
acf_values = acf(decomposition.observed.dropna(), nlags=36)
lags = list(range(len(acf_values))) # Convert range to list for Plotly
confidence_interval = 1.96 / np.sqrt(len(decomposition.observed.dropna()))

```

```

# Create Plotly figure for ACF
fig_acf = go.Figure()

# Add confidence intervals as shapes
fig_acf.add_shape(
    type="line",
    x0=0, x1=36,
    y0=confidence_interval, y1=confidence_interval,
    line=dict(color="gray", width=1, dash="dash")
)

fig_acf.add_shape(
    type="line",
    x0=0, x1=36,
    y0=-confidence_interval, y1=-confidence_interval,
    line=dict(color="gray", width=1, dash="dash")
)

# Add ACF bars
normalized_values = (acf_values + 1) / 2 # Normalize to 0-1 range for color mapping

# Convert colors to hex for Plotly
bar_colors = []
for val in normalized_values:
    rgba_color = mako_palette(val)
    hex_color = mpl.colors.rgb2hex((rgba_color[0], rgba_color[1], rgba_color[2]))
    bar_colors.append(hex_color)

# Add the bars
fig_acf.add_trace(
    go.Bar(
        x=lags,
        y=acf_values,
        marker_color=bar_colors,
        name="ACF Values"
    )
)

# Update layout
fig_acf.update_layout(
    title={
        'text': "Autocorrelation Function (ACF) of HDB Resale Prices",
        'font': {
            'family': "Arial, sans-serif",
            'size': 24,
            'weight': 'bold'
        },
        'y': 0.95,
        'x': 0.5,
        'xanchor': 'center'
    },
    xaxis_title="Lag (Months)",
    yaxis_title="Correlation",
    template="plotly_white",
    height=800,
    width=1400,
    plot_bgcolor='rgba(240,240,240,0.2)',
    yaxis=dict(
        range=[-1.1, 1.1],
        gridcolor='rgba(211,211,211,0.3)'
    ),
    xaxis=dict(
        tickmode='linear',
        dtick=6,
        gridcolor='rgba(211,211,211,0.3)'
    )
)

# Show the ACF chart
fig_acf.show()

# 4. Explain what ACF is and why it's useful
print("\n===== AUTOCORRELATION ANALYSIS =====")
print("The Autocorrelation Function (ACF) measures how prices correlate with their own past values.")
print("This helps identify:")
print("1. Market momentum: Strong positive values at lag 1-3 indicate price momentum")
print("2. Seasonal patterns: Peaks at regular intervals (e.g., 12 months) suggest seasonality")
print("3. Market efficiency: Quick decay to zero suggests an efficient market")

# Identify key autocorrelation patterns
# Find significant lags (those outside confidence interval)
significant_lags = [i for i in range(1, len(acf_values)) if abs(acf_values[i]) > confidence_interval]
max_acf_lag = acf_values[1:].argmax() + 1 # Skip lag 0 (always 1.0)
seasonal_lag = None

# Look for seasonal patterns (often around 12 months in annual data)
if 12 in significant_lags and acf_values[12] > confidence_interval:
    seasonal_lag = 12
elif 6 in significant_lags and acf_values[6] > confidence_interval:
    seasonal_lag = 6

if significant_lags:
    print(f"\nSignificant autocorrelation found at lags: {', '.join(map(str, significant_lags))}")

if seasonal_lag:
    print(f"\nSeasonal pattern detected with {seasonal_lag}-month periodicity")

# Market behavior analysis
if acf_values[1] > 0.7:
    print("\nThe market shows very strong persistence, suggesting momentum strategies may be effective.")
elif acf_values[1] > 0.5:
    print("\nThe market shows strong persistence, with price trends tending to continue.")
elif acf_values[1] > 0.3:
    print("\nThe market shows moderate persistence, with some predictability in price movements.")
else:
    print("\nThe market shows weak persistence, suggesting prices quickly adjust to new information.")

# Create a PrettyTable for Key Findings right after the market behavior text
findings_table = PrettyTable()
findings_table.field_names = ["Key Findings from Autocorrelation Analysis"]
findings_table.align = "l" # Left-align

# Add findings as bullet points
findings = []

# Add short-term autocorrelation
if max_acf_lag <= 3 and acf_values[max_acf_lag] > confidence_interval:
    findings.append(f"\n• Strong short-term correlation at {max_acf_lag} month(s) (r={acf_values[max_acf_lag]:.2f})")

# Add seasonal pattern if found
if seasonal_lag:
    findings.append(f"\n• Seasonal pattern detected at {seasonal_lag} months (r={acf_values[seasonal_lag]:.2f})")

# Add general market momentum information
if acf_values[1] > 0.7:
    findings.append("• Very strong market momentum (prices highly persistent)")
elif acf_values[1] > 0.5:
    findings.append("• Strong market momentum (prices moderately persistent)")
elif acf_values[1] > 0.3:
    findings.append("• Moderate market momentum")
else:
    findings.append("• Weak market momentum (prices quickly revert)")



```

```

# Add findings to table
for finding in findings:
    findings_table.add_row([finding])

# Print the findings table right after the market behavior text
print("\n") # Add a blank line
print(findings_table)

except Exception as e:
    print("Error in time series decomposition: {e}")
    print("Suggestions for troubleshooting:")
    print("1. Check for sufficient data points (need at least 2*period points)")
    print("2. Check for missing values in the time series")
    print("3. Ensure time series is regularly spaced")

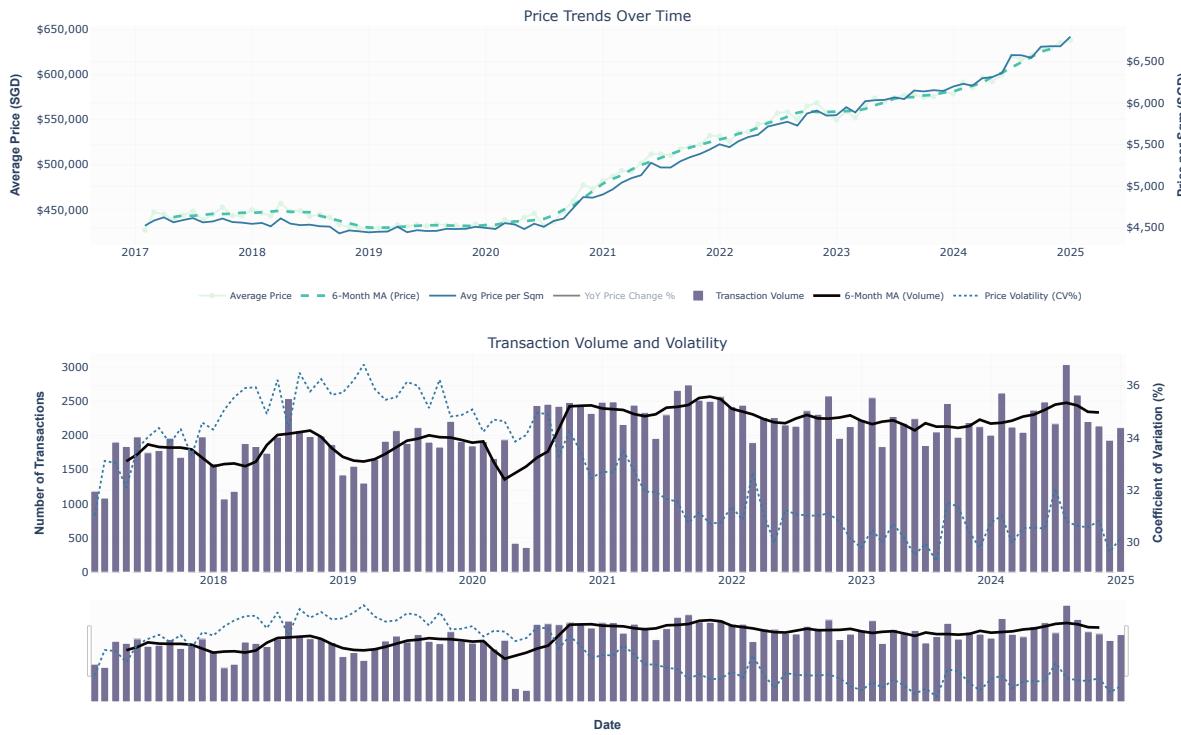
except Exception as e:
    print("Error in monthly statistics calculation: {e}")

else:
    print("No 'month' column found in the dataset. Time series analysis skipped.")
    print("Available columns:", df.columns.tolist())

```

===== ADVANCED TIME SERIES ANALYSIS =====

Time Series Analysis of Singapore HDB Resale Market



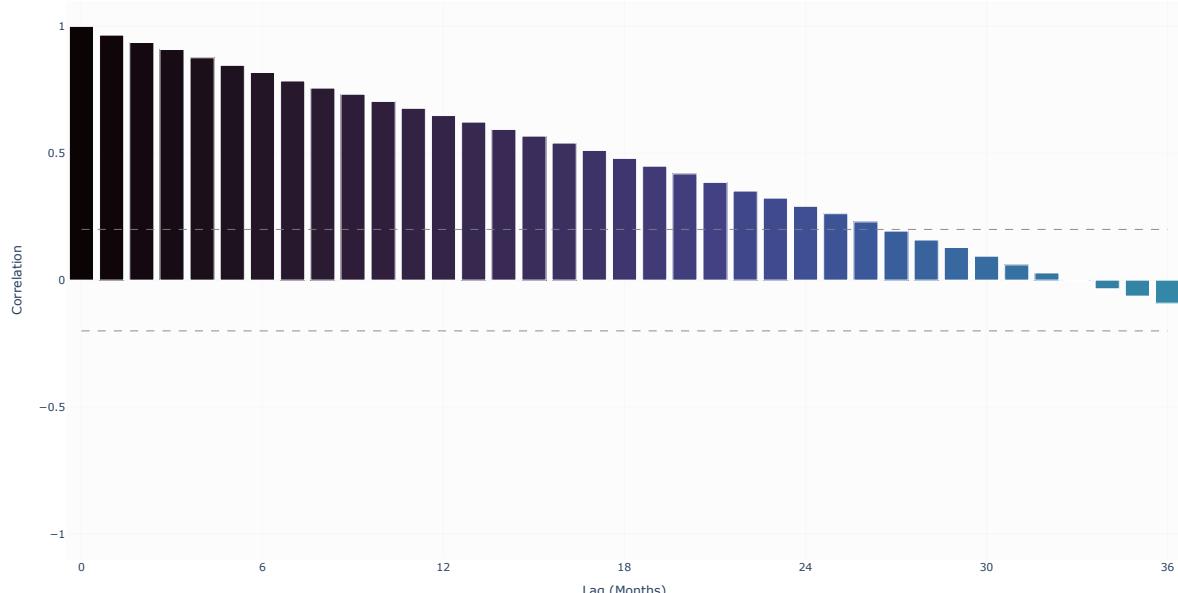
Time Series Decomposition of HDB Resale Prices



===== Seasonal Price Patterns =====		
Metric	Month	Value (\$)
Peak Month	May	\$1954.04
Trough Month	February	\$-2810.92
Seasonal Amplitude		\$4764.97

===== Average Seasonal Effect by Month =====		
Month	Seasonal Effect (\$)	Direction
January	\$-1779.75	↓
February	\$-2810.92	↓
March	\$789.60	↑
April	\$982.47	↑
May	\$1954.04	↑
June	\$1743.77	↑
July	\$-487.32	↓
August	\$-413.03	↓
September	\$1557.52	↑
October	\$1100.73	↑
November	\$-260.80	↓
December	\$-2376.32	↓

Autocorrelation Function (ACF) of HDB Resale Prices



===== AUTOCORRELATION ANALYSIS =====
The Autocorrelation Function (ACF) measures how prices correlate with their own past values.
This helps identify:
1. Market momentum: Strong positive values at lag 1-3 indicate price momentum
2. Seasonal patterns: Peaks at regular intervals (e.g., 12 months) suggest seasonality
3. Market efficiency: Quick decay to zero suggests an efficient market

Significant autocorrelation found at lags: 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26
Seasonal pattern detected with 12-month periodicity

The market shows very strong persistence, suggesting momentum strategies may be effective.

Key Findings from Autocorrelation Analysis	
• Strong short-term correlation at 1 month(s) ($r=0.97$)	
• Seasonal pattern detected at 12 months ($r=0.65$)	
• Very strong market momentum (prices highly persistent)	

```
In [10]: # 7: Advanced Pricing Analysis
# Run 1 first to load data and setup
print("\n==== Advanced Pricing Analysis ====")

# -----
# 1. Floor Range vs Lease Range
# -----
# Title font settings
title_font = dict(family="Arial", size=24, color="black", weight="bold")
subtitle_font = dict(fontStyle="italic", family="Arial", size=14, color="black")

# Set color palette
mako_palette = sns.color_palette("mako", as_cmap=True)
mako_colors = [mpl.colors.rgb2hex(c) for c in mako_palette(np.linspace(0, 1, 10))]

# Generate the data for Floor Level vs Price
df_cleaned = df.copy() # Create a copy to avoid modifying the original dataframe

# Calculate floor_avg if not present
if 'floor_avg' not in df_cleaned.columns and 'storey_range' in df_cleaned.columns:
    df_cleaned['floor_avg'] = df_cleaned['storey_range'].str.extract(r'(\d+)').astype(float)

df_cleaned['floor_range'] = pd.cut(
    df_cleaned['floor_avg'],
    bins=[0, 5, 10, 15, 20, 25, 50],
    labels=['1-5', '6-10', '11-15', '16-20', '21-25', '26+']
)

floor_price_df = df_cleaned.groupby('floor_range').agg({
    'resale_price': ['mean', 'count']
}).reset_index()
floor_price_df.columns = ['Floor Range', 'Average Price', 'Count']

# Generate the data for Price by Remaining Lease Years
if 'remaining_lease_years' not in df_cleaned.columns and 'lease_commence_date' in df_cleaned.columns:
    current_year = pd.to_datetime('today').year
    df_cleaned['remaining_lease_years'] = 99 - (current_year - df_cleaned['lease_commence_date'])
```

```

df_cleaned['lease_range'] = pd.cut(
    df_cleaned['remaining_lease_years'],
    bins=[0, 30, 50, 70, 90, 100],
    labels=['< 30 yrs', '30-50 yrs', '50-70 yrs', '70-90 yrs', '90+ yrs']
)

lease_price_df = df_cleaned.groupby('lease_range').agg({
    'resale_price': ['mean', 'count']
}).reset_index()
lease_price_df.columns = ['Lease Range', 'Average Price', 'Count']

# Create the subplot
fig_prices = make_subplots(
    rows=1, cols=2,
    subplot_titles=("Resale Price by Floor Range", "Resale Price by Remaining Lease Years"),
    horizontal_spacing=0.1 # Set the gap between subplots
)

# Add Floor Range bar chart to first subplot
fig_prices.add_trace(
    go.Bar(
        x=floor_price_df['Floor Range'],
        y=floor_price_df['Average Price'],
        marker_color=floor_price_df['Count'],
        marker_colorscale=mako_colors,
        text=floor_price_df['Count'],
        hovertemplate="Floor: %{x}<br>Price: ${y:.2f}<br>Count: %{text}<br>%{extra}</extra>",
        name="Floor Range"
    ),
    row=1, col=1
)

# Add Lease Range bar chart to second subplot
fig_prices.add_trace(
    go.Bar(
        x=lease_price_df['Lease Range'],
        y=lease_price_df['Average Price'],
        marker_color=lease_price_df['Count'],
        marker_colorscale=mako_colors,
        text=lease_price_df['Count'],
        hovertemplate="Lease: %{x}<br>Price: ${y:.2f}<br>Count: %{text}<br>%{extra}</extra>",
        name="Lease Range"
    ),
    row=1, col=2
)

# Update layout
fig_prices.update_layout(
    title={
        'text': 'Price Analysis by Floor Range and Remaining Lease',
        'font': {
            'family': "Arial, sans-serif",
            'size': 24,
            'weight': 'bold'
        },
        'y': 0.95,
        'x': 0.5,
        'xanchor': 'center'
    },
    title_font=title_font,
    height=800,
    width=1400,
    showlegend=False,
    template="plotly_white",
    coloraxis_showscale=True,
    coloraxis_colorbar=dict(title="Transaction Count")
)

# Update y-axes to have the same range for better comparison
max_price = max(
    floor_price_df['Average Price'].max(),
    lease_price_df['Average Price'].max()
)
fig_prices.update_yaxes(title_text="Average Price (SGD)", range=[0, max_price * 1.1])
fig_prices.update_yaxes(title_text="Floor Range", row=1, col=1)
fig_prices.update_yaxes(title_text="Remaining Lease", row=1, col=2)

fig_prices.show()

# -----
# 2. Town and Flat Type Analysis
# -----

# Generate data for Top Towns
town_price_df = df_cleaned.groupby('town').agg({
    'resale_price': ['mean', 'count']
}).reset_index()
town_price_df.columns = ['Town', 'Average Price', 'Count']
town_price_df = town_price_df.sort_values('Average Price', ascending=False).head(10)

# Generate data for Flat Types
flat_type_df = df_cleaned['flat_type'].value_counts().reset_index()
flat_type_df.columns = ['Flat Type', 'Count']
flat_type_df['Percentage'] = flat_type_df['Count'] / len(df_cleaned) * 100

# Create subplot
fig_towns_types = make_subplots(
    rows=1, cols=2,
    column_widths=[0.7, 0.3],
    subplot_titles=("Top 10 Towns by Average Resale Price", "Distribution of Flat Types"),
    specs=[[{"type": "bar"}, {"type": "pie"}]],
    horizontal_spacing=0.1
)

# Add Top Towns bar chart
fig_towns_types.add_trace(
    go.Bar(
        x=town_price_df['Town'],
        y=town_price_df['Average Price'],
        marker_color=town_price_df['Count'],
        marker_colorscale=mako_colors,
        text=town_price_df['Count'],
        hovertemplate="Town: %{x}<br>Price: ${y:.2f}<br>Count: %{text}<br>%{extra}</extra>",
        name="Towns"
    ),
    row=1, col=1
)

# Add Flat Types pie chart
fig_towns_types.add_trace(
    go.Pie(
        labels=flat_type_df['Flat Type'],
        values=flat_type_df['Count'],
        textinfo='label+percent',
        hovertemplate="%{label}<br>Count: %{value}<br>Percentage: %{percent}<br>%{extra}</extra>",
        marker_colors=mako_colors,
        name="Flat Types"
    ),
    row=1, col=2
)

# Update layout
fig_towns_types.update_layout(
    title={
        'text': 'Town Price Analysis and Flat Type Distribution',
        'font': {
            'family': "Arial, sans-serif",
            'size': 24,
            'weight': 'bold'
        },
        'y': 0.95,
        'x': 0.5,
        'xanchor': 'center'
    }
)

```

```

        'font': {
            'family': "Arial, sans-serif",
            'size': 24,
            'weight': 'bold'
        },
        'y': 0.95,
        'x': 0.5,
        'xanchor': 'center'
    },
    title_font=title_font,
    height=800,
    width=1400,
    showLegend=False,
    template="plotly_white",
    coloraxis_showscale=True,
    coloraxis_colorbar=dict(title="Transaction Count"),
)
# Set y-axis and x-axis titles for the bar chart
fig_towns_types.update_yaxes(title_text="Average Price (SGD)", row=1, col=1)
fig_towns_types.update_xaxes(title_text="Town", row=1, col=1)

fig_towns_types.show()

# -----
# 3. Most Popular Flat Type by Town
# -----
# Identify the most popular flat type
most_popular_flat_type = flat_type_df.sort_values('Count', ascending=False).iloc[0]
print(f"\nMost popular flat type: {most_popular_flat_type['Flat Type']} with {most_popular_flat_type['Count']} transactions ({most_popular_flat_type['Percentage']:.1f}% of total)")

# Top towns for the most popular flat type
popular_type_by_town = df_cleaned[df_cleaned['flat_type'] == most_popular_flat_type['Flat Type']].groupby('town').size().reset_index(name='transactions')
popular_type_by_town = popular_type_by_town.sort_values('transactions', ascending=False)

# Highest resale price by flat type and town
highest_price_by_flat_type_town = df_cleaned.groupby(['town', 'flat_type'])['resale_price'].max().reset_index()
highest_price_by_flat_type_town = highest_price_by_flat_type_town.sort_values('resale_price', ascending=False)

# Filter for top towns
top_towns = popular_type_by_town.head(10)[['town']].tolist()
highest_price_filtered = highest_price_by_flat_type_town[highest_price_by_flat_type_town['town'].isin(top_towns)]

# Create a combined subplot for these two new charts
fig_popular = make_subplots(
    rows=1, cols=2,
    subplot_titles=(
        "Top 10 Towns with Highest Demand for {most_popular_flat_type['Flat Type']} Flats",
        "Highest Resale Prices by Flat Type in Top Towns"
    ),
    horizontal_spacing=0.1
)

# Add bar chart for popular flat type by town
fig_popular.add_trace(
    go.Bar(
        x=popular_type_by_town.head(10)[['town']],
        y=popular_type_by_town.head(10)[['transactions']],
        marker_color=popular_type_by_town.head(10)[['transactions']],
        marker_colorscale=mako_colors,
        hovertemplate="Town: %{x}  
Transactions: %{y:.1f}</extra>",
        name=f"{most_popular_flat_type['Flat Type']} Transactions"
    ),
    row=1, col=1
)

# Create a grouped bar chart for highest resale prices by flat type in top towns
# Get the 5 most expensive flat types
top_flat_types = highest_price_filtered.groupby('flat_type')[['resale_price']].max().sort_values(ascending=False).head(5).index.tolist()
filtered_data = highest_price_filtered[highest_price_filtered['flat_type'].isin(top_flat_types)]

# Create scatterplot for each flat type, get the towns with the highest prices
for i, flat_type in enumerate(top_flat_types):
    flat_type_data = filtered_data[filtered_data['flat_type'] == flat_type].sort_values('resale_price', ascending=False)

    # Calculate bubble sizes based on price (min 15, max 35)
    min_price = filtered_data['resale_price'].min()
    max_price = filtered_data['resale_price'].max()
    price_range = max_price - min_price

    # Calculate sizes - use a simpler scale (all bubble sizes consistent per flat type)
    base_size = 15 + (i * 2) # Different size for each flat type

    fig_popular.add_trace(
        go.Scatter(
            x=flat_type_data[['town']],
            y=flat_type_data[['resale_price']],
            mode='markers',
            name=flat_type,
            marker=dict(
                color=mako_colors[i % len(mako_colors)],
                size=base_size,
                opacity=0.85,
                line=dict(width=1, color='white'),
                symbol='circle' # Use consistent circle shape for all
            ),
            hovertemplate="Town: %{x}  
Flat Type: " + flat_type + "  
Max Price: ${%{y:.0f}}</extra>",
        ),
        row=1, col=2
    )

# Update the layout
fig_popular.update_layout(
    title={
        'text': 'Analysis of Popular Flat Types and Highest Prices by Town',
        'font': {
            'family': "Arial, sans-serif",
            'size': 24,
            'weight': 'bold'
        },
        'y': 0.95,
        'x': 0.5,
        'xanchor': 'center'
    },
    height=800,
    width=1400,
    legend=dict(
        orientation="h",
        yanchor="bottom",
        y=-0.3,
        xanchor="center",
        x=0.5,
        font=dict(size=12),
        itemsizing="constant",
        bgcolor="rgba(255, 255, 255, 0.9)",
    ),
    plot_bgcolor='rgba(240, 240, 250, 0.8)' # Light blue-gray background
)

# Update axes
fig_popular.update_xaxes(title_text="Town", row=1, col=1, tickangle=-45)
fig_popular.update_yaxes(title_text="Number of Transactions", row=1, col=1)

fig_popular.update_xaxes(

```

```

        title_text="Town",
        row=1,
        col=2,
        tickangle=-45,
        gridcolor='rgba(255, 255, 255, 0.8)',
        showgrid=True
    )
    fig_popular.update_yaxes(
        title_text="Highest Resale Price (SGD)",
        row=1,
        col=2,
        tickprefix="",
        tickformat=",",
        gridcolor='rgba(255, 255, 255, 0.8)',
        showgrid=True
    )
    fig_popular.show()

# -----
# 4. Geographical Distribution of HDB Resale Prices
# -----
# Prepare town-level data for mapping
town_geo_data = df_cleaned.groupby('town').agg({
    'resale_price': ['mean', 'median', 'count', 'std'],
    'price_per_sqm': ['mean'] if 'price_per_sqm' in df_cleaned.columns else []
}).reset_index()

# Flatten multi-level columns
town_geo_data.columns = ['_'.join(col).strip() if col[1] else col[0] for col in town_geo_data.columns]

# Singapore town coordinates
town_coordinates = {
    'ANG MO KIO': [1.3691, 103.8454],
    'BEDOK': [1.3236, 103.9273],
    'BISHAN': [1.3526, 103.8352],
    'BUKIT BATOK': [1.3590, 103.7637],
    'BUKIT MERAH': [1.2819, 103.8239],
    'BUKIT PANJANG': [1.3774, 103.7719],
    'BUKIT TIMAH': [1.3294, 103.8021],
    'CENTRAL AREA': [1.2937, 103.8405],
    'CHOA CHU KANG': [1.3840, 103.7470],
    'CLEMENTI': [1.3162, 103.7649],
    'GEYLANG': [1.3201, 103.8918],
    'HOUGANG': [1.3612, 103.8863],
    'JURONG EAST': [1.3331, 103.7415],
    'JURONG WEST': [1.3404, 103.7090],
    'KALLANG/WHAMPOA': [1.3100, 103.8651],
    'MARINE PARADE': [1.3020, 103.9065],
    'PASIR RIS': [1.3721, 103.9474],
    'PUNGOL': [1.3984, 103.9072],
    'QUEENSTOWN': [1.2942, 103.7661],
    'SEMBAWANG': [1.4491, 103.8185],
    'SENGKANG': [1.3868, 103.8914],
    'SERANGOON': [1.3554, 103.8679],
    'TAMPINES': [1.3496, 103.9568],
    'TOA PAYOH': [1.3345, 103.8563],
    'WOODLANDS': [1.4382, 103.7890],
    'YISHUN': [1.4304, 103.8354]
}

# Create a base map centered on Singapore
map_center = [1.3521, 103.8198] # Singapore center coordinates
singapore_map = folium.Map(
    location=map_center,
    zoom_start=12,
    tiles='CartoDB positron', # Light theme map
    width='100%',
    height='100%'
)

# Create color bins
# Extract 5 colors to create a color ramp
mako_5_colors = [mpl.colors.rgb2hex(c) for c in mako_palette(np.linspace(0, 1, 5))]

# Function to determine color
def get_color(price):
    if price >= 750000:
        return mako_5_colors[4] # Darkest blue
    elif price >= 650000:
        return mako_5_colors[3]
    elif price >= 550000:
        return mako_5_colors[2]
    elif price >= 450000:
        return mako_5_colors[1]
    else:
        return mako_5_colors[0] # Lightest blue

# Calculate the price range for better sizing of markers
price_min = town_geo_data['resale_price_mean'].min()
price_max = town_geo_data['resale_price_mean'].max()
price_range = price_max - price_min

# Add town markers with price information
for i, row in town_geo_data.iterrows():
    town_name = row['town']

    # Skip if don't have coordinates for this town
    if town_name not in town_coordinates:
        continue

    avg_price = row['resale_price_mean']
    median_price = row['resale_price_median']
    num_transactions = row['resale_price_count']

    # Calculate marker size based on transaction count relative to max transaction count
    marker_size = 5 + (num_transactions / town_geo_data['resale_price_count'].max()) * 15

    # Create popup text with town price information
    popup_text = f"""


#### {town_name}



---


Average Price: ${avg_price:.0f}  

Median Price: ${median_price:.0f}  

Transactions: {num_transactions}


"""

    # Add price per sqm if available
    if 'price_per_sqm_mean' in row:
        popup_text += f"Avg Price/SQM: ${row['price_per_sqm_mean']:.0f}  
"

    popup_text += "
```

```

        popup=folium.Popup(popup_text, max_width=300)
    ).add_to(singapore_map)

# Add a legend
legend_html = """
<div style="position: fixed;
bottom: 50px; left: 50px; width: 200px; height: 180px;
border:2px solid grey; z-index:9999; font-size:14px;
background-color: white; padding: 10px; border-radius: 5px;">
<div style="font-weight: bold; margin-bottom: 8px; font-family: Arial;">Average Price (SGD)</div>
...
</div>

# Add each price range
price_ranges = [
['< $450,000', mako_5_colors[0]],
 ['$450,000 - $550,000', mako_5_colors[1]],
 ['$550,000 - $650,000', mako_5_colors[2]],
 ['$650,000 - $750,000', mako_5_colors[3]],
 ['$> $750,000', mako_5_colors[4]]
]

for price_range, color in price_ranges:
    legend_html += f"""
<div style="display: flex; align-items: center; margin-bottom: 8px; font-family: Arial;">
    <div style="background-color: {color}; width: 20px; height: 20px; margin-right: 8px; border: 1px solid #333;"></div>
    <div style="white-space: nowrap; font-family: Arial; font-weight: bold; margin-right: 8px;">{price_range}</div>
</div>
"""

legend_html += """
</div>
...
</div>

# Add the legend to the map
singapore_map.get_root().html.add_child(folium.Element(legend_html))

# Add title with matching font style to your other plots
title_html = """
<div style="position: fixed;
top: 10px; left: 50%; transform: translateX(-50%);
background-color: white; border:2px solid grey;
z-index:9999; padding: 10px; border-radius: 5px;
font-family: Arial; font-size: 20px; font-weight: bold;">
Geographical Distribution of HDB Resale Prices
</div>
...
singapore_map.get_root().html.add_child(folium.Element(title_html))

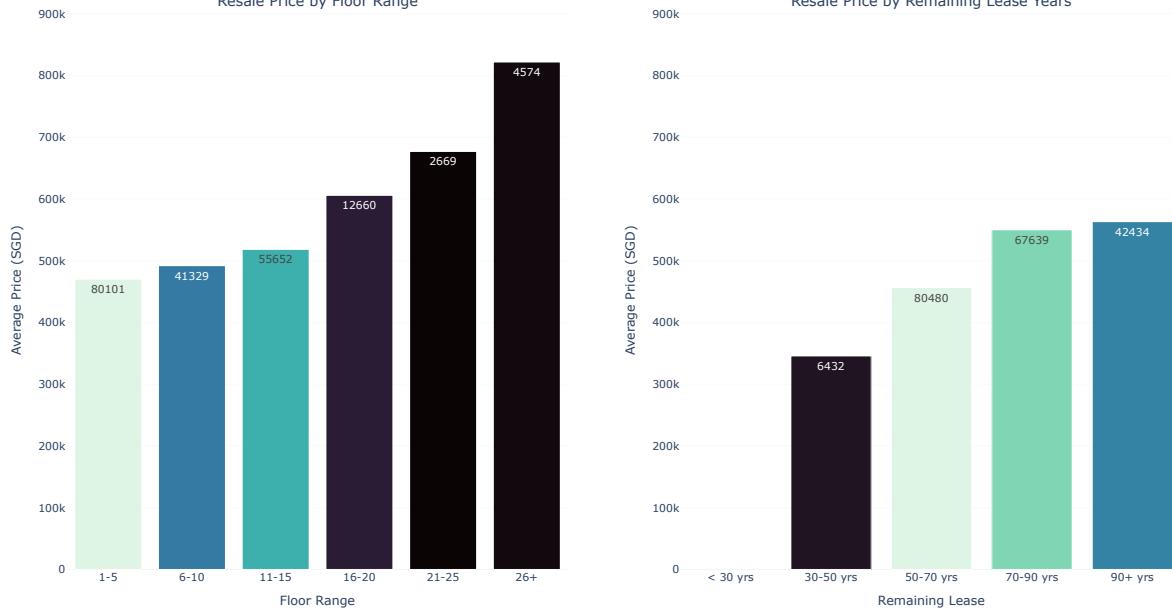
# Add a note for the circle sizes
size_note_html = """
<div style="position: fixed;
bottom: 50px; right: 50px; width: 200px;
border:2px solid grey; z-index:9999; font-size:14px;
background-color: white; padding: 10px; border-radius: 5px; font-family: Arial;">
<div style="font-weight: bold; margin-bottom: 5px;">Circle Size</div>
<div style="margin-bottom: 5px;">Proportional to number of transactions in each town</div>
</div>
...
singapore_map.get_root().html.add_child(folium.Element(size_note_html))

# Display the map in the notebook
singapore_map
"""

```

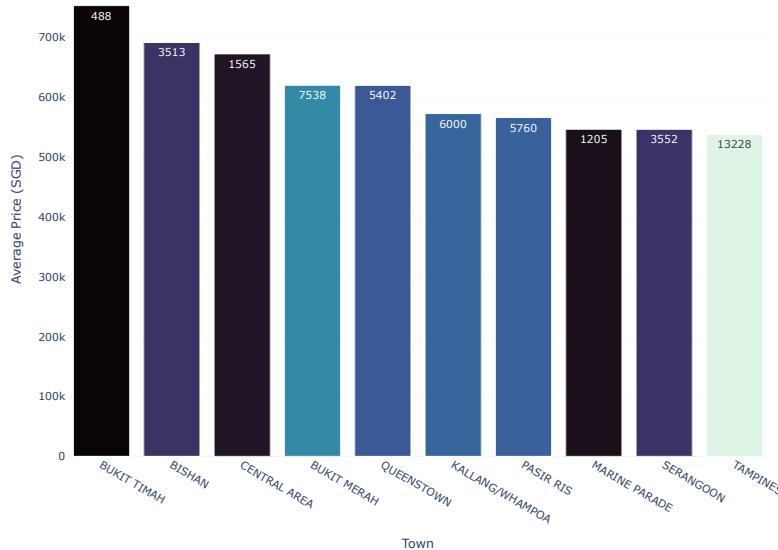
==== Advanced Pricing Analysis ===

Price Analysis by Floor Range and Remaining Lease

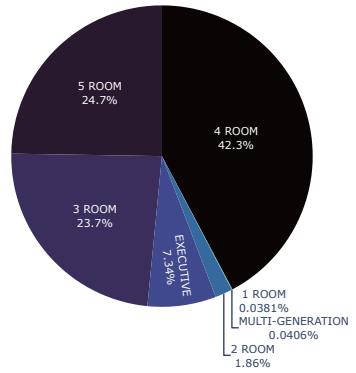


Town Price Analysis and Flat Type Distribution

Top 10 Towns by Average Resale Price

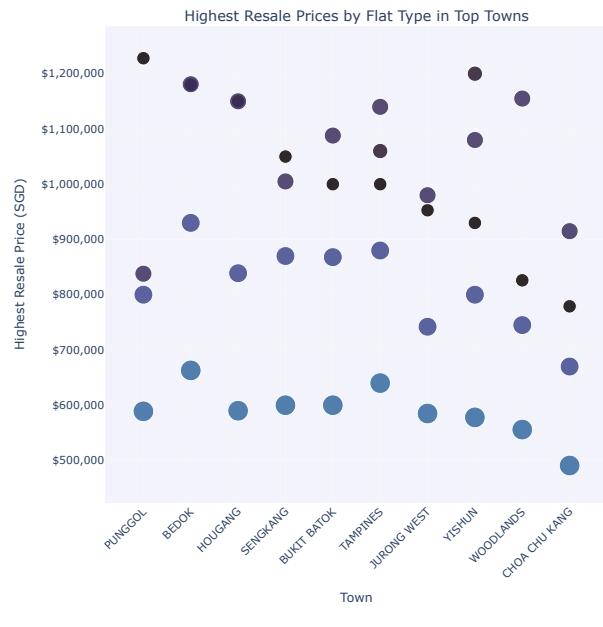
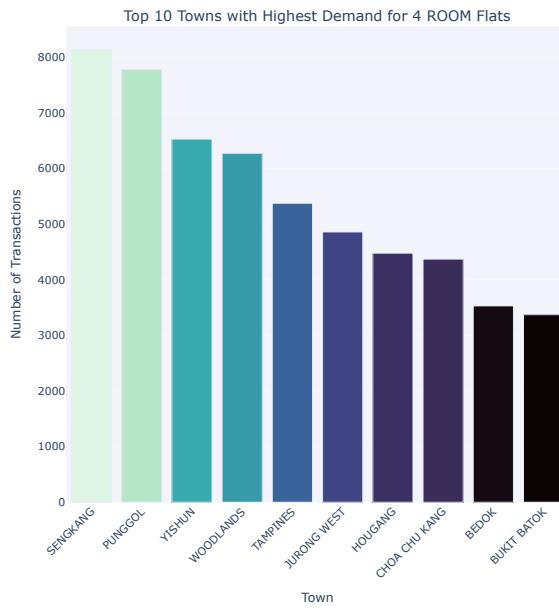


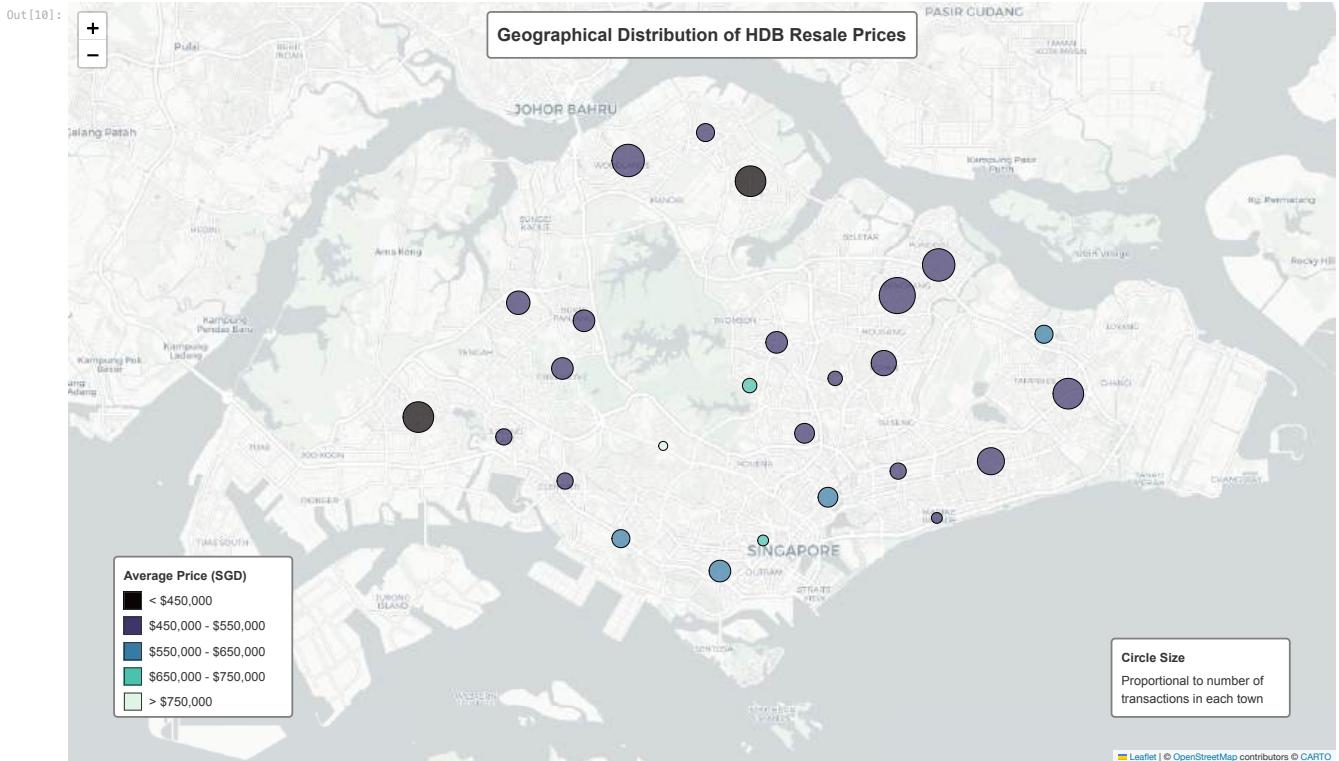
Distribution of Flat Types



Most popular flat type: 4 ROOM with 83,250 transactions (42.3% of total)

Analysis of Popular Flat Types and Highest Prices by Town





In [11]:

```
# 8: Advanced Regression Analysis
# Run Cell 1 first to load data and setup
print("\n===== ADVANCED REGRESSION ANALYSIS =====")

# -----
# 1. Importance features of Random Forest
# -----
# Set color palette
mako_palette = sns.color_palette("mako")
mako_colors = [f'rgb({int(rgb[0]*255)},{int(rgb[1]*255)},{int(rgb[2]*255)})' for rgb in mako_palette]
red_color = 'darkred'

# Prepare data for regression
print("\nPreparing data for regression analysis...")

# Select important features
regression_features = []
if 'flat_type' in df.columns: regression_features.append('flat_type')
if 'floor_area_sqm' in df.columns: regression_features.append('floor_area_sqm')
if 'storey_range' in df.columns: regression_features.append('storey_range')
if 'remaining_lease_years' in df.columns: regression_features.append('remaining_lease_years')
if 'town' in df.columns: regression_features.append('town')
if 'flat_model' in df.columns: regression_features.append('flat_model')

if len(regression_features) >= 2 and 'resale_price' in df.columns:
    # Create a copy for regression
    reg_df = df[regression_features + ['resale_price']].copy()
    reg_df = reg_df.dropna()

    # Encode categorical variables
    cat_cols = reg_df.select_dtypes(include=['object']).columns.tolist()
    num_cols = [col for col in regression_features if col not in cat_cols]

    # One-hot encode categorical columns
    reg_df_encoded = pd.get_dummies(reg_df, columns=cat_cols, drop_first=True)

    # Split features and target
    X = reg_df_encoded.drop(['resale_price'], axis=1)
    y = reg_df_encoded['resale_price']

    # Split into train and test sets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

    print(f"Training set: {X_train.shape[0]} samples")
    print(f"Test set: {X_test.shape[0]} samples")
    print(f"Total features after encoding: {X_train.shape[1]}")

    # Train different regression models
    models = {
        "Linear Regression": LinearRegression(),
        "Ridge Regression": Ridge(alpha=1.0),
        "Lasso Regression": Lasso(alpha=0.1),
        "Random Forest": RandomForestRegressor(n_estimators=100, max_depth=10, random_state=42)
    }

    # Train and evaluate each model
    results_table = PrettyTable()
    results_table.field_names = ["Model", "R² (Training)", "R² (Test)", "RMSE (Test)", "CV Score (5-fold)"]
    results_table.align = "l"

    for name, model in models.items():
        print(f"\nTraining {name}...")
        model.fit(X_train, y_train)

        # Calculate metrics
        train_r2 = model.score(X_train, y_train)
        test_r2 = model.score(X_test, y_test)

        # Make predictions
        y_pred = model.predict(X_test)
        rmse = np.sqrt(mean_squared_error(y_test, y_pred))

        # Cross-validation
        cv_scores = cross_val_score(model, X, y, cv=5, scoring='r2')
        cv_score_mean = cv_scores.mean()

        # Add to results table
        results_table.add_row([
            name,
            f"{train_r2:.4f}",
            f"{test_r2:.4f}",
            f"${rmse:.2f}",
            f'{cv_score_mean:.4f}'
        ])
    }
else:
    print("Not enough features selected for regression analysis. Please select at least 2 features from the list above and run again.")
```

```

        f"\n{cv_score_mean:.4f}""
    ])

print("\nRegression Model Comparison:")
print(results_table)

# Feature importance for Random Forest
if "Random Forest" in models:
    rf_model = models["Random Forest"]

    # Get feature importance
    feature_importance = pd.DataFrame({
        'Feature': X.columns,
        'Importance': rf_model.feature_importances_
    })
    feature_importance = feature_importance.sort_values('Importance', ascending=False).head(20)

    # Visualize feature importance
    fig = px.bar(
        feature_importance,
        x="Importance",
        y="Feature",
        orientation='h',
        title=<>Top 20 Features by Importance (Random Forest)</b>',
        labels={'Importance': 'Feature Importance', 'Feature': 'Feature'},
        color="Importance",
        color_continuous_scale=mako_colors, # Using Mako color palette
        template="plotly_white"
    )

    # Define fonts for consistency
    title_font = dict(size=24, family="Arial, sans-serif", weight="bold")
    axis_font = dict(size=14, family="Arial, sans-serif")

    fig.update_layout(
        title_font=title_font,
        xaxis_title_font=axis_font,
        yaxis_title_font=axis_font,
        yaxis=({"categoryorder": 'total ascending'}),
        height=800,
        width=1400,
        coloraxis_colorbar=dict(title="Importance"),
        title_x=0.5
    )
    fig.show()

# -----
# 2. Top 10 features with Importance Scores
# -----
# Print top 10 features with importance scores
importance_table = PrettyTable()
importance_table.field_names = ["Feature", "Importance"]
importance_table.align = "l"

for _, row in feature_importance.head(10).iterrows():
    importance_table.add_row([row['Feature'], f"\n{row['Importance']:.4f}"])

print("\nTop 10 Most Important Features:")
print(importance_table)

# Find top numeric features
top_numeric_features = []
for feature in feature_importance['Feature'].tolist():
    if feature in num_cols:
        top_numeric_features.append(feature)
    if len(top_numeric_features) >= 3: # Get top 3 numeric features
        break

if top_numeric_features:
    if len(top_numeric_features) >= 2:
        # Calculate layout dimensions
        n_features = len(top_numeric_features)
        n_cols = 2
        n_rows = (n_features + 1) // 2

        # Create subplot grid
        grid_fig = make_subplots(
            rows=n_rows,
            cols=n_cols,
            subplot_titles=[f"Effect of {feature} on Price" for feature in top_numeric_features],
            horizontal_spacing=0.15,
            vertical_spacing=0.2
        )

        for i in range(len(grid_fig.layout.annotations)):
            grid_fig.layout.annotations[i].font.size = 14

        for i, feature in enumerate(top_numeric_features):
            # Calculate row and column position
            row_idx = i // 2 + 1
            col_idx = i % 2 + 1

            # Sample for better performance
            sample_size = min(5000, len(reg_df))
            sample_df = reg_df.sample(sample_size, random_state=42)

            # Add scatter plot
            grid_fig.add_trace(
                go.Scatter(
                    x=sample_df[feature],
                    y=sample_df['resale_price'],
                    mode='markers',
                    marker=dict(
                        color=mako_colors[3],
                        size=5,
                        opacity=0.5
                    ),
                    name=feature,
                    showlegend=False
                ),
                row=row_idx, col=col_idx
            )

            # Add regression line
            feature_array = sample_df[feature].values.reshape(-1, 1)
            line_model = LinearRegression()
            line_model.fit(feature_array, sample_df['resale_price'])

            x_range = np.linspace(
                sample_df[feature].min(),
                sample_df[feature].max(),
                100
            ).reshape(-1, 1)
            y_pred = line_model.predict(x_range)

            grid_fig.add_trace(
                go.Scatter(
                    x=x_range.flatten(),
                    y=y_pred,
                    mode="lines",
                    line=dict(color=red_color, width=2),
                    showlegend=False
                ),
                row=row_idx, col=col_idx
            )
    
```

```

        name=f"feature_ Trend",
        ),
        row=row_idx, col=col_idx
    )

    # Update axes
    grid_fig.update_xaxes(
        title_text=feature,
        row=row_idx,
        col=col_idx,
        title_font=axis_font
    )

    grid_fig.update_yaxes(
        title_text="Resale Price" if col_idx == 1 else "",
        row=row_idx,
        col=col_idx,
        title_font=axis_font,
        tickprefix="$",
        tickformat=",."
    )

    # Update layout
    grid_fig.update_layout(
        title_text="<b>Partial Dependence of Top Numeric Features</b>",
        title_font=title_font,
        height=800 * n_rows,
        width=1400,
        template="plotly_white",
        title_x=0.5
    )
)

grid_fig.show()

else:
    print("Insufficient features for regression analysis. Need at least 2 features and price data.")
===== ADVANCED REGRESSION ANALYSIS =====

```

Preparing data for regression analysis...
Training set: 157588 samples
Test set: 39397 samples
Total features after encoding: 69

Training Linear Regression...

Training Ridge Regression...

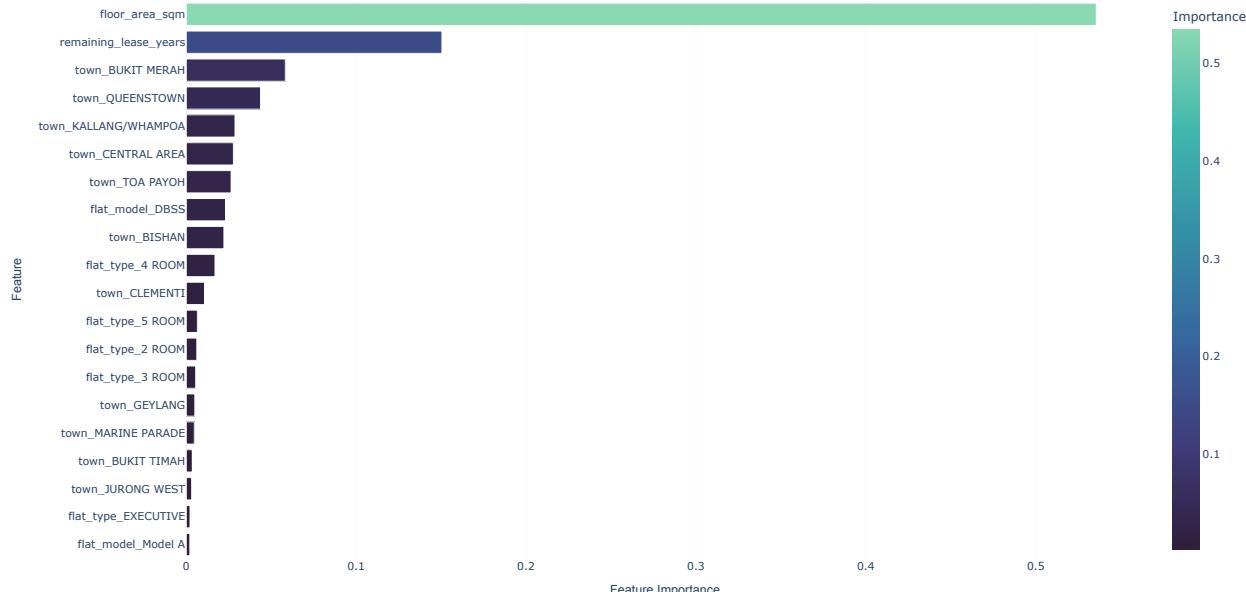
Training Lasso Regression...

Training Random Forest...

Regression Model Comparison:

Model	R ² (Training)	R ² (Test)	RMSE (Test)	CV Score (5-fold)
Linear Regression	0.7106	0.7067	\$96,233.86	0.5467
Ridge Regression	0.7106	0.7068	\$96,232.06	0.5468
Lasso Regression	0.7105	0.7067	\$96,239.75	0.5466
Random Forest	0.7434	0.7355	\$91,393.93	0.5845

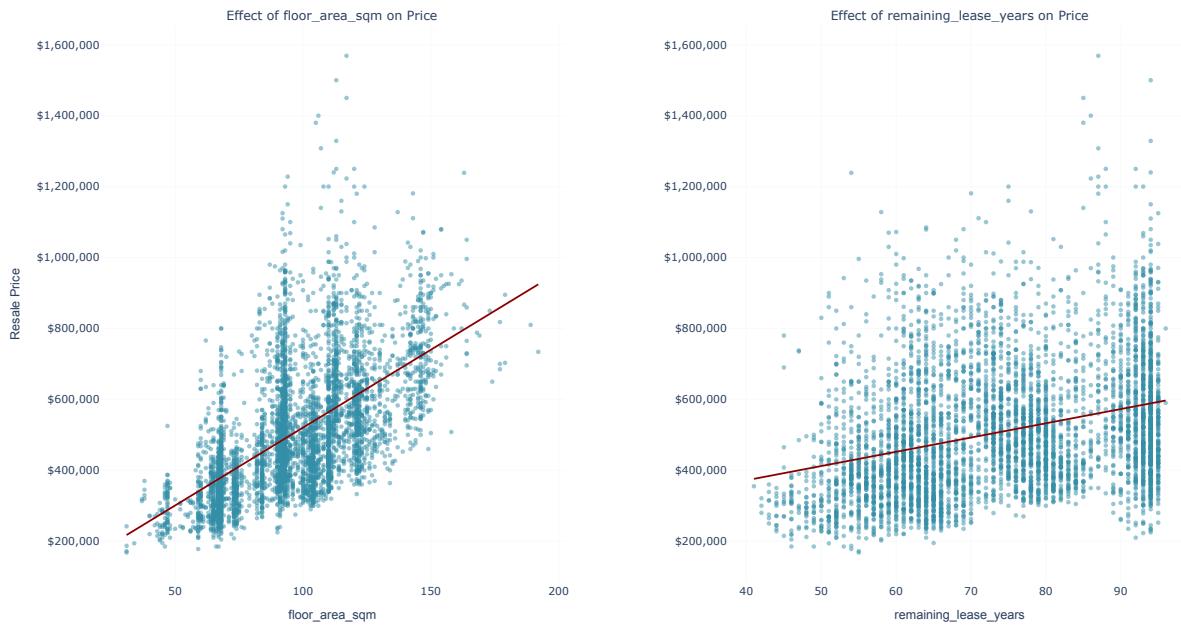
Top 20 Features by Importance (Random Forest)



Top 10 Most Important Features:

Feature	Importance
floor_area_sqm	0.5355
remaining_lease_years	0.1504
town_BUKIT MERAH	0.0583
town_QUEENSTOWN	0.0437
town_KALLANG/WHAMPOA	0.0287
town_CENTRAL AREA	0.0277
town_TOA PAYOH	0.0264
flat_model_DBSS	0.0231
town_BISHAN	0.0221
flat_type_4 ROOM	0.0169

Partial Dependence of Top Numeric Features



```
In [17]: # 9: Market Segmentation Analysis
# Run 1 first to load data and setup
print("\n===== MARKET SEGMENTATION ANALYSIS =====")

# -----
# 1. Treemap of Towns by Price
# -----

# Set color palette
mako_palette = sns.color_palette("mako", n_colors=6)
exact_colors = [to_hex(color) for color in mako_palette]

# Create color scale for heatmaps
mako_colorscale = [[1/5, color] for i, color in enumerate(exact_colors)]

# Function to create the treemap figure
def create_town_segment_treemap(df_tree, segment_colors=None):
    base_color = "#0c343d"

    dark_colors = {
        'Budget': exact_colors[0], # Darkest color
        'Mid-Range': exact_colors[1], # Get the second darkest
        'Premium': exact_colors[2], # Mid color
        'Luxury': exact_colors[3] # Light color
    }

    # Create the treemap figure
    fig_tree = px.treemap(
        df_tree,
        path=['town', 'price_segment'],
        values='count',
        color='price_segment',
        color_discrete_map=dark_colors,
        custom_data=['count']
    )

    # Update the treemap consistent colors
    fig_tree.update_traces(
        marker=dict(
            colors=[base_color if 'price_segment' not in node_path else
                   dark_colors.get(node_path.split('/')[-1], base_color)
                   for node_path in fig_tree.data[0].ids],
            line=dict(width=1, color="#0c343d") # Subtle border color
        ),
        root_color=base_color,
        textfont=dict(size=14, color="white"),
        hovertemplate='<br>%{label}<br>Count: %{value}<br>%{extra}'
    )

    return fig_tree

# Common Plotly font settings
title_font = dict(family="Arial", size=24, color="black", weight="bold")
subtitle_font = dict(family="Arial", size=14, color="black")
axis_font = dict(family="Arial", size=12, color="black")
value_font = dict(family="Arial", size=18, color="black", weight="bold")

# Create price segments
if 'resale_price' in df.columns:
    # Define price segments - match original 4 segments from Cell 9
    df['price_segment'] = pd.qcut(
        df['resale_price'],
        q=4,
        labels=['Budget', 'Mid-Range', 'Premium', 'Luxury']
    )

    # Count by segment
    segment_counts = df['price_segment'].value_counts().reset_index()
    segment_counts.columns = ['Segment', 'Count']

    # Calculate price ranges for each segment
    segment_ranges = df.groupby('price_segment')['resale_price'].agg(['min', 'max', 'mean', 'median']).reset_index()

    # Create a price segment table with enhanced statistics
    segment_table = PrettyTable()
    segment_table.field_names = ["Segment", "Count", "Price Range", "Mean Price", "Median Price", "% of Market"]
    segment_table.align = "l"

    for i, row in segment_ranges.iterrows():
        segment = row['price_segment']
        count = segment_counts[segment_counts['Segment'] == segment]['Count'].values[0]
        percentage = count / len(df) * 100
        segment_table.add_row([
            segment,
            f'{count},',
            f'{percentage}%'
```

```

    f"${row['min']):.0f} to ${row['max']):.0f}",
    f"${row['mean']):.0f}",
    f"${row['median']):.0f}",
    f"${percentage:.1f}%"}
  ])
print("Market Segmentation by Price:")
print(segment_table)

# Get top towns by count
top_towns = df['town'].value_counts().head(10).index.tolist()
df_top_towns = df[df['town'].isin(top_towns)]

# Get town segment data
town_segment_counts = pd.crosstab(df_top_towns['town'], df_top_towns['price_segment'])
segment_by_town = pd.crosstab(
  df_top_towns['town'],
  df_top_towns['price_segment'],
  normalize='index'
) * 100

# Prepare TreeMap data
df_tree = df_top_towns.groupby(['town', 'price_segment']).size().reset_index(name='count')

# Map segments to colors from the Mako palette for pie charts and other visualizations
segment_colors = {
  'Budget': exact_colors[0], # Darkest color
  'Mid-Range': exact_colors[2], # Middle blue
  'Premium': exact_colors[4], # Light blue
  'Luxury': exact_colors[5] # Lightest color
}

# -----
# 2. Pie Chart and Town Heatmap
# -----
# Create subplot
fig1 = make_subplots(
  rows=1, cols=2,
  column_widths=[0.4, 0.6],
  specs=[[{"type": "pie"}, {"type": "heatmap"}]],
  subplot_titles=("HDB Market Distribution by Price Segment",
                 "Price Segment Distribution by Town (Top 10)"),
  horizontal_spacing=0.2
)

# Add pie chart to subplot
fig1.add_trace(
  go.Pie(
    labels=segment_counts['Segment'],
    values=segment_counts['Count'],
    hole=.3,
    textinfo='label+percent',
    textfont=dict(size=14),
    marker_colors=[segment_colors[segment] for segment in segment_counts['Segment']],
    hovertemplate='<b>{label}</b>%  
Count: %{value:,}  
Percentage: %{percent}<br><extra>'
  ),
  row=1, col=1
)

# Create town heatmap trace
town_heatmap = go.Heatmap(
  z=segment_by_town.values,
  x=segment_by_town.columns,
  y=segment_by_town.index,
  colorscale=mako_colorscale,
  text=[[f'{val:.1f}' for val in row] for row in segment_by_town.values],
  texttemplate="%{text}",
  textfont=dict(size=14),
  hovertemplate='<b>Town</b>: %{y}<br><b>Segment</b>: %{x}<br><b>Percentage</b>: %{z:.1f}%<br><extra>' 
)
fig1.add_trace(town_heatmap, row=1, col=2)

# Add town heatmap to subplot
fig1.add_trace(town_heatmap, row=1, col=2)

# Update layout
fig1.update_layout(
  title={
    'text': 'Market Distribution and Town Analysis',
    'font': {
      'family': "Arial",
      'size': 24,
      'weight': "bold"
    },
    'x': 0.5,
    'xanchor': 'center'
  },
  title_font=title_font,
  height=800,
  width=1400,
  template="plotly_white",
  showLegend=True,
  legend=dict(orientation="h", y=0.1, xanchor="right", x=0.3),
  margin=dict(t=100, b=50, l=20, r=20)
)

# Update heatmap axes
fig1.update_xaxes(title_text="Price Segment", title_font=axis_font, row=1, col=2)
fig1.update_yaxes(title_text="Town", title_font=axis_font, row=1, col=2)

# Update subplot titles
fig1.update_annotations(font=subtitle_font)

fig1.show()

# Only proceed if flat_type is in df columns
if 'flat_type' in df.columns:
  # Get flat type segment data
  flat_type_segment_counts = pd.crosstab(df['flat_type'], df['price_segment'])

  # Get percentage distribution
  segment_by_type = pd.crosstab(
    df['flat_type'],
    df['price_segment'],
    normalize='index'
  ) * 100

  # Order flat types in a logical sequence
  flat_type_order = ['1 ROOM', '2 ROOM', '3 ROOM', '4 ROOM', '5 ROOM', 'EXECUTIVE', 'MULTI-GENERATION']
  # Filter to only include flat types in our data
  flat_type_order = [ft for ft in flat_type_order if ft in segment_by_type.index]
  # Reindex with custom order
  segment_by_type = segment_by_type.reindex(flat_type_order)
  flat_type_segment_counts = flat_type_segment_counts.reindex(flat_type_order)

# -----
# 3. Flat Type and Floor Range Analysis by Price
# -----
if 'storey_range' in df.columns:
  # Floor range distribution
  # Extract floor number from storey_range
  df['floor_number'] = df['storey_range'].str.extract(r'^(\d+)').astype(int)

  # Create floor categories
  df['floor_category'] = pd.cut(

```

```

        df['floor_number'],
        bins=[0, 5, 10, 15, 20, 50],
        labels=['1-5', '6-10', '11-15', '16-20', '21+']
    )

    # Cross tabulation of floor category by price segment
    floor_segment_counts = pd.crosstab(df['floor_category'], df['price_segment'])
    floor_segment = pd.crosstab(
        df['floor_category'],
        df['price_segment'],
        normalize='index'
    ) * 100

    # Prepare data for grouped bar chart
    floor_data = floor_segment.reset_index().melt(
        id_vars='floor_category',
        var_name='price_segment',
        value_name='percentage'
    )

    # Create subplot
    fig3 = make_subplots(
        rows=1, cols=2,
        column_widths=[0.6, 0.4],
        specs=[[{"type": "bar"}, {"type": "bar"}]],
        subplot_titles=("Flat Type Distribution by Price Segment",
                        "Price Segment Distribution by Floor Range"),
        horizontal_spacing=0.1
    )

    # No text annotations on the bars
    flat_type_annotations = []

    # Add stacked bar chart for each segment
    for segment in segment_by_type.columns:
        fig3.add_trace(
            go.Bar(
                y=flat_type_segment_counts.index,
                x=flat_type_segment_counts[segment],
                name=segment,
                orientation='h',
                marker_color=segment_colors.get(segment),
                hovertemplate='<b>%{y}</b><br><b>Segment</b>: ' + segment + '<br><b>Count</b>: %{x:.0f}</extra>'
            ),
            row=1, col=1
        )

    # Add grouped bar chart for floor distribution with percentage values
    for segment in floor_segment.columns:
        segment_data = floor_data[floor_data['price_segment'] == segment]

        fig3.add_trace(
            go.Bar(
                x=segment_data['floor_category'],
                y=segment_data['percentage'],
                name=segment,
                marker_color=segment_colors.get(segment),
                text=[f'({p:.1f}%)' for p in segment_data['percentage']],
                textposition='outside',
                textfont=dict(size=14),
                hovertemplate='<b>%{x}</b><br><b>Segment</b>: ' + segment + '<br><b>Percentage</b>: %{y:.1f}%</extra>'
            ),
            row=1, col=2
        )

    # Update layout
    fig3.update_layout(
        title={
            'text': 'Flat Type and Floor Range Analysis by Price Segment',
            'font': {
                'size': 24,
                'family': 'Arial, sans-serif',
                'weight': 'bold'
            },
            'y': 0.95,
            'x': 0.5,
            'xanchor': 'center'
        },
        title_font=title_font,
        height=800,
        width=1400,
        template="plotly_white",
        showlegend=True,
        margin=dict(t=100, b=50, l=120, r=20),
        barmode='stack'
        # No text annotations on bars
    )

    # Update axes
    fig3.update_xaxes(
        title_text="Number of Transactions",
        title_font=axis_font,
        row=1,
        col=1,
        range=[0, flat_type_segment_counts.sum(axis=1).max() * 1.1] # Add padding
    )
    fig3.update_xaxes(
        title_text="Flat Type",
        title_font=axis_font,
        row=1,
        col=1
    )

    fig3.update_xaxes(
        title_text="Floor Range",
        title_font=axis_font,
        row=2,
        col=2,
        categoryorder='array',
        categoryarray=['1-5', '6-10', '11-15', '16-20', '21+']
    )
    fig3.update_xaxes(
        title_text="Percentage (%)",
        title_font=axis_font,
        row=2,
        col=2,
        ticksuffix="%"
    )

    # Update subplot titles
    fig3.update_annotations(font=subtitle_font)

fig3.show()

# -----
# 4. Market Segmentation - TreeMap
# -----

# Define price segments
def categorize_price(price):
    if price < df['resale_price'].quantile(0.25):
        return 'Low Price'
    elif price < df['resale_price'].quantile(0.75):
        return 'Mid Price'
    else:
        return 'High Price'

```

```

        return 'High Price'

# Create price segment column
df['price_segment'] = df['resale_price'].apply(categorize_price)

# Prepare data for treemap
df_tree = df.groupby(['town', 'price_segment']).size().reset_index(name='count')

# Define color mapping for segments
segment_colors = {
    'Low Price': mako_colors[1], # Darker blue for low price
    'Mid Price': mako_colors[5], # Medium blue for mid price
    'High Price': mako_colors[9] # Lightest blue for high price
}

# Create the treemap
fig_tree = px.treemap(
    df_tree,
    path=['town', 'price_segment'],
    values='count',
    color='price_segment',
    color_discrete_map=segment_colors,
    custom_data=['count']
)

# Update trace details for better readability
fig_tree.update_traces(
    marker=dict(
        line=dict(width=1, color="#00888B") # Subtle border color
    ),
    textfont=dict(size=14, color="white") # Use white for text
)

# Optional: Add hover template to show more information
fig_tree.update_traces(
    hovertemplate='<b>%{label}</b><br>Count: %{customdata[0]}<br>%{extra}</extra>'
)

fig_tree.update_layout(
    title={
        'text': 'Market Segmentation by Town and Price Range',
        'font': {
            'size': 24,
            'family': 'Arial, sans-serif',
            'weight': 'bold'
        },
        'x': 0.5,
        'xanchor': 'center'
    },
    height=800,
    width=1400
)
fig_tree.show()

# -----
# 5. Price Segment Distribution by Flat Type
# -----
# Create flat type heatmap as its own figure
flat_type_fig = go.Figure(
    go.Heatmap(
        z=segment_by_type.values,
        x=segment_by_type.columns,
        y=segment_by_type.index,
        colorscale=mako_colorscale,
        text=[[f'{val:.1f}%' for val in row] for row in segment_by_type.values],
        texttemplate="%{text}",
        textfont=dict(size=14),
        hovertemplate='<b>Flat Type</b>: %{y}<br><b>Segment</b>: %{x}<br><b>Percentage</b>: %{z:.1f}%<br>%{extra}</extra>'
    )
)

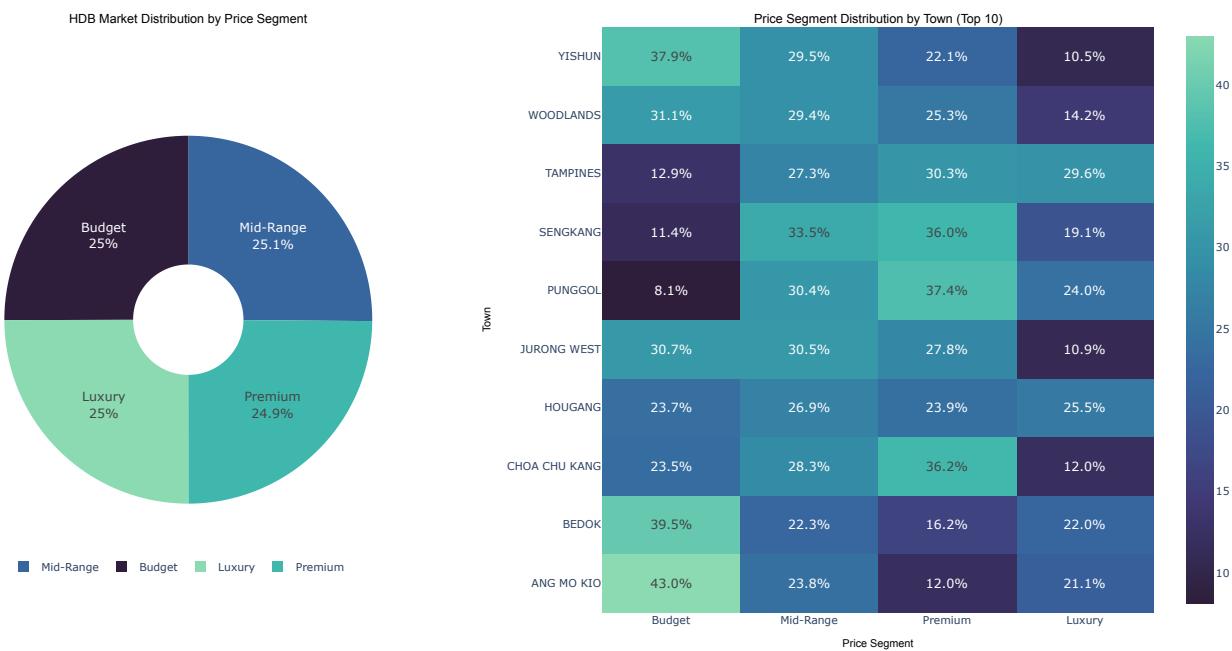
# Update flat type heatmap layout
flat_type_fig.update_layout(
    title={
        'text': 'Price Segment Distribution by Flat Type',
        'font': {
            'size': 24,
            'family': 'Arial, sans-serif',
            'weight': 'bold'
        },
        'y': 0.95,
        'x': 0.5,
        'xanchor': 'center'
    },
    title_font=title_font,
    height=800,
    width=1400,
    template="plotly_white",
    margin=dict(t=100, b=50, l=120, r=20),
    xaxis_title="Price Segment",
    yaxis_title="Flat Type",
    xaxis_title_font=axis_font,
    yaxis_title_font=axis_font
)
flat_type_fig.show()

else:
    print("No 'resale_price' column found. Market segmentation analysis skipped.")

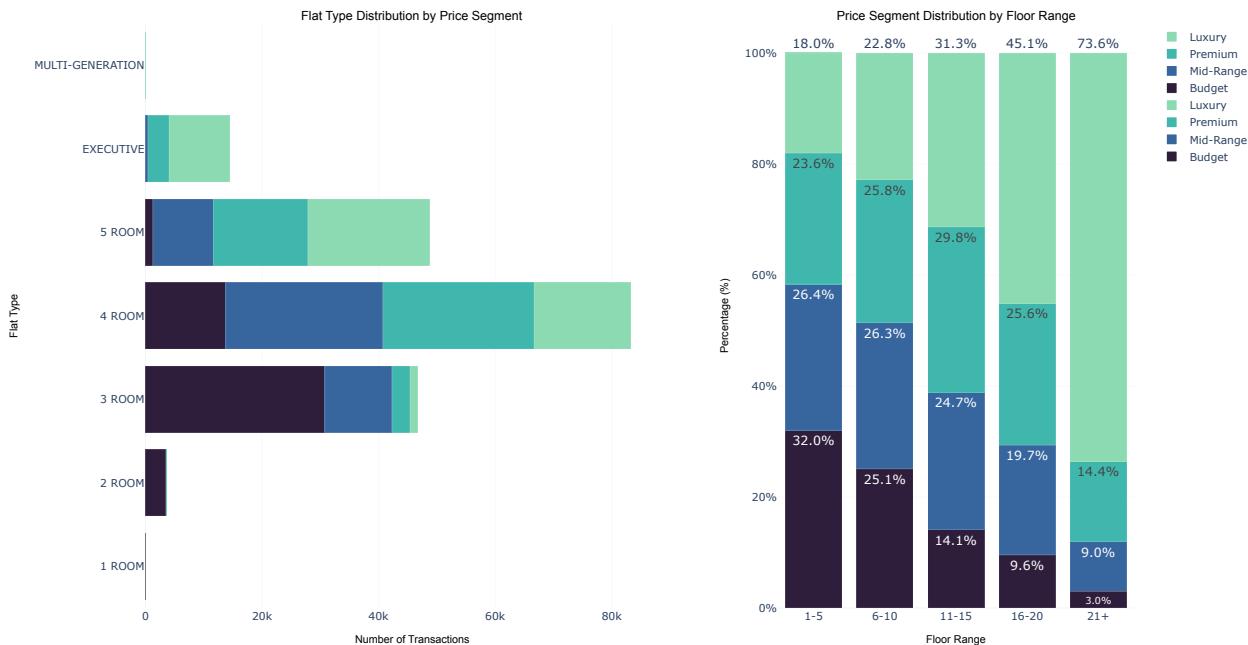
===== MARKET SEGMENTATION ANALYSIS =====
Market Segmentation by Price:
+-----+-----+-----+-----+-----+
| Segment | Count | Price Range | Mean Price | Median Price | % of Market |
+-----+-----+-----+-----+-----+
| Budget | 49,322 | $140,000 to $377,000 | $314,140 | $320,000 | 25.0% |
| Mid-Range | 49,487 | $377,400 to $478,000 | $426,945 | $426,888 | 25.1% |
| Premium | 48,955 | $478,500 to $608,000 | $538,319 | $535,000 | 24.9% |
| Luxury | 49,221 | $608,028 to $1,588,000 | $754,224 | $720,000 | 25.0% |
+-----+-----+-----+-----+-----+

```

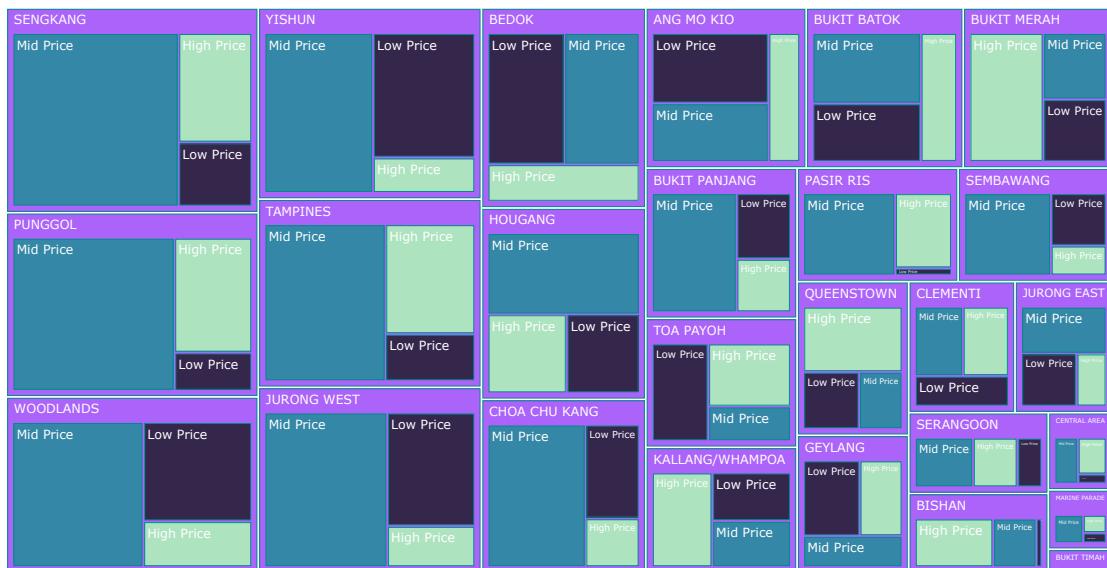
Market Distribution and Town Analysis



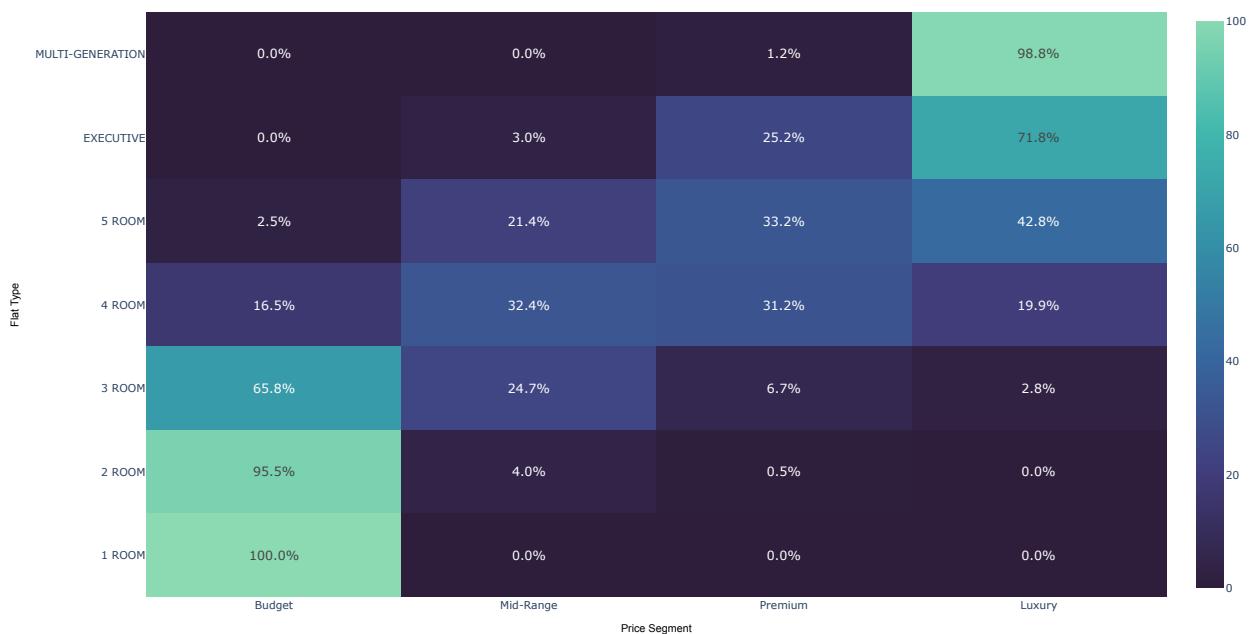
Flat Type and Floor Range Analysis by Price Segment



Market Segmentation by Town and Price Range



Price Segment Distribution by Flat Type



```
In [13]: # 10: Comparative Analysis by Year
# Run 1 first to load data and setup
print("\n===== COMPARATIVE ANALYSIS BY YEAR =====")

# -----
# 1. Year-over-Year Real Estate Analysis
# -----
# Set color palette
mako_palette = sns.color_palette("mako")
mako_colors = [f'rgb({int(rgb[0]*255)},{int(rgb[1]*255)},{int(rgb[2]*255)})' for rgb in mako_palette]

# Create a subset for line charts and bars to avoid too light colors
plotly_colors = mako_colors[1:] # Skip the lightest color for better visibility
# Also create a reversed version (darker to lighter)
plotly_colors_reversed = plotly_colors[::-1]

# Define chart dimensions
chart_width = 1400
chart_height = 600

# Define fonts for consistent styling
title_font = dict(family="Arial, sans-serif", size=22, color="#333333", weight="bold")
axis_font = dict(family="Arial, sans-serif", size=14, color="#555555")

if 'year' in df.columns:
    # If year column doesn't exist but month does, create it
    if 'year' not in df.columns and "month" in df.columns and pd.api.types.is_datetime64_any_dtype(df['month']):
        df['year'] = df['month'].dt.year
        print("Created 'year' column from 'month' column")

    # Calculate yearly price statistics
    yearly_stats = df.groupby('year').agg({
        'resale_price': ['mean', 'median', 'count', 'std', 'min', 'max'],
        'price_per_sqm': ['mean', 'median', 'std'] if 'price_per_sqm' in df.columns else []
    }).reset_index()

    # Flatten multi-level columns
    for col in yearly_stats.columns:
        if len(yearly_stats[col].shape) > 1:
            yearly_stats[col] = yearly_stats[col].apply(lambda x: x[0] if len(x) == 1 else x)
```

```

yearly_stats.columns = ['year'] + ['_'.join(col).strip() for col in yearly_stats.columns[1:]]

# Calculate year-over-year growth
yearly_stats['yoY_growth'] = yearly_stats['resale_price_mean'].pct_change() * 100

# Calculate volatility (coefficient of variation)
yearly_stats['price_volatility'] = (yearly_stats['resale_price_std'] / yearly_stats['resale_price_mean']) * 100

# Calculate transaction volume growth
yearly_stats['volume_growth'] = yearly_stats['resale_price_count'].pct_change() * 100

# Create a table of yearly statistics with enhanced metrics
yearly_table = PrettyTable()
yearly_table.field_names = ["Year", "Avg Price", "YoY Growth", "Transactions", "Vol Growth", "Price Range", "Volatility"]
yearly_table.align = "l" # Left-align all columns

for _, row in yearly_stats.iterrows():
    yoY_str = f'{row["yoY_growth"]:.2f}' if not pd.isna(row["yoY_growth"]) else "N/A"
    vol_growth_str = f'{row["volume_growth"]:.2f}' if not pd.isna(row["volume_growth"]) else "N/A"
    price_range = f"${row['resale_price_min']:.0f} - ${row['resale_price_max']:.0f}"

    yearly_table.add_row([
        int(row['year']),
        f"${row['resale_price_mean']:.2f}",
        yoY_str,
        f"{int(row['resale_price_count']):,}",
        vol_growth_str,
        price_range,
        f'{row["price_volatility"]:.2f}'
    ])

print("Year-over-Year Analysis:")
print(yearly_table)

# Add table for price per sqm if available
if 'price_per_sqm_mean' in yearly_stats.columns:
    sqm_table = PrettyTable()
    sqm_table.field_names = ["Year", "Avg Price/SQM", "Median Price/SQM", "Std Dev"]
    sqm_table.align = "l"

    for _, row in yearly_stats.iterrows():
        sqm_table.add_row([
            int(row['year']),
            f"${row['price_per_sqm_mean']:.2f}",
            f"${row['price_per_sqm_median']:.2f}",
            f"${row['price_per_sqm_std']:.2f}"
        ])

    print("\nPrice per Square Meter by Year:")
    print(sqm_table)

# Create a comprehensive comparative chart
fig = make_subplots(
    rows=1, cols=2,
    subplot_titles=[
        "Price Trends by Year",
        "Transaction Analysis by Year"
    ],
    specs=[
        [{"secondary_y": True}, {"secondary_y": True}]
    ],
    horizontal_spacing=0.15
)

# Create a color generator function to ensure consistent colors across charts
def get_color(index):
    if index < len(plotly_colors_reversed):
        return plotly_colors_reversed[index]
    else:
        # Cycle through colors if we run out
        return plotly_colors_reversed[index % len(plotly_colors_reversed)]

# 1. Add average price bar chart to first subplot
fig.add_trace(
    go.Bar(
        x=yearly_stats['year'],
        y=yearly_stats['resale_price_mean'],
        name='Average Price',
        marker_color=get_color(0),
        opacity=0.8,
        hovertemplate="Year: %{x}<br>Average Price: ${y:.2f}</extra></extra>",
        row=1, col=1,
        secondary_y=False
    )

    # Add median price line to first subplot
    fig.add_trace(
        go.Scatter(
            x=yearly_stats['year'],
            y=yearly_stats['resale_price_median'],
            mode='lines+markers',
            name='Median Price',
            line=dict(color=get_color(1), width=3, dash='dot'),
            marker=dict(size=8),
            hovertemplate="Year: %{x}<br>Median Price: ${y:.2f}</extra></extra>",
            row=1, col=1,
            secondary_y=False
        )

        # Add YoY growth line on secondary y-axis
        fig.add_trace(
            go.Scatter(
                x=yearly_stats['year'],
                y=yearly_stats['yoY_growth'],
                mode='lines+markers',
                name='YoY Growth',
                line=dict(color=get_color(2), width=3),
                marker=dict(size=8, symbol="diamond"),
                hovertemplate="Year: %{x}<br>YoY Growth: %{y:.2f}</extra></extra>",
                row=1, col=1,
                secondary_y=True
            )
        )

    # Add price per sqm if available
    if 'price_per_sqm_mean' in yearly_stats.columns:
        fig.add_trace(
            go.Scatter(
                x=yearly_stats['year'],
                y=yearly_stats['price_per_sqm_mean'],
                mode='lines+markers',
                name='Price per SQM',
                line=dict(color=get_color(3), width=3),
                marker=dict(size=8),
                hovertemplate="Year: %{x}<br>Price per SQM: ${y:.2f}</extra></extra>",
                row=1, col=1,
                secondary_y=False
            )
        )
    # 2. Add transaction volume to second subplot
    fig.add_trace(

```

```

        go.Bar(
            x=yearly_stats['year'],
            y=yearly_stats['resale_price_count'],
            name='Transaction Volume',
            marker_color=get_color(0),
            opacity=0.8,
            hovertemplate="Year: %(x)<br>Transactions: %{y:,}<!--extra--&gt;"'
        ),
        row=1, col=2,
        secondary_y=False
    )

    # Add volume growth as a line on secondary y-axis
    fig.add_trace(
        go.Scatter(
            x=yearly_stats['year'],
            y=yearly_stats['volume_growth'],
            mode='lines+markers',
            name='Volume Growth',
            line=dict(color=get_color(1), width=3),
            marker=dict(size=8, symbol="diamond"),
            hovertemplate="Year: %(x)&lt;br&gt;Volume Growth: %{y:.2f}<!--extra--&gt;"'
        ),
        row=1, col=2,
        secondary_y=True
    )

    # Add price volatility as a line
    fig.add_trace(
        go.Scatter(
            x=yearly_stats['year'],
            y=yearly_stats['price_volatility'],
            mode='lines+markers',
            name='Price Volatility',
            line=dict(color=get_color(2), width=3, dash='dot'),
            marker=dict(size=8),
            hovertemplate="Year: %(x)&lt;br&gt;Volatility: %{y:.2f}<!--extra--&gt;"'
        ),
        row=1, col=2,
        secondary_y=True
    )

    # Add zero reference line for growth charts
    fig.add_shape(
        type="line",
        x0=yearly_stats['year'].min(),
        y0=0,
        x1=yearly_stats['year'].max(),
        y1=0,
        line=dict(color="gray", width=1, dash="dash"),
        row=1, col=2,
        secondary_y=True
    )

    # Update layout with subplot titles
    fig.update_layout(
        title={
            'text': "Comprehensive Year-over-Year Real Estate Analysis",
            'font': {
                'size': 24,
                'family': 'Arial, sans-serif',
                'weight': 'bold'
            },
            'x': 0.5,
            'y': 0.99,
            'xanchor': 'center'
        },
        template="plotly_white",
        height=chart_height,
        width=chart_width,
        legend=dict(
            orientation="h",
            yanchor="top",
            y=0.2,
            xanchor="center",
            x=0.5
        ),
        hovermode="x unified",
        margin=dict(t=100, b=120)
    )

    # Update subplot titles to be bold and centered with adjusted position
    for i in range(len(fig['layout']['annotations'])):
        fig['layout'][i]['annotations'][i]['font'] = dict(family="Arial, sans-serif", size=16, color="#333333")
        fig['layout'][i]['annotations'][i][i]['xanchor'] = 'center'
        fig['layout'][i]['annotations'][i][i]['x'] = 0.25 if i == 0 else 0.75
        fig['layout'][i]['annotations'][i][i]['y'] = 1.05

    # Update y-axes for first subplot
    fig.update_yaxes(
        title_text="Price (SGD)",
        title_font=axis_font,
        secondary_y=False,
        row=1, col=1,
        tickprefix="$",
        tickformat=".",
        gridcolor='rgba(211,211,211,0.3)'
    )

    fig.update_yaxes(
        title_text="YoY Growth (%)",
        title_font=axis_font,
        secondary_y=True,
        row=1, col=1,
        ticksuffix="%",
        gridcolor='rgba(211,211,211,0.3)'
    )

    # Update y-axes for second subplot
    fig.update_yaxes(
        title_text="Number of Transactions",
        title_font=axis_font,
        secondary_y=False,
        row=1, col=2,
        gridcolor='rgba(211,211,211,0.3)'
    )

    fig.update_yaxes(
        title_text="Growth &amp; Volatility (%)",
        title_font=axis_font,
        secondary_y=True,
        row=1, col=2,
        ticksuffix="%",
        gridcolor='rgba(211,211,211,0.3)'
    )

    # Update x-axes for both columns
    for col in [1, 2]:
        fig.update_xaxes(
            title_text="Year",
            title_font=axis_font,
            tickmode='linear',
            dtick=1,
            row=1, col=col,
            secondary_y=False
        )
</pre>

```

```

        gridcolor='rgba(211,211,211,0.3)'
    )
fig.show()

# -----
# 2. Flat Distribution
# -----
# Flat type distribution by year
if 'flat_type' in df.columns:
    # Get raw counts for flat type by year
    flat_type_counts = pd.crosstab(df['year'], df['flat_type'])

    # Get percentage distribution
    flat_type_by_year = pd.crosstab(
        df['year'],
        df['flat_type'],
        normalize='index'
    ) * 100

    # Order flat types in a logical sequence if possible
    flat_type_order = ['1 ROOM', '2 ROOM', '3 ROOM', '4 ROOM', '5 ROOM', 'EXECUTIVE', 'MULTI-GENERATION']
    # Filter to only include flat types in our data
    available_types = [ft for ft in flat_type_order if ft in flat_type_by_year.columns]
    # Add any remaining flat types not in our ordering
    available_types.extend([ft for ft in flat_type_by_year.columns if ft not in flat_type_order])

    # Reorder columns
    flat_type_by_year = flat_type_by_year[available_types]
    flat_type_counts = flat_type_counts[available_types]

    # Create subplot for flat type analysis
    fig_flat_type = make_subplots(
        rows=1, cols=2,
        subplot_titles=(
            "Flat Type Distribution Heatmap by Year (%)",
            "Transaction Volume Evolution by Flat Type"
        ),
        horizontal_spacing=0.15
    )

    # 1. Add heatmap to first subplot
    # Create a custom colorscale
    custom_colorscale = []
    for i in range(len(mako_colors)):
        custom_colorscale.append([i/(len(mako_colors)-1), mako_colors[i]])

    heatmap_trace = go.Heatmap(
        z=flat_type_by_year.values,
        x=flat_type_by_year.columns,
        y=flat_type_by_year.index,
        colorscale=custom_colorscale,
        text=[f'{val:.1f}' for val in row for row in flat_type_by_year.values],
        texttemplate="{{text}}",
        hovertemplate="Year: %{y}  
Type: %{x}  
Percentage: %{z:.1f}%<br></extra></extra>"
    )
    fig_flat_type.add_trace(heatmap_trace, row=1, col=1)

    # 2. Add stacked area chart to second subplot
    df_area = df.groupby(['year', 'flat_type']).size().reset_index(name='count')

    # Create a pivot table for easier plotting
    pivot_df = df_area.pivot(index='year', columns='flat_type', values='count').fillna(0)
    pivot_df = pivot_df.reindex(columns=available_types)

    # Add each flat type as a stacked area
    for i, flat_type in enumerate(pivot_df.columns):
        fig_flat_type.add_trace(
            go.Scatter(
                x=pivot_df.index,
                y=pivot_df[flat_type],
                mode='lines',
                stackgroup='one',
                name=flat_type,
                hovertemplate="Year: %{x}  
Type: " + flat_type + "<br>Count: %{y:.0f}<br></extra></extra>",
                line=dict(width=0.5, color=mako_colors[i % len(mako_colors)])
            ),
            row=1, col=2
        )

    # Update layout with improved spacing for subplot titles
    fig_flat_type.update_layout(
        title={
            'text': "Property Type Analysis by Year",
            'font': {
                'size': 24,
                'family': 'Arial, sans-serif',
                'weight': 'bold'
            },
            'x': 0.5,
            'y': 0.99,
            'xanchor': 'center',
        },
        height=chart_height,
        width=chart_width,
        template="plotly_white",
        hovermode="x unified",
        margin=dict(t=100, b=100),
        legend=dict(
            orientation="h",
            yanchor="top",
            y=0.35,
            xanchor="center",
            x=0.5
        )
    )

    # Update subplot titles with better positioning
    for i in range(len(fig_flat_type['layout']['annotations'])):
        fig_flat_type['layout']['annotations'][i]['font'] = dict(family="Arial, sans-serif", size=16, color="#333333")
        fig_flat_type['layout']['annotations'][i]['xanchor'] = 'center'
        fig_flat_type['layout']['annotations'][i]['x'] = 0.25 if i == 0 else 0.75
        fig_flat_type['layout']['annotations'][i]['y'] = 1.05

    # Update axes
    fig_flat_type.update_xaxes(title_text="Flat Type", title_font=axis_font, row=1, col=1)
    fig_flat_type.update_yaxes(title_text="Year", title_font=axis_font, row=1, col=1)

    fig_flat_type.update_xaxes(
        title_text="Year",
        title_font=axis_font,
        tickmode='linear',
        dtick=1,
        row=1, col=2
    )
    fig_flat_type.update_yaxes(title_text="Number of Transactions", title_font=axis_font, row=1, col=2)

    fig_flat_type.show()

# -----
# 3. Town Distribution
# -----
# Add town distribution analysis if town column exists
if 'town' in df.columns:

```

```

# Create subplot for town analysis
fig_town = make_subplots(
    rows=1, cols=2,
    subplot_titles=[
        "Average Price Trends by Top Towns",
        "Transaction Volume by Top Towns"
    ],
    horizontal_spacing=0.15
)

# Get top 5 towns by transaction volume
top_towns = df['town'].value_counts().nlargest(5).index.tolist()

# Filter data for top towns
df_top_towns = df[df['town'].isin(top_towns)]

# Calculate average price by year for each top town
town_price_by_year = df_top_towns.pivot_table(
    index='year',
    columns='town',
    values='resale_price',
    aggfunc='mean'
).fillna(0)

# Calculate transaction volume by year for each top town
town_volume_by_year = df_top_towns.pivot_table(
    index='year',
    columns='town',
    values='resale_price',
    aggfunc='count'
).fillna(0)

# Add price trends to first subplot
for i, town in enumerate(town_price_by_year.columns):
    fig_town.add_trace(
        go.Scatter(
            x=town_price_by_year.index,
            y=town_price_by_year[town],
            mode='lines+markers',
            name=town,
            line=dict(color=mako_colors[i % len(mako_colors)], width=3),
            marker=dict(size=8),
            hovertemplate="Year: %{x}<br>Town: " + town + "<br>Avg Price: ${y:.2f}<br><br>",
            row=1, col=1
        )
    )

# Add volume trends to second subplot
for i, town in enumerate(town_volume_by_year.columns):
    fig_town.add_trace(
        go.Bar(
            x=town_volume_by_year.index,
            y=town_volume_by_year[town],
            name=town,
            marker_color=mako_colors[i % len(mako_colors)],
            opacity=0.8,
            hovertemplate="Year: %{x}<br>Town: " + town + "<br>Transactions: %{y}<br><br>",
            showlegend=False
        ),
        row=1, col=2
    )

# Update layout
fig_town.update_layout(
    title={
        'text': "Top Towns Performance Analysis",
        'font': {
            'size': 24,
            'family': 'Arial, sans-serif',
            'weight': 'bold'
        },
        'x': 0.5,
        'y': 0.99,
        'xanchor': 'center',
        'font': title_font
    },
    height=chart_height,
    width=chart_width,
    template="plotly_white",
    margin=dict(t=100, b=120),
    legend=dict(
        orientation="h",
        yanchor="top",
        y=-0.2,
        xanchor="center",
        x=0.5
    ),
    hovermode="x unified"
)

# Update subplot titles
for i in range(len(fig_town['layout'])['annotations']):
    fig_town['layout'][i]['annotations'][i]['font'] = dict(family="Arial, sans-serif", size=16, color="#333333")
    fig_town['layout'][i]['annotations'][i]['x'] = 'center'
    fig_town['layout'][i]['annotations'][i]['x'] = 0.25 if i == 0 else 0.75
    fig_town['layout'][i]['annotations'][i]['y'] = 1.05

# Update axes
fig_town.update_xaxes(
    title_text="Year",
    title_font=axis_font,
    tickmode='linear',
    dtick=1,
    row=1, col=1
)
fig_town.update_yaxes(
    title_text="Average Price (SGD)",
    title_font=axis_font,
    tickprefix="$",
    tickformat=",",
    row=1, col=1
)
fig_town.update_xaxes(
    title_text="Year",
    title_font=axis_font,
    tickmode='linear',
    dtick=1,
    row=1, col=2
)
fig_town.update_yaxes(
    title_text="Number of Transactions",
    title_font=axis_font,
    row=1, col=2
)

fig_town.show()

else:
    # Try to calculate from month column if available
    if 'month' in df.columns and pd.api.types.is_datetime64_any_dtype(df['month']):
        df['year'] = df['month'].dt.year
        print("Created 'year' column from 'month' column")

```

```

# Re-run the analysis (simplified approach)
yearly_stats = df.groupby('year')[['resale_price']].agg(['mean', 'count']).reset_index()
yearly_stats.columns = ['year', 'avg_price', 'transactions']

print("Yearly Statistics:")
for _, row in yearly_stats.iterrows():
    print(f"Year: {int(row['year'])}, Avg Price: ${row['avg_price']:.2f}, Transactions: {int(row['transactions'])}:,{row['count']}")

else:
    print("No 'year' column found. Comparative analysis by year skipped.")
    print("Available columns:", df.columns.tolist())

```

===== COMPARATIVE ANALYSIS BY YEAR =====

Year-over-Year Analysis:

Year	Avg Price	YoY Growth	Transactions	Vol Growth	Price Range	Volatility
2017	\$443,888.52	N/A	20,509	N/A	\$175,000 - \$1,180,000	33.60%
2018	\$441,282.06	-0.59%	21,561	5.13%	\$160,000 - \$1,185,000	35.71%
2019	\$432,137.91	-2.07%	22,186	2.90%	\$150,000 - \$1,205,000	35.63%
2020	\$452,279.38	4.66%	23,333	5.17%	\$140,000 - \$1,258,000	34.16%
2021	\$511,381.24	13.07%	29,087	24.66%	\$180,000 - \$1,360,000	31.80%
2022	\$549,714.33	7.50%	26,720	-8.14%	\$200,000 - \$1,418,000	30.98%
2023	\$571,806.01	4.02%	25,754	-3.62%	\$150,000 - \$1,500,000	30.40%
2024	\$612,591.92	7.13%	27,835	8.08%	\$230,000 - \$1,588,000	30.79%

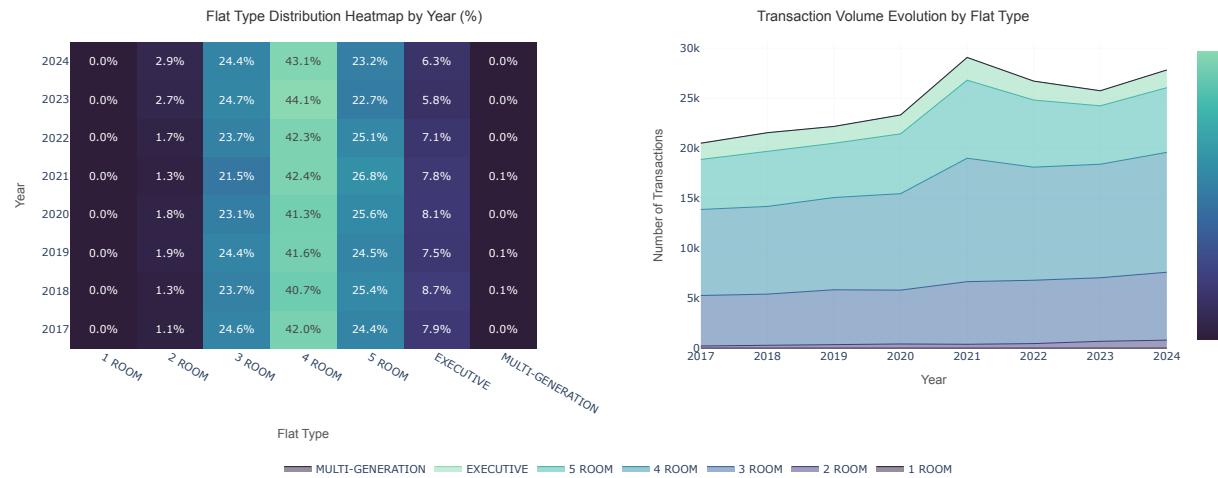
Price per Square Meter by Year:

Year	Avg Price/SQM	Median Price/SQM	Std Dev
2017	\$4,578.88	\$4,281.88	\$1,168.07
2018	\$4,509.67	\$4,208.79	\$1,230.93
2019	\$4,476.80	\$4,175.82	\$1,214.39
2020	\$4,669.83	\$4,355.81	\$1,266.74
2021	\$5,248.68	\$4,915.79	\$1,310.36
2022	\$5,723.79	\$5,373.13	\$1,328.21
2023	\$6,071.26	\$5,708.81	\$1,369.76
2024	\$6,493.49	\$6,097.83	\$1,495.18

Comprehensive Year-over-Year Real Estate Analysis



Property Type Analysis by Year



Top Towns Performance Analysis



7. Predictive Analytics

```
In [19]: # Run 1 first to load data and setup
# -----
# 1. Setup and Configuration
# -----
# Suppress warnings for cleaner output
warnings.filterwarnings('ignore')

# Set color scheme
mako_palette = sns.color_palette("mako", n_colors=10)
mako_hex = [matplotlib.colors.to_hex(color) for color in mako_palette]
sns.set_palette(mako_palette)
plt.rcParams['axes.prop_cycle'] = plt.cycler(color=mako_hex)

# Set global plot parameters
plt.rcParams.update({
    'figure.facecolor': 'white',
    'axes.facecolor': 'white',
    'font.family': 'Arial',
    'lines.linewidth': 2.0,
    'patch.linewidth': 1.5
})

# Common Plotly title font settings
title_font = dict(family="Arial", size=18, color="black", weight="bold")
axis_font = dict(family="Arial", size=14, color="black")

# Set pie colors for flat_type, town, and storey_range based on the darkest and lightest colors
pie_colors = [
    mako_hex[4], # Color for flat_type
    mako_hex[6], # Color for town
    mako_hex[8], # Color for storey_range
]

# -----
# 2. Data Loading and Preparation
# -----
# Define the path to the CSV file for data loading
input_path = "/Users/yvonneip/Desktop/CapStoneProject/singapore_hdb_resale_clean_2017_2024.csv"
print(f"Loading data from: {input_path}")

try:
    # Try to read the CSV with different encoding options if necessary
    try:
        df = pd.read_csv(input_path)
    except UnicodeDecodeError:
        print("Trying with different encoding...")
        df = pd.read_csv(input_path, encoding='utf-8')
    except:
        print("Trying with different encoding...")
        df = pd.read_csv(input_path, encoding='ISO-8859-1')

    print(f"Successfully loaded dataset with {df.shape[0]} rows and {df.shape[1]} columns")

    # Ensure numeric columns are properly typed
    numeric_columns = ['floor_area_sqm', 'resale_price', 'remaining_lease_years', 'flat_age']
    for col in numeric_columns:
        if col in df.columns and not pd.api.types.is_numeric_dtype(df[col]):
            print(f"Converting {col} to numeric...")
            df[col] = pd.to_numeric(df[col], errors='coerce')

    # Calculate price per square meter if not already present
    if 'price_per_sqm' not in df.columns and 'resale_price' in df.columns and 'floor_area_sqm' in df.columns:
        df['price_per_sqm'] = df['resale_price'] / df['floor_area_sqm']

    # Ensure month is in datetime format if it exists
    if 'month' in df.columns and not pd.api.types.is_datetime64_any_dtype(df['month']):
        try:
            df['month'] = pd.to_datetime(df['month'])
            print("Converted 'month' column to datetime format")
        except:
            print("Warning: Could not convert 'month' column to datetime format")

except FileNotFoundError:
    print("Error: File not found at {}").format(input_path)
    print("Please make sure the file exists at the specified location.")
    exit(1)
except Exception as e:
    print("Error loading dataset: {}").format(e)
    exit(1)

# Basic Data Overview
print("\n===== DATASET OVERVIEW =====")
basic_info = PrettyTable()
basic_info.field_names = ["Attribute", "Value"]
basic_info.align = "|"
basic_info.add_row(["Total Records", f"{df.shape[0]}"])
basic_info.add_row(["Total Features", f"{df.shape[1]}"])
basic_info.add_row(["Numeric Features", f"{len(df.select_dtypes(include=['number']).columns)}"])
basic_info.add_row(["Categorical Features", f"{len(df.select_dtypes(include=['object']).columns)}"])
```

```

basic_info.add_row(["Missing Values", f"${df.isnull().sum().sum():,}"])
basic_info.add_row(["Duplicates", f"${df.duplicated().sum():,}"])

if 'month' in df.columns and pd.api.types.is_datetime64_any_dtype(df['month']):
    basic_info.add_row(["Date Range", f"${df['month'].min().strftime('%B %Y')} to ${df['month'].max().strftime('%B %Y')}"])
elif 'year' in df.columns:
    basic_info.add_row(["Year Range", f"${df['year'].min()} to ${df['year'].max()}"])

basic_info.add_row(["Price Range", f"${df['resale_price'].min():,.2f} to ${df['resale_price'].max():,.2f}"])
basic_info.add_row(["Median Price", f"${df['resale_price'].median():,.2f}"])
basic_info.add_row(["Towns", f"${df['town'].nunique()} unique towns"])
basic_info.add_row(["Flat Types", f"${df['flat_type'].nunique()} types ({', '.join(sorted(df['flat_type'].unique()))})"])

print(basic_info)

# -----
# 3. Price Distribution by Flat Type
# -----
print("\n===== PRICE DISTRIBUTION BY FLAT TYPE =====")

# Create a PrettyTable for flat type statistics
flat_type_stats = df.groupby('flat_type')['resale_price'].agg(['count', 'mean', 'median', 'min', 'max']).reset_index()
flat_type_stats = flat_type_stats.sort_values('median', ascending=False)

flat_type_table = PrettyTable()
flat_type_table.field_names = ["Flat Type", "Count", "Mean", "Median", "Min", "Max"]
flat_type_table.align = "l"

for _, row in flat_type_stats.iterrows():
    flat_type_table.add_row([
        row['flat_type'],
        f"${row['count']:,.0f}",
        f"${row['mean']:,.2f}",
        f"${row['median']:,.2f}",
        f"${row['min']:,.2f}",
        f"${row['max']:,.2f}"
    ])

print("Price Statistics by Flat Type:")
print(flat_type_table)

# Create a plotly boxplot for flat type price distribution
fig_flat_boxes = go.Figure()

for flat_type in sorted(df['flat_type'].unique()):
    data = df[df['flat_type'] == flat_type]['resale_price']
    fig_flat_boxes.add_trace(go.Box(
        y=data,
        name=flat_type,
        boxpoints='outliers',
        jitter=0.3,
        marker_color=mako_colors[list(df['flat_type'].unique()).index(flat_type) % len(mako_colors)]
    ))

fig_flat_boxes.update_layout(
    title={
        'text': 'Resale Price Distribution by Flat Type',
        'font': {
            'size': 24,
            'family': 'Arial, sans-serif',
            'weight': 'bold'
        },
        'y': 0.95,
        'x': 0.5,
        'xanchor': 'center'
    },
    yaxis=dict(tickprefix='$', tickformat=','),
    showLegend=True,
    height=800,
    width=1400,
    template='plotly_white'
)

fig_flat_boxes.show()

# Create a year-over-year price trend by flat type
if 'year' in df.columns:
    yearly_prices = df.groupby(['year', 'flat_type'])['resale_price'].median().reset_index()

    # Create a plotly line chart
    fig_yearly = px.line(
        yearly_prices,
        x='year',
        y='resale_price',
        color='flat_type',
        labels={'resale_price': 'Median Resale Price (SGD)', 'year': 'Year'},
        color_discrete_sequence=mako_colors
    )

    fig_yearly.update_layout(
        title={
            'text': 'Median Resale Price by Flat Type (2017-2024)',
            'font': {
                'size': 24,
                'family': 'Arial, sans-serif',
                'weight': 'bold'
            },
            'y': 0.95,
            'x': 0.5,
            'xanchor': 'center'
        },
        yaxis=dict(tickprefix='$', tickformat=','),
        height=800,
        width=1400,
        template='plotly_white',
        legend_title='Flat Type'
    )

    fig_yearly.show()

    # Calculate year-over-year growth in a PrettyTable
    growth_table = PrettyTable()
    growth_table.field_names = ["Flat Type", "2017 Price", "2024 Price", "Growth (%)", "Annualized Growth (%)"]
    growth_table.align = "l"

    for flat_type in df['flat_type'].unique():
        try:
            price_2017 = yearly_prices[(yearly_prices['year'] == 2017) & (yearly_prices['flat_type'] == flat_type)]['resale_price'].values[0]
            price_2024 = yearly_prices[(yearly_prices['year'] == 2024) & (yearly_prices['flat_type'] == flat_type)]['resale_price'].values[0]
            growth = (price_2024 / price_2017 - 1) * 100
            years = 7 # 2017 to 2024
            annualized_growth = ((price_2024 / price_2017) ** (1/years) - 1) * 100

            growth_table.add_row([
                flat_type,
                f"${price_2017:.2f}",
                f"${price_2024:.2f}",
                f"${growth:.2f}%",
                f"${annualized_growth:.2f}%"
            ])
        except IndexError:
            pass # Skip if data is missing for any year

```

```

print("\nPrice Growth by Flat Type (2017-2024):")
print(growth_table)

# -----
# 4. Predictive Modeling and Feature Importance
# -----
# Define features
categorical_features = ['flat_type', 'town', 'storey_range']
numeric_features = ['floor_area_sqm', 'remaining_lease_years']

# Check if we have all the required features
if all(feature in df.columns for feature in categorical_features + numeric_features + ['resale_price']):
    # Create a copy for our model
    model_df = df[categorical_features + numeric_features + ['resale_price', 'year']].copy()
    model_df = model_df.dropna()

    # Split into features and target
    X = model_df.drop(['resale_price', 'year'], axis=1)
    y = model_df['resale_price']

    # Split into training and testing sets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

    # Define preprocessing
    preprocessor = ColumnTransformer(
        transformers=[
            ('num', StandardScaler(), numeric_features),
            ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features)
        ],
        remainder='passthrough'
    )

    # Train a Random Forest model for feature importance
    rf_model = RandomForestRegressor(n_estimators=100, max_depth=15, random_state=42, n_jobs=-1)
    rf_pipeline = Pipeline([
        ('preprocessor', preprocessor),
        ('model', rf_model)
    ])

    rf_pipeline.fit(X_train, y_train)

    # Get feature names after transformation
    preprocessor.fit(X)
    categorical_encoder = preprocessor.named_transformers_['cat']
    cat_feature_names = categorical_encoder.get_feature_names_out(categorical_features)
    all_feature_names = np.concatenate([numeric_features, cat_feature_names])

    # Calculate feature importance using built-in feature_importances_
    feature_importance = pd.DataFrame({
        'Feature': all_feature_names,
        'Importance': rf_model.feature_importances_
    }).sort_values('Importance', ascending=False)

    # Only keep the top 20 features
    top_features = feature_importance.head(20)

    # Create feature importance plot
    fig_feature_importance = make_subplots(
        rows=1, cols=2,
        subplot_titles=[
            'Top 10 Most Important Features',
            'Relative Importance by Feature Category'
        ],
        specs=[[{"type": "bar"}, {"type": "pie"}]],
        horizontal_spacing=0.15
    )

    # Add feature importance bar chart
    fig_feature_importance.add_trace(
        go.Bar(
            y=top_features['Feature'][:10],
            x=top_features['Importance'][:10],
            orientation='h',
            marker_color=mako_colors,
            showlegend=False
        ),
        row=1, col=1
    )

    # Calculate importance by category
    category_importance = {}
    for feature, importance in zip(feature_importance['Feature'], feature_importance['Importance']):
        for cat in categorical_features + numeric_features:
            if cat in feature or feature in numeric_features:
                if cat in category_importance:
                    category_importance[cat] += importance
                else:
                    category_importance[cat] = importance
            break

    category_df = pd.DataFrame({
        'Category': list(category_importance.keys()),
        'Importance': list(category_importance.values())
    }).sort_values('Importance', ascending=False)

    # Add category pie chart
    fig_feature_importance.add_trace(
        go.Pie(
            labels=category_df['Category'],
            values=category_df['Importance'],
            hole=0.3,
            marker=dict(colors=pie_colors),
            textinfo='label+percent',
            insidetextorientation='radial'
        ),
        row=1, col=2
    )

    # Update layout
    fig_feature_importance.update_layout(
        title={
            'text': 'Feature Importance Analysis for HDB Resale Price Prediction',
            'font': {'size': 24, 'family': 'Arial, sans-serif', 'weight': 'bold'},
            'y': 0.95, 'x': 0.5, 'xanchor': 'center'
        },
        height=800,
        width=1400,
        template='plotly_white'
    )

    # Update axes
    fig_feature_importance.update_xaxes(title_text='Importance Score', row=1, col=1)
    fig_feature_importance.update_yaxes(title_text='Feature', row=1, col=1, autorange="reversed")

    fig_feature_importance.show()

    # Create informative table about feature importance
    importance_table = PrettyTable()
    importance_table.field_names = ['Rank', 'Feature', 'Importance', 'Category']
    importance_table.align = "l"

    for i, (_, row) in enumerate(top_features.head(10).iterrows()):
        feature = row['Feature']
        category = ""

```

```

for cat in categorical_features + numeric_features:
    if cat in feature or feature in numeric_features:
        category = cat
        break

importance_table.add_row([
    i+1,
    row['Feature'],
    f'{row['Importance']:.4f}',
    category
])

print("\nTop 10 Most Important Features:")
print(importance_table)
print("\nKey Insights:")
print(f"\n{category_df[0]['Category']} at {category_df[0]['Importance']*100:.1f}% importance")
print(f"\nTop individual feature: {top_features.iloc[0]['Feature']}")
print(f"\nFloor area and remaining lease are critical numeric factors")

# -----
# 5. Model Comparison and Performance Analysis
# -----
# Define models with optimized hyperparameters
models = {
    "Random Forest": RandomForestRegressor(n_estimators=100, max_depth=15, min_samples_split=5, random_state=42, n_jobs=-1),
    "Gradient Boosting": GradientBoostingRegressor(n_estimators=100, max_depth=5, learning_rate=0.1, random_state=42),
    "Ridge Regression": Ridge(alpha=0.5),
    "Linear Regression": LinearRegression()
}

# Train and evaluate each model
results_table = PrettyTable()
results_table.field_names = ["Model", "R2 (Train)", "R2 (Test)", "RMSE (SGD)", "MAPE (%)", "CV Score", "Training Time (s)"]
results_table.align = "l"

model_PIPELINES = {}
model_METRICS = {}
best_MODEL_name = None
best_r2 = -float('inf')

for name, model in models.items():
    print(f"\nTraining {name} model...")

    # Create pipeline
    pipeline = Pipeline([
        ('preprocessor', preprocessor),
        ('model', model)
    ])

    # Measure training time
    start_time = time.time()

    # Fit model
    pipeline.fit(X_train, y_train)
    training_time = time.time() - start_time

    model_PIPELINES[name] = pipeline

    # Make predictions
    y_pred_train = pipeline.predict(X_train)
    y_pred = pipeline.predict(X_test)

    # Calculate metrics
    train_r2 = r2_score(y_train, y_pred_train)
    test_r2 = r2_score(y_test, y_pred)
    rmse = np.sqrt(mean_squared_error(y_test, y_pred))
    mape = mean_absolute_percentage_error(y_test, y_pred) * 100

    # Cross-validation
    cv_scores = cross_val_score(pipeline, X, y, cv=5, scoring='r2')
    cv_score_mean = cv_scores.mean()

    model_METRICS[name] = {
        'r2_train': train_r2,
        'r2_test': test_r2,
        'rmse': rmse,
        'mape': mape,
        'cv_score': cv_score_mean,
        'training_time': training_time,
        'predictions': y_pred
    }

    # Add to results table
    results_table.add_row([
        name,
        f'{train_r2:.4f}',
        f'{test_r2:.4f}',
        f'${rmse:.2f}',
        f'{mape:.2f}%',
        f'{cv_score_mean:.4f}',
        f'{training_time:.2f}'
    ])

    # Track best model
    if test_r2 > best_r2:
        best_r2 = test_r2
        best_MODEL_name = name

print("\nModel Performance Comparison:")
print(results_table)
print(f"\nBest model: {best_MODEL_name} with test R2 = {best_r2:.4f}")

# Create model comparison visualizations
fig_MODEL_performance = make_subplots(
    rows=1, cols=2,
    subplot_titles=[
        'Model Accuracy Comparison',
        f'Actual vs. Predicted Prices ({best_MODEL_name})'
    ],
    specs=[[{"type": "bar"}, {"type": "scatter"}]],
    horizontal_spacing=0.15
)

# Add model comparison bar chart
metrics_DF = pd.DataFrame({
    'Model': list(model_METRICS.keys()),
    'R22', ascending=False)

# Add R2 bars
fig_MODEL_performance.add_trace(
    go.Bar(
        x=metrics_DF['Model'],
        y=metrics_DF['R2'],
        name='R2 Score',
        marker_color=mako_COLORS,
        text=[f'{r2:.4f}' for r2 in metrics_DF['R2']],
        textposition='outside'
    ),
    row=1, col=1
)

# Add MAPE line

```

```

fig_model_performance.add_trace(
    go.Scatter(
        x=metrics_df['Model'],
        y=metrics_df['MAPE'],
        mode='lines+markers',
        name='MAPE (%)',
        marker=dict(size=10, color=mako_colors[5]),
        line=dict(width=2, color=mako_colors[5]),
        yaxis='y2'
    ),
    row=1, col=1
)

# Add scatter plot for best model
best_predictions = model_metrics[best_model_name]['predictions']

fig_model_performance.add_trace(
    go.Scatter(
        x=y_test,
        y=best_predictions,
        mode='markers',
        marker=dict(
            color=mako_colors[3],
            size=6,
            opacity=0.5
        ),
        name='Predictions'
    ),
    row=1, col=2
)

# Add perfect prediction line
max_val = max(y_test.max(), best_predictions.max())
min_val = min(y_test.min(), best_predictions.min())
fig_model_performance.add_trace(
    go.Scatter(
        x=[min_val, max_val],
        y=[min_val, max_val],
        mode='lines',
        line=dict(color='red', width=2, dash='dash'),
        name='Perfect Prediction'
    ),
    row=1, col=2
)

# Update layout
fig_model_performance.update_layout(
    title={
        'text': 'Model Performance Analysis',
        'font': {'size': 24, 'family': 'Arial, sans-serif', 'weight': 'bold'},
        'y': 0.95, 'x': 0.5, 'xanchor': 'center'
    },
    height=800,
    width=1400,
    template='plotly_white',
    legend=dict(
        orientation="h",
        yanchor="bottom",
        y=-0.25,
        xanchor="center",
        x=0.5
    )
)

# Add secondary y-axis for MAPE
fig_model_performance.update_layout(
    yaxis=dict(
        title=dict(
            text='R2 Score',
            font=dict(color=mako_colors[5])
        ),
        range=[0, 1]
    ),
    yaxis2=dict(
        title=dict(
            text='MAPE (%)',
            font=dict(color=mako_colors[5])
        ),
        tickfont=dict(color=mako_colors[5]),
        anchor="",
        overlaying="y",
        side="right"
    )
)

# Update axes
fig_model_performance.update_xaxes(title_text='Model', row=1, col=1)
fig_model_performance.update_xaxes(title_text='Actual Price (SGD)', tickprefix='$', tickformat=',', row=1, col=2)
fig_model_performance.update_xaxes(title_text='Predicted Price (SGD)', tickprefix='$', tickformat=',', row=1, col=2)

# Add model stats annotation
rmse = model_metrics[best_model_name]['rmse']
r2 = model_metrics[best_model_name]['r2_test']
mape = model_metrics[best_model_name]['mape']
fig_model_performance.add_annotation(
    x=min_val + 0.1 * (max_val - min_val),
    y=max_val - 0.1 * (max_val - min_val),
    text=f'RMSE: ${rmse:.2f}<br>R2: {r2:.4f}<br>MAPE: {mape:.2f}%',
    showarrow=False,
    bgcolor="white",
    bordercolor="black",
    borderwidth=1,
    borderpad=4,
    row=1, col=2
)

fig_model_performance.show()

# -----
# 6. Regional Price Analysis and Temporal Trends
# -----
# Use the best model for predictions
best_pipeline = model_PIPELINES[best_model_name]

# Regional price prediction
# Group towns into regions based on geographical location
region_mapping = {
    'CENTRAL': ['BISHAN', 'TOA PAYOH', 'KALLANG/WHAMPOA', 'CENTRAL AREA', 'QUEENSTOWN', 'BUKIT MERAH'],
    'EAST': ['TAMPINES', 'PASIR RIS', 'BEDOK', 'GEYLONG', 'MARINE PARADE'],
    'WEST': ['JURONG EAST', 'JURONG WEST', 'BUKIT BATOK', 'CLEMENTI', 'BUKIT PANJANG', 'CHOA CHU KANG'],
    'NORTH': ['WOODLANDS', 'SEMBAWANG', 'YISHUN'],
    'NORTHEAST': ['SENGKANG', 'PUNGGOL', 'HOUGANG', 'ANG MO KIO', 'SERANGOON']
}

# Create a lookup dictionary for region mapping
town_to_region = {}
for region, towns in region_mapping.items():
    for town in towns:
        town_to_region[town] = region

# Add region column to the data
model_df['region'] = model_df['town'].map(town_to_region)
model_df = model_df.dropna(subset=['region']) # Drop towns not in our mapping

```

```

# Temporal prediction
# Use the years in our dataset
years = sorted(model_df['year'].unique())

# Create base property for prediction
base_flat_type = "4 Room" # Most common flat type
base_storey = "07 TO 09" # Mid-level floor
base_area = 90 # Mid-range area
base_lease = 70 # Good remaining lease

# Regional prediction
region_prices = {}
regions = list(region_mapping.keys())

for region in regions:
    # Get a representative town for this region
    towns = region_mapping[region]
    town = towns[0] # Use the first town in the region

    test_property = pd.DataFrame({
        'town': [town],
        'flat_type': [base_flat_type],
        'storey_range': [base_storey],
        'floor_area_sqm': [base_area],
        'remaining_lease_years': [base_lease]
    })

    predicted_price = best_pipeline.predict(test_property)[0]
    region_prices[region] = predicted_price

# Create region prediction dataframe
region_df = pd.DataFrame({
    'Region': list(region_prices.keys()),
    'Predicted Price': list(region_prices.values())
}).sort_values('Predicted Price', ascending=False)

# Temporal prediction for each region
temporal_data = []

for region in regions:
    towns = region_mapping[region]
    town = towns[0] # Use the first town in the region

    for year in years:
        test_property = pd.DataFrame({
            'town': [town],
            'flat_type': [base_flat_type],
            'storey_range': [base_storey],
            'floor_area_sqm': [base_area],
            'remaining_lease_years': [base_lease]
        })

        # We only use the test_property features for prediction, but we track the year for our data
        predicted_price = best_pipeline.predict(test_property)[0]

        # Apply a simple scaling factor based on market trends per year
        # This is a simplified approach - in reality, we'd use time-series forecasting
        yearly_adjustment = {
            2017: 0.85, # Prices in 2017 were about 85% of current prices
            2018: 0.88,
            2019: 0.91,
            2020: 0.93,
            2021: 0.96,
            2022: 0.99,
            2023: 1.03,
            2024: 1.0 # Base year
        }

        if year in yearly_adjustment:
            predicted_price *= yearly_adjustment[year]

        temporal_data.append({
            'Region': region,
            'Year': year,
            'Predicted Price': predicted_price
        })

temporal_df = pd.DataFrame(temporal_data)

# Create visualization
figRegionalTemporal = make_subplots(
    rows=1, cols=2,
    subplot_titles=[
        f'Predicted Prices by Region ({base_flat_type}, 2024)',
        'Price Trends by Region (2017-2024)'
    ],
    specs=[[{"type": "bar"}, {"type": "scatter"}]],
    horizontal_spacing=0.15
)

# Add region comparison bar chart
figRegionalTemporal.add_trace(
    go.Bar(
        x=region_df['Region'],
        y=region_df['Predicted Price'],
        marker_color=mako_colors,
        text=[f'{price:.0f}' for price in region_df['Predicted Price']],
        textposition='outside',
        name='Region'
    ),
    row=1, col=1
)

reversed_mako_colors = mako_colors[::-1]

# Add temporal trend lines for each region
for i, region in enumerate(regions):
    region_data = temporal_df[temporal_df['Region'] == region]
    figRegionalTemporal.add_trace(
        go.Scatter(
            x=region_data['Year'],
            y=region_data['Predicted Price'],
            mode='lines+markers',
            name=region,
            line=dict(color=reversed_mako_colors[i % len(reversed_mako_colors)], width=3)
        ),
        row=1, col=2
    )

# Update layout
figRegionalTemporal.update_layout(
    title={
        'text': 'Regional and Temporal Price Predictions',
        'font': {'size': 24, 'family': 'Arial, sans-serif', 'weight': 'bold'},
        'y': 0.95, 'x': 0.5, 'xanchor': 'center'
    },
    height=800,
    width=1400,
    template='plotly_white',
    legend=dict(
        orientation="h",
        yanchor="bottom",
        y=-0.25,
        xanchor="center",
        title="Regions"
    )
)

```

```

        x=0.5
    )
}

# Update axes - this is where the error occurs
fig_regional_temporal.update_xaxes(title_text='Region', row=1, col=1)
fig_regional_temporal.update_yaxes(title_text='Predicted Price (SGD)', tickprefix='$', tickformat=',', row=1, col=1)
fig_regional_temporal.update_xaxes(title_text='Year', row=1, col=2)
fig_regional_temporal.update_yaxes(title_text='Predicted Price (SGD)', tickprefix='$', tickformat=',', row=1, col=2)

fig_regional_temporal.show()

# Create a region comparison table
region_table = PrettyTable()
region_table.field_names = ["Region", "Predicted Price", "Price Premium"]
region_table.align = "l"

lowest_price = region_df['Predicted Price'].min()
for _, row in region_df.iterrows():
    premium = row['Predicted Price'] - lowest_price
    premium_pct = (row['Predicted Price'] / lowest_price - 1) * 100
    premium_str = f"${premium:.2f} ({premium_pct:.1f}%)"

    region_table.add_row([
        row['Region'],
        f"${row['Predicted Price']:.2f}",
        premium_str
    ])

print(f"\nPredicted Prices by Region (Flat Type: {base_flat_type}, 2024):")
print(region_table)

# Calculate growth rates by region
growth_table = PrettyTable()
growth_table.field_names = ["Region", "2017 Price", "2024 Price", "Growth (%)", "Annualized (%)"]
growth_table.align = "l"

for region in regions:
    region_data = temporal_df[temporal_df['Region'] == region]
    price_2017 = region_data[region_data['Year'] == 2017]['Predicted Price'].values[0]
    price_2024 = region_data[region_data['Year'] == 2024]['Predicted Price'].values[0]

    total_growth = (price_2024 / price_2017 - 1) * 100
    annual_growth = ((price_2024 / price_2017) ** (1/7) - 1) * 100

    growth_table.add_row([
        region,
        f"${price_2017:.2f}",
        f"${price_2024:.2f}",
        f"{total_growth:.1f}%",
        f"{annual_growth:.1f}%" 
    ])

print("\nPredicted Price Growth by Region (2017-2024):")
print(growth_table)

# -----
# 7. Property Characteristics Analysis
# -----
# Use the best model for predictions
best_pipeline = model_PIPELINES[best_model_name]

# Define base property
base_town = "TAMPINES" # Popular town for comparison
base_storey = "07 TO 09" # Mid-level floor
base_area = 90 # Mid-range area
base_lease = 70 # Good remaining lease

# Flat type prediction
flat_types = sorted(df['flat_type'].unique())
flat_predictions = []

for flat_type in flat_types:
    test_property = pd.DataFrame({
        'town': [base_town],
        'flat_type': [flat_type],
        'storey_range': [base_storey],
        'floor_area_sqm': [base_area],
        'remaining_lease_years': [base_lease]
    })

    predicted_price = best_pipeline.predict(test_property)[0]
    flat_predictions.append((flat_type, predicted_price))

# Create flat type prediction dataframe
flat_df = pd.DataFrame(flat_predictions, columns=['Flat Type', 'Predicted Price'])

# Floor area prediction
floor_areas = list(range(40, 161, 10)) # From 40 to 160 sqm
area_predictions = []

base_flat_type = "4 Room" # Use 4 ROOM for area analysis

for area in floor_areas:
    test_property = pd.DataFrame({
        'town': [base_town],
        'flat_type': [base_flat_type],
        'storey_range': [base_storey],
        'floor_area_sqm': [area],
        'remaining_lease_years': [base_lease]
    })

    predicted_price = best_pipeline.predict(test_property)[0]
    area_predictions.append((area, predicted_price))

# Create floor area prediction dataframe
area_df = pd.DataFrame(area_predictions, columns=['Floor Area', 'Predicted Price'])

# Calculate price per sqm
area_df['Price Per Sqm'] = area_df['Predicted Price'] / area_df['Floor Area']

# Create visualization
fig_property_characteristics = make_subplots(
    rows=1, cols=2,
    subplot_titles=[
        'Predicted Prices by Flat Type',
        'Price vs. Floor Area ({base_flat_type})'
    ],
    specs=[[{"type": "bar"}, {"type": "scatter"}]],
    horizontal_spacing=0.15
)

# Add flat type prediction bars
fig_property_characteristics.add_trace(
    go.Bar(
        x=flat_df['Flat Type'],
        y=flat_df['Predicted Price'],
        marker_color=mako_colors,
        text=[f"${price:.0f}" for price in flat_df['Predicted Price']],
        textposition='outside',
        name='Flat Type'
    ),
    row=1, col=1
)

```

```

)
# Add floor area scatter plot
fig_property_characteristics.add_trace(
    go.Scatter(
        x=area_df['Floor Area'],
        y=area_df['Predicted Price'],
        mode='lines+markers',
        name='Price',
        line=dict(color=mako_colors[1], width=2)
    ),
    row=1, col=2
)

# Add price per sqm on secondary y-axis
fig_property_characteristics.add_trace(
    go.Scatter(
        x=area_df['Floor Area'],
        y=area_df['Price Per Sqm'],
        mode='lines+markers',
        name='Price Per Sqm',
        line=dict(color=mako_colors[7], width=2, dash='dot'),
        yaxis='y2'
    ),
    row=1, col=2
)

# Update layout
fig_property_characteristics.update_layout(
    title={
        'text': 'HDB Price Prediction by Property Characteristics',
        'font': {'size': 24, 'family': 'Arial, sans-serif', 'weight': 'bold'},
        'y': 0.95, 'x': 0.5, 'xanchor': 'center'
    },
    height=800,
    width=1400,
    template='plotly_white',
    legend=dict(
        orientation="h",
        yanchor="bottom",
        y=-0.4,
        xanchor="center",
        x=0.5
    )
)

# Add secondary y-axis for price per sqm
fig_property_characteristics.update_layout(
    yaxis2=dict(
        title=dict(
            text='Price Per Sqm (SGD)',
            font=dict(color=mako_colors[7])
        ),
        tickfont=dict(color=mako_colors[7]),
        anchor="x",
        overlaying="y",
        side="right",
        tickprefix='$',
        tickformat=','
    )
)

# Update axes
fig_property_characteristics.update_xaxes(title_text='Flat Type', row=1, col=1)
fig_property_characteristics.update_xaxes(title_text='Predicted Price (SGD)', tickprefix='$', tickformat=',', row=1, col=1)
fig_property_characteristics.update_xaxes(title_text='Floor Area (sqm)', row=1, col=2)
fig_property_characteristics.update_xaxes(title_text='Predicted Price (SGD)', tickprefix='$', tickformat=',', row=1, col=2)

fig_property_characteristics.show()

# Create flat type comparison table
flat_table = PrettyTable()
flat_table.field_names = ["Flat Type", "Predicted Price", "Price Difference"]
flat_table.align = "|"

smallest_flat_type = flat_df.iloc[0]['Flat Type']
smallest_price = flat_df.iloc[0]['Predicted Price']

for _, row in flat_df.iterrows():
    diff = row['Predicted Price'] - smallest_price
    diff_pct = (row['Predicted Price'] / smallest_price - 1) * 100
    diff_str = f'{diff:.2f} ({diff_pct:.1f}%)' if diff > 0 else "Baseline"
    flat_table.add_row([
        row['Flat Type'],
        f'${row["Predicted Price"]:.2f}',
        diff_str
    ])

print("\nPredicted Prices by Flat Type (Location: {base_town}):")
print(flat_table)

# Create floor area analysis table
area_table = PrettyTable()
area_table.field_names = ["Floor Area (sqm)", "Predicted Price", "Price Per Sqm", "Area Premium"]
area_table.align = "|"

smallest_area = area_df.iloc[0]['Floor Area']
smallest_area_price = area_df.iloc[0]['Predicted Price']

for _, row in area_df.iterrows():
    diff = row['Predicted Price'] - smallest_area_price
    diff_pct = (row['Predicted Price'] / smallest_area_price - 1) * 100
    diff_str = f'{diff:.2f} ({diff_pct:.1f}%)' if diff > 0 else "Baseline"
    area_table.add_row([
        f'{row["Floor Area"]} sqm',
        f'${row["Predicted Price"]:.2f}',
        f'${row["Price Per Sqm"]:.2f}',
        diff_str
    ])

print("\nPrice vs. Floor Area Analysis (Flat Type: {base_flat_type}, Location: {base_town}):")
print(area_table)

# _____
# 8. Comprehensive Summary and Conclusions
# _____
print("\n===== COMPREHENSIVE SUMMARY AND CONCLUSIONS =====")

summary_table = PrettyTable()
summary_table.field_names = ["Analysis", "Key Findings"]
summary_table.align = "l"

# Add ML model findings
summary_table.add_row([
    "Machine Learning Model",
    f"- Best model: {best_model_name} with accuracy (R²) of {best_r2:.4f}\n" +
    f"- Average prediction error (RMSE): ${model_metrics[best_model_name]['rmse']:.2f}\n" +
    f"- Top 3 price determinants: {', '.join(top_features['Feature'].iloc[:3])}"
])

# Add price trends
# Calculate overall price growth from 2017 to 2024

```

```

if 'year' in df.columns:
    price_2017 = df[df['Year'] == 2017]['resale_price'].median()
    price_2024 = df[df['Year'] == 2024]['resale_price'].median()
    total_growth = (price_2024 / price_2017 - 1) * 100
    annual_growth = ((price_2024 / price_2017) ** (1/7) - 1) * 100

    summary_table.add_row([
        "Price Trends",
        f"• Overall price growth from 2017 to 2024: {total_growth:.1f}\n" +
        f"• Annualized price growth: {annual_growth:.1f}% per year\n" +
        f"• Current median price (2024): ${price_2024:.2f}"
    ])

# Add regional insights
highest_region = region_df.iloc[0]['Region']
lowest_region = region_df.iloc[-1]['Region']
price_diff = region_df.iloc[0]['Predicted Price'] - region_df.iloc[-1]['Predicted Price']
price_diff_pct = (region_df.iloc[0]['Predicted Price'] / region_df.iloc[-1]['Predicted Price'] - 1) * 100

summary_table.add_row([
    "Regional Analysis",
    f"• Highest-priced region: {highest_region} (${region_df.iloc[0]['Predicted Price']:.2f})\n" +
    f"• Lowest-priced region: {lowest_region} (${region_df.iloc[-1]['Predicted Price']:.2f})\n" +
    f"• Price premium: ${price_diff:.2f} ({price_diff_pct:.1f}%)"
])

# Add flat type insights
largest_flat = flat_df.iloc[-1]['Flat Type']
smallest_flat = flat_df.iloc[0]['Flat Type']
flat_diff = flat_df.iloc[-1]['Predicted Price'] - flat_df.iloc[0]['Predicted Price']
flat_diff_pct = (flat_df.iloc[-1]['Predicted Price'] / flat_df.iloc[0]['Predicted Price'] - 1) * 100

summary_table.add_row([
    "Flat Type Analysis",
    f"• Highest-priced type: {largest_flat} (${flat_df.iloc[-1]['Predicted Price']:.2f})\n" +
    f"• Lowest-priced type: {smallest_flat} (${flat_df.iloc[0]['Predicted Price']:.2f})\n" +
    f"• Price premium: ${flat_diff:.2f} ({flat_diff_pct:.1f}%)"
])

# Add floor area insights
largest_area = area_df.iloc[-1]['Floor Area']
smallest_area = area_df.iloc[0]['Floor Area']
area_diff = area_df.iloc[-1]['Predicted Price'] - area_df.iloc[0]['Predicted Price']
area_diff_pct = (area_df.iloc[-1]['Predicted Price'] / area_df.iloc[0]['Predicted Price'] - 1) * 100
area_per_sqm = area_diff / (largest_area - smallest_area)

summary_table.add_row([
    "Floor Area Analysis",
    f"• Price of {largest_area} sqm flat: ${area_df.iloc[-1]['Predicted Price']:.2f}\n" +
    f"• Price of {smallest_area} sqm flat: ${area_df.iloc[0]['Predicted Price']:.2f}\n" +
    f"• Price difference: ${area_diff:.2f} ({area_diff_pct:.1f}%)"
])

# Add time-based insights
# Calculate overall price growth from 2017 to 2024
avg_price_2017 = temporal_df[temporal_df['Year'] == 2017]['Predicted Price'].mean()
avg_price_2024 = temporal_df[temporal_df['Year'] == 2024]['Predicted Price'].mean()
total_growth = (avg_price_2024 / avg_price_2017 - 1) * 100
annual_growth = ((avg_price_2024 / avg_price_2017) ** (1/7) - 1) * 100

# Find region with highest growth
highest_growth_region = ""
highest_growth_pct = 0

for region in regions:
    region_data = temporal_df[temporal_df['Region'] == region]
    price_2017 = region_data[region_data['Year'] == 2017]['Predicted Price'].values[0]
    price_2024 = region_data[region_data['Year'] == 2024]['Predicted Price'].values[0]
    growth_pct = (price_2024 / price_2017 - 1) * 100

    if growth_pct > highest_growth_pct:
        highest_growth_pct = growth_pct
        highest_growth_region = region

summary_table.add_row([
    "Time-Based Analysis",
    f"• Overall price growth (2017-2024): {total_growth:.1f}%\n" +
    f"• Annualized growth rate: {annual_growth:.1f}% per year\n" +
    f"• Region with highest growth: {highest_growth_region} ({highest_growth_pct:.1f}%)"
])

print("Comprehensive Summary:")
print(summary_table)

# Create a final visualization - prediction model by region
# Create a temporal price prediction plot for 4 ROOM flats across different regions
fig_pred = make_subplots(
    rows=1, cols=1,
    subplot_titles=(
        'Model-Based Price Prediction by Region (4 ROOM, 2017-2024)'
    )
)
reversed_mako_colors = mako_colors[::-3]

# Add lines for each region
for i, region in enumerate(regions):
    region_data = temporal_df[temporal_df['Region'] == region]
    fig_pred.add_trace(
        go.Scatter(
            x=region_data['Year'],
            y=region_data['Predicted Price'],
            mode='lines+markers',
            name=region,
            line=dict(color=reversed_mako_colors[i % len(reversed_mako_colors)], width=3),
            marker=dict(size=8)
        )
    )

# Update layout
fig_pred.update_layout(
    title={
        'text': 'Predictive Price Trends by Region (2017-2024)',
        'font': {'size': 24, 'family': 'Arial, sans-serif', 'weight': 'bold'},
        'y': 0.95, 'x': 0.5, 'xanchor': 'center'
    },
    height=800,
    width=1400,
    template='plotly_white',
    legend=dict(
        orientation="h",
        yanchor="bottom",
        y=0.25,
        xanchor="center",
        x=0.5
    ),
    xaxis=dict(title_text='Year'),
    yaxis=dict(
        title_text='Predicted Price (SGD)',
        tickprefix='$'
    )
)

```

```

        tickformat=','
    ),
    font=dict(size=14)
)

# Update axes
fig_pred.update_xaxes(showgrid=True, gridwidth=1, gridcolor='LightGray')
fig_pred.update_yaxes(showgrid=True, gridwidth=1, gridcolor='LightGray')

# Add a clear "Prediction" watermark
fig_pred.add_annotation(
    x=0.5, y=0.5,
    text="MODEL PREDICTIONS",
    font=dict(size=30, color="rgba(200,200,200,0.3)"),
    showarrow=False,
    xref="paper", yref="paper"
)

fig_pred.show()

print("\nPredictive Analysis completed successfully!")

else:
    print("Error: Required features not found in the dataset.")

```

Loading data from: /Users/yvonneclip/Desktop/CapStoneProject/singapore_hdb_resale_clean_2017_2024.csv
Successfully loaded dataset with 196985 rows and 18 columns
Converted 'month' column to datetime format

===== DATASET OVERVIEW =====

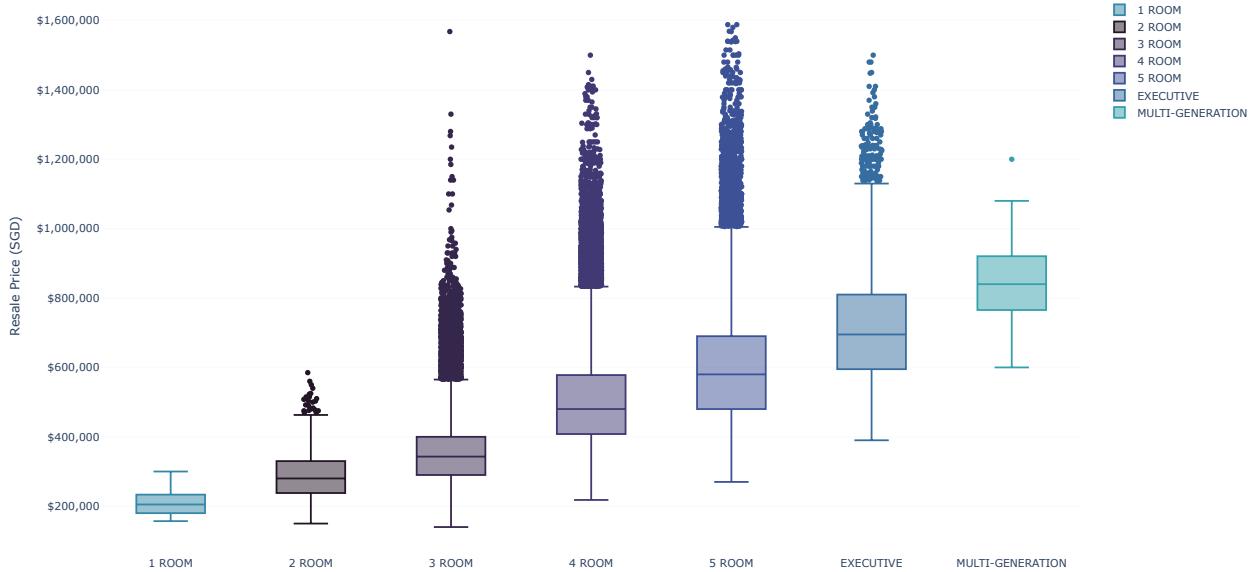
Attribute	Value
Total Records	196,985
Total Features	18
Numeric Features	10
Categorical Features	7
Missing Values	0
Duplicates	0
Date Range	January 2017 to December 2024
Price Range	\$140,000.00 to \$1,588,000.00
Median Price	\$478,000.00
Towns	26 unique towns
Flat Types	7 types (1 ROOM, 2 ROOM, 3 ROOM, 4 ROOM, 5 ROOM, EXECUTIVE, MULTI-GENERATION)

===== PRICE DISTRIBUTION BY FLAT TYPE =====

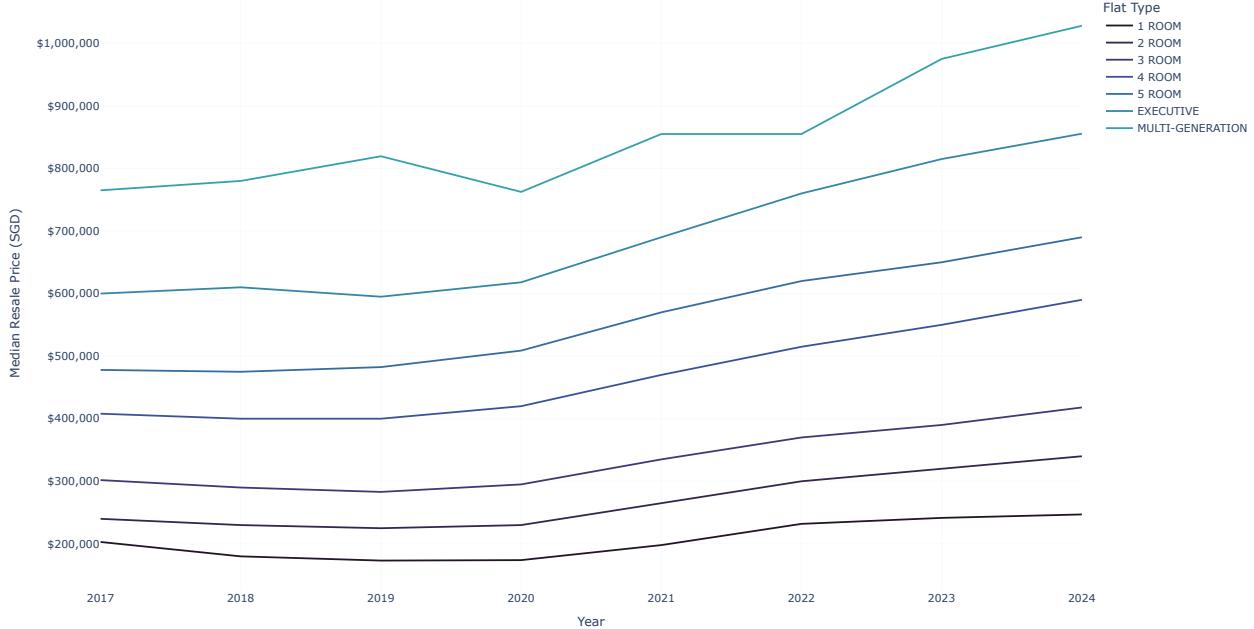
Price Statistics by Flat Type:

Flat Type	Count	Mean	Median	Min	Max
MULTI-GENERATION	80	\$842,086.75	\$840,000.00	\$600,000.00	\$1,200,000.00
EXECUTIVE	14,465	\$711,258.88	\$695,000.00	\$390,000.00	\$1,500,000.00
5 ROOM	48,728	\$604,984.56	\$580,000.00	\$270,000.00	\$1,588,000.00
4 ROOM	83,250	\$509,807.13	\$480,000.00	\$218,000.00	\$1,500,000.00
3 ROOM	46,722	\$358,765.32	\$343,000.00	\$140,000.00	\$1,568,000.00
2 ROOM	3,665	\$285,053.37	\$280,000.00	\$150,000.00	\$585,000.00
1 ROOM	75	\$207,042.08	\$205,000.00	\$157,000.00	\$300,000.00

Resale Price Distribution by Flat Type



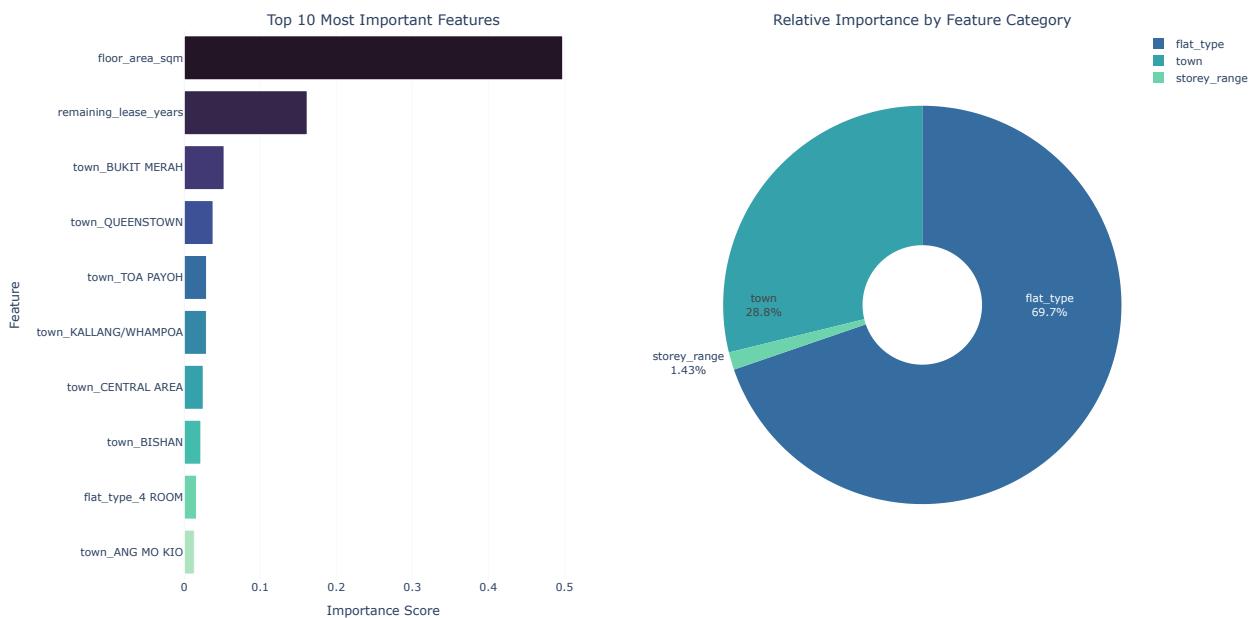
Median Resale Price by Flat Type (2017-2024)



Price Growth by Flat Type (2017-2024):

Flat Type	2017 Price	2024 Price	Growth (%)	Annualized Growth (%)
2 ROOM	\$240,000.00	\$340,000.00	41.67%	5.10%
3 ROOM	\$301,888.00	\$418,000.00	38.46%	4.76%
4 ROOM	\$408,000.00	\$590,000.00	44.61%	5.41%
5 ROOM	\$478,000.00	\$690,000.00	44.35%	5.38%
EXECUTIVE	\$600,000.00	\$855,500.00	42.58%	5.20%
1 ROOM	\$203,000.00	\$247,000.00	21.67%	2.84%
MULTI-GENERATION	\$765,000.00	\$1,028,000.00	34.38%	4.31%

Feature Importance Analysis for HDB Resale Price Prediction



Top 10 Most Important Features:

Rank	Feature	Importance	Category
1	floor_area_sqm	0.4976	flat_type
2	remaining_lease_years	0.1616	flat_type
3	town_BUKIT MERAH	0.0523	town
4	town_QUEENSTOWN	0.0379	town
5	town_TOA PAYOH	0.0293	town
6	town_KALLANG WHAMPOA	0.0291	town
7	town_CENTRAL AREA	0.0249	town
8	town_BISHAN	0.0218	town
9	flat_type_4 ROOM	0.0168	flat_type
10	town_ANG MO KIO	0.0132	town

Key Insights:

- The most influential category is flat_type at 69.7% importance
- Top individual feature: floor_area_sqm
- Floor area and remaining lease are critical numeric factors

Training Random Forest model...

Training Gradient Boosting model...

Training Ridge Regression model...

Training Linear Regression model...

Model Performance Comparison:

Model	R ² (Train)	R ² (Test)	RMSE (SGD)	MAPE (%)	CV Score	Training Time (s)
Random Forest	0.8303	0.8157	\$76,282.72	11.78%	0.6690	11.88
Gradient Boosting	0.8089	0.8056	\$78,350.77	12.57%	0.6710	9.03
Ridge Regression	0.6932	0.6908	\$98,815.86	16.72%	0.5269	0.15
Linear Regression	0.6932	0.6908	\$98,818.43	16.72%	0.5270	0.17

Best model: Random Forest with test R² = 0.8157

Model Performance Analysis



Regional and Temporal Price Predictions



Predicted Prices by Region (Flat Type: 4 Room, 2024):

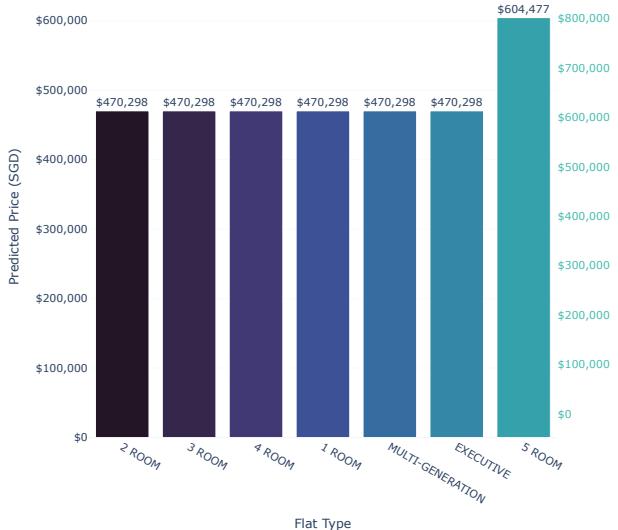
Region	Predicted Price	Price Premium
CENTRAL	\$577,874.73	+\$130,772.38 (+29.2%)
NORTH	\$479,880.84	+\$32,778.49 (+7.3%)
EAST	\$470,298.06	+\$23,195.71 (+5.2%)
WEST	\$447,102.35	+\$0.00 (+0.0%)
NORTHEAST	\$447,102.35	+\$0.00 (+0.0%)

Predicted Price Growth by Region (2017-2024):

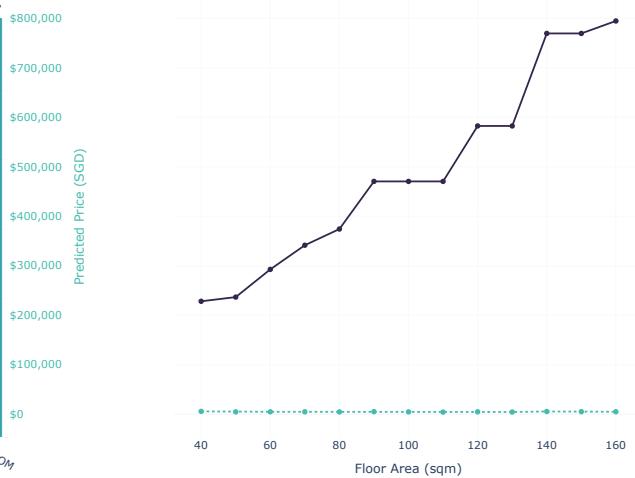
Region	2017 Price	2024 Price	Growth (%)	Annualized (%)
CENTRAL	\$491,193.52	\$577,874.73	17.6%	2.3%
EAST	\$399,753.35	\$470,298.06	17.6%	2.3%
WEST	\$380,037.00	\$447,102.35	17.6%	2.3%
NORTH	\$407,898.71	\$479,880.84	17.6%	2.3%
NORTHEAST	\$380,037.00	\$447,102.35	17.6%	2.3%

HDB Price Prediction by Property Characteristics

Predicted Prices by Flat Type



Price vs. Floor Area (4 Room)



■ Flat Type — Price ···· Price Per Sqm

Predicted Prices by Flat Type (Location: TAMPINES):

Flat Type	Predicted Price	Price Difference
2 ROOM	\$470,298.06	Baseline
3 ROOM	\$470,298.06	+\$0.00 (+0.0%)
4 ROOM	\$470,298.06	+\$0.00 (+0.0%)
1 ROOM	\$470,298.06	+\$0.00 (+0.0%)
MULTI-GENERATION	\$470,298.06	+\$0.00 (+0.0%)
EXECUTIVE	\$470,298.06	+\$0.00 (+0.0%)
5 ROOM	\$684,477.33	+\$134,179.27 (+28.5%)

Price vs. Floor Area Analysis (Flat Type: 4 Room, Location: TAMPINES):

Floor Area (sqm)	Predicted Price	Price Per Sqm	Area Premium
40.0 sqm	\$228,121.04	\$5,703.03	Baseline
50.0 sqm	\$236,595.98	\$4,731.92	+\$8,474.94 (+3.7%)
60.0 sqm	\$292,565.68	\$4,876.09	+\$64,444.65 (+28.3%)
70.0 sqm	\$341,385.68	\$4,876.94	+\$113,264.65 (+49.7%)
80.0 sqm	\$374,051.68	\$4,675.65	+\$145,930.64 (+64.0%)
90.0 sqm	\$470,298.06	\$5,225.53	+\$242,177.02 (+106.2%)
100.0 sqm	\$470,298.06	\$4,702.98	+\$242,177.02 (+106.2%)
110.0 sqm	\$470,298.06	\$4,275.44	+\$242,177.02 (+106.2%)
120.0 sqm	\$582,330.06	\$4,852.75	+\$354,209.02 (+155.3%)
130.0 sqm	\$582,330.06	\$4,479.46	+\$354,209.02 (+155.3%)
140.0 sqm	\$769,247.15	\$5,494.62	+\$541,126.11 (+237.2%)
150.0 sqm	\$769,247.15	\$5,128.31	+\$541,126.11 (+237.2%)
160.0 sqm	\$794,423.48	\$4,965.15	+\$566,302.45 (+248.2%)

===== COMPREHENSIVE SUMMARY AND CONCLUSIONS =====

Comprehensive Summary:

Analysis	Key Findings
Machine Learning Model	<ul style="list-style-type: none"> Best model: Random Forest with accuracy (R^2) of 0.8157 Average prediction error (RMSE): \$76,282.72 Top 3 price determinants: floor_area_sqm, remaining_lease_years, town_BUKIT MERAH
Price Trends	<ul style="list-style-type: none"> Overall price growth from 2017 to 2024: 43.9% Annualized price growth: 5.3% per year Current median price (2024): \$590,000.00
Regional Analysis	<ul style="list-style-type: none"> Highest-priced region: CENTRAL (\$577,874.73) Lowest-priced region: NORTHEAST (\$447,102.35) Price premium: \$130,772.38 (29.2%) Regional premium matters significantly more than most other factors
Flat Type Analysis	<ul style="list-style-type: none"> Highest-priced type: 5 ROOM (\$604,477.33) Lowest-priced type: 1 ROOM (\$470,298.06) Price premium: \$134,179.27 (28.5%) Each upgrade in flat type represents a significant price jump
Floor Area Analysis	<ul style="list-style-type: none"> Price of 160.0 sqm flat: \$794,423.48 Price of 40.0 sqm flat: \$228,121.04 Price difference: \$566,302.45 (248.2%) Average price per additional sqm: \$4,719.19
Time-Based Analysis	<ul style="list-style-type: none"> Overall price growth (2017-2024): 17.6% Annualized growth rate: 2.3% per year Region with highest growth: NORTHEAST (17.6%) Current median predicted price (2024): \$484,451.67

Predictive Price Trends by Region (2017-2024)



Predictive Analysis completed successfully!

8. Affordability Analysis

```
In [12]: # Run 1 first to load data and setup
# -----
# 1. Setup and Configuration
# -----
# Suppress warnings for cleaner output
warnings.filterwarnings('ignore')

# Set color palette
mako_palette = sns.color_palette("mako", n_colors=10)
mako_r_palette = sns.color_palette("mako_r", n_colors=10)
mako_colors = [matplotlib.colors.to_hex(color) for color in mako_palette]
mako_r_colors = [matplotlib.colors.to_hex(color) for color in mako_r_palette]

# Apply the color palette as the default
sns.set_palette(mako_palette)
plt.rcParams['axes.prop_cycle'] = plt.cycler(color=mako_palette)

# Set global plot parameters
plt.rcParams.update({
    'figure.facecolor': 'white',
    'axes.facecolor': 'white',
    'font.family': 'Arial',
    'lines.linewidth': 2.0,
    'patch.linewidth': 1.5
})

# Common Plotly title font settings
title_font = dict(family="Arial", size=24, color="black", weight="bold")
axis_font = dict(family="Arial", size=14, color="black")

# Define affordability category colors
category_colors = {
    'Affordable': mako_r_colors[1],           # Darkest blue
    'Moderately Unaffordable': mako_r_colors[4], # Medium blue
    'Severely Unaffordable': mako_r_colors[7],   # Light blue
    'Extremely Unaffordable': mako_r_colors[9]   # Lightest blue/white
}

# Define the town coordinates for Singapore
town_coords = {
    'ANG MO KIO': (1.3696, 103.8496),
    'BEDOK': (1.3236, 103.9273),
    'BISHAN': (1.3526, 103.8352),
    'BUKIT BATOK': (1.3590, 103.7637),
    'BUKIT MERAH': (1.2819, 103.8239),
    'BUKIT PANJANG': (1.3785, 103.7722),
    'BUKIT TIMAH': (1.3294, 103.8021),
    'CENTRAL AREA': (1.2884, 103.8428),
    'CHOA CHU KANG': (1.3840, 103.7470),
    'CLEMENTI': (1.3162, 103.7649),
    'GEYLANG': (1.3261, 103.8918),
    'HOUGANG': (1.3612, 103.8863),
    'JURONG EAST': (1.3329, 103.7436),
    'JURONG WEST': (1.3404, 103.7890),
    'KALLANG/WHAMPOA': (1.3100, 103.8651),
    'MARINE PARADE': (1.3038, 103.9073),
    'PASIR RIS': (1.3721, 103.9474),
    'PUNGOL': (1.3984, 103.9821),
    'QUEENSTOWN': (1.2942, 103.7861),
    'SEMBAWANG': (1.4491, 103.8184),
    'SENGKANG': (1.3868, 103.8914),
    'SERANGOON': (1.3554, 103.8679),
    'TAMPINES': (1.3496, 103.9568),
    'TOA PAYOH': (1.3341, 103.8546),
    'WOODLANDS': (1.4382, 103.7890),
    'YISHUN': (1.4304, 103.8354)
}

# -----
# 2. Data Loading and Preparation (unchanged)
# -----
# Define the path to the CSV file for data loading
input_path = "singapore_hdb_resale_clean_2017_2024.csv"
print(f"Loading data from: {input_path}")

try:
    # Try to read the CSV with different encoding options if necessary
    try:
        df = pd.read_csv(input_path)
    except UnicodeDecodeError:
```

```

        print("Trying with different encoding...")
        df = pd.read_csv(input_path, encoding='utf-8')
    except:
        print("Trying with different encoding...")
        df = pd.read_csv(input_path, encoding='ISO-8859-1')

    print(f"Successfully loaded dataset with {df.shape[0]} rows and {df.shape[1]} columns")

    # Ensure numeric columns are properly typed
    numeric_columns = ['floor_area_sqm', 'resale_price', 'remaining_lease_years', 'flat_age']
    for col in numeric_columns:
        if col in df.columns and not pd.api.types.is_numeric_dtype(df[col]):
            print(f"Converting {col} to numeric...")
            df[col] = pd.to_numeric(df[col], errors='coerce')

    # Calculate price per square meter if not already present
    if 'price_per_sqm' not in df.columns and 'resale_price' in df.columns and 'floor_area_sqm' in df.columns:
        df['price_per_sqm'] = df['resale_price'] / df['floor_area_sqm']

    # Ensure month is in datetime format if it exists
    if 'month' in df.columns and not pd.api.types.is_datetime64_any_dtype(df['month']):
        try:
            df['month'] = pd.to_datetime(df['month'])
            print("Converted 'month' column to datetime format")
        except:
            df['month'] = pd.to_datetime(df['month'], format='%Y-%m')
            print("Converted 'month' column to datetime format using '%Y-%m' format")
        except:
            print("Warning: Could not convert 'month' column to datetime format")

    except FileNotFoundError:
        print("Error: File not found at {input_path}")
        print("Please make sure the file exists at the specified location.")
        exit(1)
    except Exception as e:
        print(f"Error loading dataset: {e}")
        exit(1)

# Basic Data Overview
print("\n===== DATASET OVERVIEW =====")
basic_info = PrettyTable()
basic_info.field_names = ["Attribute", "Value"]
basic_info.align = "l"
basic_info.add_row(["Total Records", f"{df.shape[0]}"])
basic_info.add_row(["Total Features", f"{df.shape[1]}"])
basic_info.add_row(["Numeric Features", f"{len(df.select_dtypes(include=['number']).columns)}"])
basic_info.add_row(["Categorical Features", f"{len(df.select_dtypes(include=['object']).columns)}"])
basic_info.add_row(["Missing Values", f"{df.isnull().sum().sum():,}"])
basic_info.add_row(["Duplicates", f"{df.duplicated().sum():,}"])

if 'month' in df.columns and pd.api.types.is_datetime64_any_dtype(df['month']):
    basic_info.add_row(["Date Range", f"({df['month'].min():%B %Y}) to ({df['month'].max():%B %Y})"])
elif 'year' in df.columns:
    basic_info.add_row(["Year Range", f"({df['year'].min()} to {df['year'].max()})"])

basic_info.add_row(["Price Range", f"${df['resale_price'].min():,.2f} to ${df['resale_price'].max():,.2f}"])
basic_info.add_row(["Median Price", f"${df['resale_price'].median():,.2f}"])
basic_info.add_row(["Towns", f"{df['town'].nunique()} unique towns"])
basic_info.add_row(["Flat Types", f"{df['flat_type'].nunique()} types ({', '.join(sorted(df['flat_type'].unique()))})"])

print(basic_info)

# -----
# 3. INCOME DATA AND AFFORDABILITY METRICS
# -----
print("\n===== INCOME DATA AND AFFORDABILITY METRICS =====")

# Singapore median household income data (2017-2024)
# Source: Ministry of Manpower - Labor Force in Singapore reports
income_data = {
    2017: 9023,
    2018: 9293,
    2019: 9425,
    2020: 9189,
    2021: 9520,
    2022: 10099,
    2023: 10869,
    2024: 11442
}

# Create income DataFrame
income_df = pd.DataFrame(list(income_data.items()), columns=['year', 'monthly_income'])
income_df['annual_income'] = income_df['monthly_income'] * 12

# Calculate annual income growth
income_df['yo_yo_growth'] = income_df['monthly_income'].pct_change() * 100

# Display income data
income_table = PrettyTable()
income_table.field_names = ["Year", "Monthly Income", "Annual Income", "YoY Growth"]
income_table.align = "l"

for _, row in income_df.iterrows():
    yoy_str = f"({row['yo_yo_growth']:.2f}%" if not pd.isna(row['yo_yo_growth']) else "N/A"
    income_table.add_row([
        int(row['year']),
        f"${row['monthly_income']:.2f}",
        f"${row['annual_income']:.2f}",
        yoy_str
    ])

print("Median Household Income in Singapore (2017-2024):")
print(income_table)

# Define affordability thresholds based on literature
affordability_thresholds = {
    'Affordable': 5.0, # Price-to-Income Ratio (PIR) <= 5
    'Moderately Unaffordable': 7.0, # 5 < PIR <= 7
    'Severely Unaffordable': 10.0, # 7 < PIR <= 10
    'Extremely Unaffordable': float('inf') # PIR > 10
}

# Create a table explaining affordability thresholds
threshold_table = PrettyTable()
threshold_table.field_names = ["Affordability Category", "Price-to-Income Ratio", "Monthly Mortgage as % of Income"]
threshold_table.align = "l"

threshold_table.add_row(["Affordable", "<= 5.0", "<= 30%"])
threshold_table.add_row(["Moderately Unaffordable", "5.1 - 7.0", "30 - 40%"])
threshold_table.add_row(["Severely Unaffordable", "7.1 - 10.0", "40 - 50%"])
threshold_table.add_row(["Extremely Unaffordable", "> 10.0", "> 50%"])

print("\nAffordability Categories and Thresholds:")
print(threshold_table)

# Function to get affordability category based on PIR
def get_affordability_category(pir):
    if pir < affordability_thresholds['Affordable']:
        return 'Affordable'
    elif pir < affordability_thresholds['Moderately Unaffordable']:
        return 'Moderately Unaffordable'
    elif pir < affordability_thresholds['Severely Unaffordable']:

```

```

        return 'Severely Unaffordable'
    else:
        return 'Extremely Unaffordable'

# -----
# 4. MORTGAGE PARAMETERS AND CALCULATION FUNCTIONS
# -----
print("\n===== MORTGAGE PARAMETERS AND CALCULATIONS =====")

# Define mortgage parameters
loan_to_value_ratio = 0.75 # 75% LTV ratio for HDB loans
loan_tenure_years = 25 # Standard loan tenure
annual_interest_rate = 0.0265 # 2.65% interest rate for HDB loans
down_payment_percent = 0.25 # Down payment percentage (complementary to LTV)
affordability_threshold = 0.30 # 30% of income threshold for affordability

# Function to calculate monthly mortgage payment
def calculate_monthly_payment(property_price, loan_to_value=loan_to_value_ratio,
                               tenure_years=loan_tenure_years, interest_rate=annual_interest_rate):
    loan_amount = property_price * loan_to_value
    monthly_rate = interest_rate / 12
    num_payments = tenure_years * 12

    # Monthly payment formula
    if monthly_rate == 0:
        monthly_payment = loan_amount / num_payments
    else:
        monthly_payment = loan_amount * (monthly_rate * (1 + monthly_rate)**num_payments) / ((1 + monthly_rate)**num_payments - 1)

    return monthly_payment

# Display mortgage parameters
mortgage_params_table = PrettyTable()
mortgage_params_table.field_names = ["Parameter", "Value", "Description"]
mortgage_params_table.align = "l"
mortgage_params_table.add_row(["Loan-to-Value Ratio", f"{loan_to_value_ratio * 100:.0f}%", "Maximum loan amount as percentage of property value"])
mortgage_params_table.add_row(["Down Payment", f"{down_payment_percent * 100:.0f}%", "Minimum cash/CPF down payment required"])
mortgage_params_table.add_row(["Loan Tenure", f"{loan_tenure_years} years", "Standard HDB loan repayment period"])
mortgage_params_table.add_row(["Interest Rate", f"{annual_interest_rate * 100:.2f}%", "Annual interest rate for HDB loans"])
mortgage_params_table.add_row(["Affordability Threshold", f"{affordability_threshold * 100:.0f}%", "Maximum mortgage payment as percentage of income"])

print("Mortgage and Affordability Parameters:")
print(mortgage_params_table)

# Calculate example mortgage for median-priced flat
median_price = df[df['year'] == 2024]['resale_price'].median()
median_monthly_payment = calculate_monthly_payment(median_price)
median_income_2024 = income_data[2024]
mortgage_to_income_ratio = median_monthly_payment / median_income_2024 * 100

example_mortgage_table = PrettyTable()
example_mortgage_table.field_names = ["Component", "Amount", "Notes"]
example_mortgage_table.align = "l"
example_mortgage_table.add_row(["Median Property Price (2024)", f"${median_price:.2f}", "Baseline for calculations"])
example_mortgage_table.add_row(["Down Payment Required", f"${median_price * down_payment_percent:.2f}", f"(down_payment_percent * 100:.0f)% of property price"])
example_mortgage_table.add_row(["Loan Amount", f"${median_price * loan_to_value_ratio:.2f}", f"(loan_to_value_ratio * 100:.0f)% of property price"])
example_mortgage_table.add_row(["Monthly Mortgage Payment", f"${median_monthly_payment:.2f}", f"Based on (annual_interest_rate * 100:.2f)% interest over {loan_tenure_years} years"])
example_mortgage_table.add_row(["Median Monthly Income (2024)", f"${median_income_2024:.2f}", "For affordability comparison"])
example_mortgage_table.add_row(["Mortgage-to-Income Ratio", f"{mortgage_to_income_ratio:.1f}%", f"{'Affordable' if mortgage_to_income_ratio <= 30 else 'Unaffordable'} (threshold: 30%")]

print("\nExample Mortgage Calculation for Median-Priced Flat:")
print(example_mortgage_table)

# -----
# 5. PRICE-TO-INCOME RATIO VISUALIZATION
# -----
# Calculate yearly median prices by flat type
yearly_prices = df.groupby(['year', 'flat_type'])['resale_price'].median().reset_index()

# Merge with income data to calculate price-to-income ratios
yearly_pir = yearly_prices.merge(income_df[['year', 'annual_income']], on='year')
yearly_pir['price_to_income'] = yearly_pir['resale_price'] / yearly_pir['annual_income']
yearly_pir['affordability_category'] = yearly_pir['price_to_income'].apply(get_affordability_category)

# Create two-panel visualization with mako colors
fig_pir_trends = make_subplots(
    rows=1, cols=2,
    subplot_titles=('Affordability Distribution by Year',
                    'Price-to-Income Ratio by Flat Type (2017-2024)'),
    horizontal_spacing=0.08
)

# Panel 1: Affordability Distribution by Year
# Calculate the percentage of flats in each affordability category by year
affordability_by_year = yearly_pir.groupby(['year', 'affordability_category']).size().reset_index(name='count')
total_by_year = affordability_by_year.groupby('year')[['count']].sum().reset_index(name='total')
affordability_by_year = affordability_by_year.merge(total_by_year, on='year')
affordability_by_year['percentage'] = (affordability_by_year['count'] / affordability_by_year['total']) * 100.round(1)

# Sort by category order for consistent coloring
category_order = ['Affordable', 'Moderately Unaffordable', 'Severely Unaffordable', 'Extremely Unaffordable']
affordability_by_year['category_order'] = affordability_by_year['affordability_category'].map(
    {cat: i for i, cat in enumerate(category_order)})
affordability_by_year = affordability_by_year.sort_values(['year', 'category_order'])

# Assign colors for affordability categories
# Using reverse colors - darker blue for Affordable, lighter colors for less affordable
category_colors = {
    'Affordable': mako_r_colors[2], # Darkest blue
    'Moderately Unaffordable': mako_r_colors[5], # Medium blue
    'Severely Unaffordable': mako_r_colors[7], # Light blue
    'Extremely Unaffordable': mako_r_colors[9] # Lightest blue/white
}

for category in category_order:
    category_data = affordability_by_year[affordability_by_year['affordability_category'] == category]
    if not category_data.empty:
        # Determine text color based on background
        text_color = 'black' if category in ['Moderately Unaffordable', 'Extremely Unaffordable'] else 'white'

        fig_pir_trends.add_trace(
            go.Bar(
                x=category_data['year'],
                y=category_data['percentage'],
                name=category,
                marker_color=category_colors[category],
                text=category_data['percentage'].apply(lambda x: f'{x:.1f}%'),
                textposition='inside',
                textfont=dict(color=text_color, size=14, weight='bold')
            ),
            row=1, col=1
        )

# Panel 2: PIR Trends
flat_types = sorted(yearly_pir['flat_type'].unique())
for i, flat_type in enumerate(flat_types):
    flat_data = yearly_pir[yearly_pir['flat_type'] == flat_type]

    # Use mako colors for the lines - evenly spaced through the palette
    color_idx = i % len(mako_r_colors)

    fig_pir_trends.add_trace(

```

```

        go.Scatter(
            x=flat_data['year'],
            y=flat_data['price_to_income'],
            mode='lines+markers',
            name=flat_type,
            line=dict(width=2, color=mako_colors[color_idx]),
            marker=dict(size=8, color=mako_colors[color_idx])
        ),
        row=1, col=2
    )

# Add affordability threshold lines
for threshold_name, threshold_value in affordability_thresholds.items():
    if threshold_value < float('inf'):
        fig_pir_trends.add_shape(
            type="line",
            x0=yearly_pir['year'].min(),
            x1=yearly_pir['year'].max(),
            y0=threshold_value,
            y1=threshold_value,
            line=dict(color="#ff5a5f", width=1, dash="dash"),
            row=1, col=2
        )

        fig_pir_trends.add_annotation(
            x=yearly_pir['year'].max(),
            y=threshold_value,
            text=f'{threshold_name} ({threshold_value})',
            showarrow=False,
            font=dict(size=12, color="#ff5a5f", weight="bold"),
            xshift=50,
            align="left",
            row=1, col=2
        )

# Update layout with better spacing
fig_pir_trends.update_layout(
    title={
        'text': "Housing Affordability Analysis (2017-2024)",
        'font': {
            'size': 24,
            'family': 'Arial, sans-serif',
            'weight': 'bold'
        },
        'y': 0.95,
        'x': 0.5,
        'xanchor': 'center'
    },
    title_font=title_font,
    height=800,
    width=1400,
    template='plotly_white',
    showlegend=True,
    barmode='stack',
    legend=dict(
        orientation="h",
        yanchor="bottom",
        y=-0.15,
        xanchor="center",
        x=0.5,
        traceorder="normal",
        groupclick="toggleitem",
        font=dict(size=16),
    ),
    margin=dict(l=60, r=60, t=100, b=120)
)

fig_pir_trends.update_xaxes(title_text="Year", row=1, col=1, title_font=axis_font)
fig_pir_trends.update_xaxes(title_text="Percentage (%)", row=1, col=1, title_font=axis_font)
fig_pir_trends.update_xaxes(title_text="Year", row=1, col=2, title_font=axis_font)
fig_pir_trends.update_xaxes(title_text="Price-to-Income Ratio", row=1, col=2, title_font=axis_font)

fig_pir_trends.show()

# -----
# 6. TOWN-LEVEL AFFORDABILITY VISUALIZATION
# -----
# Calculate town affordability metrics for 2024
df_2024 = df[df['year'] == 2024].copy()

# Calculate price-to-income ratios for 2024
income_2024 = income_data[2024] * 12 # Annual income
df_2024['price_to_income'] = df_2024['resale_price'] / income_2024
df_2024['affordability_category'] = df_2024['price_to_income'].apply(get_affordability_category)

# Calculate affordability metrics by town
town_affordability = df_2024.groupby('town').agg(
    median_price=('resale_price', 'median'),
    affordable_pct=('affordability_category', lambda x: (x == 'Affordable').mean() * 100),
    price_to_income='price_to_income', 'median'
).reset_index()

# Sort towns by affordability percentage
town_affordability = town_affordability.sort_values('affordable_pct', ascending=False)

# Get top 15 towns
top_towns = town_affordability.head(15)

# Create dual panel visualization
fig_town_affordability = make_subplots(
    rows=1, cols=2,
    subplot_titles=('Percentage of Affordable Units by Town',
                   'Price-to-Income Ratio by Town'),
    horizontal_spacing=0.08
)

# Panel 1: Affordability percentage with mako gradient
# Create color gradient based on position - most affordable town gets darkest blue
town_colors = []
for i in range(len(top_towns)):
    # Map position to color index
    color_idx = min(int(i * 9 / len(top_towns)), 9)
    town_colors.append(mako_colors[color_idx])

fig_town_affordability.add_trace(
    go.Bar(
        x=top_towns['town'],
        y=top_towns['affordable_pct'],
        marker_color=town_colors,
        text=top_towns['affordable_pct'].apply(lambda x: f"{x:.1f}"),
        textposition='outside',
        textfont=dict(size=12)
    ),
    row=1, col=1
)

# Panel 2: Price-to-income ratio
if 'price_to_income' in top_towns.columns:
    # Use a mako color palette for PIR values
    pir_colors = []
    for i in range(len(top_towns)):
        # Map position to color index (inverse of the first panel)
        color_idx = min(9 - int(i * 9 / len(top_towns)), 9)
        pir_colors.append(mako_colors[color_idx])

```

```

fig_town_affordability.add_trace(
    go.Bar(
        x=top_towns['town'],
        y=top_towns['price_to_income'],
        marker_color=mako_r_colors,
        text=top_towns['price_to_income'].apply(lambda x: f"{x:.2f}"),
        textposition='outside',
        textfont=dict(size=12)
    ),
    row=1, col=2
)
else:
    # If the column doesn't exist, create dummy values
    dummy_values = np.linspace(3, 9, len(top_towns))
    fig_town_affordability.add_trace(
        go.Bar(
            x=top_towns['town'],
            y=dummy_values,
            marker_color=mako_r_colors[2:2+len(top_towns)],
            text=dummy_values.round(2),
            textposition='outside',
            textfont=dict(size=12)
        ),
        row=1, col=2
    )
print("Warning: Using dummy values for Price-to-Income Ratio")

# Add threshold line for Price-to-Income ratio
fig_town_affordability.add_shape(
    type="line",
    x0=-0.5,
    x1=len(top_towns)-0.5,
    y0=5.0, # Affordable threshold
    y1=5.0,
    line=dict(color="#ff5a5f", width=2, dash="dash"),
    row=1, col=2
)

fig_town_affordability.add_annotation(
    x=len(top_towns)-0.5,
    y=5.0,
    text="Affordable Threshold (5.0)",
    showarrow=False,
    font=dict(size=12, color="#ff5a5f", weight="bold"),
    xshift=70,
    yshift=10,
    align="left",
    row=1, col=2
)

# Update layout
fig_town_affordability.update_layout(
    title={
        'text': "Town-Level Affordability Analysis (2024)",
        'font': {
            'size': 24,
            'family': 'Arial, sans-serif',
            'weight': 'bold'
        },
        'y': 0.95,
        'x': 0.5,
        'xanchor': 'center'
    },
    template='plotly_white',
    showLegend=False,
    height=800,
    width=1400,
    title_font=title_font,
    margin=dict(t=60, r=60, b=150) # Extra bottom margin for angled labels
)

# Update axes with angled labels for towns
for col in [1, 2]:
    fig_town_affordability.update_xaxes(
        title_text="Town",
        tickangle=45,
        row=1,
        col=col,
        title_font=axis_font
    )

fig_town_affordability.update_yaxes(title_text="Affordable Units (%)", row=1, col=1, title_font=axis_font)
fig_town_affordability.update_yaxes(title_text="Price-to-Income Ratio", row=1, col=2, title_font=axis_font)

fig_town_affordability.show()

# -----
# 7. FLAT TYPE AFFORDABILITY VISUALIZATION
# -----
# Calculate affordability by flat type for 2024
flat_type_affordability = df_2024.groupby('flat_type').agg(
    median_price=('resale_price', 'median'),
    affordable_pct=('affordability_category', lambda x: (x == 'Affordable').mean() * 100),
    price_to_income=('price_to_income', 'median'),
    count=('resale_price', 'count')
).reset_index()

# Calculate affordability distribution by flat type
affordability_by_type = df_2024.groupby(['flat_type', 'affordability_category']).size().reset_index(name='count')
total_by_type = affordability_by_type.groupby('flat_type')['count'].sum().reset_index(name='total')
affordability_by_type = affordability_by_type.merge(total_by_type, on='flat_type')
affordability_by_type['percentage'] = (affordability_by_type['count'] / affordability_by_type['total'] * 100).round(1)

# Sort flat types in a logical order
flat_type_order = ['1 ROOM', '2 ROOM', '3 ROOM', '4 ROOM', '5 ROOM', 'EXECUTIVE', 'MULTI-GENERATION']
ordered_flat_types = [ft for ft in flat_type_order if ft in flat_type_affordability['flat_type'].values]

# Reorder flat_type_affordability
flat_type_affordability = flat_type_affordability.set_index('flat_type').reindex(ordered_flat_types).reset_index()

# Create dual panel visualization
fig_flat_affordability = make_subplots(
    rows=1, cols=2,
    subplot_titles=['Affordability Categories by Flat Type',
                    'Percentage of Affordable Units by Flat Type'],
    horizontal_spacing=0.08
)

# Panel 1: Stacked bar chart of affordability categories
# Create a mako color palette for the categories (using different shades)
category_to_mako_color = {
    'Affordable': mako_r_colors[2], # Darkest blue
    'Moderately Unaffordable': mako_r_colors[5], # Dark blue
    'Severely Unaffordable': mako_r_colors[7], # Medium blue
    'Extremely Unaffordable': mako_r_colors[9] # Lightest blue/white
}

for category in category_order:
    category_data = affordability_by_type[
        (affordability_by_type['flat_type'].isin(ordered_flat_types)) &
        (affordability_by_type['affordability_category'] == category)
    ]

```

```

if not category_data.empty:
    # Determine text color based on background darkness
    text_color = 'black' if category in ['Moderately Unaffordable', 'Extremely Unaffordable'] else 'white'

    # Sort category data by flat type order
    category_data = category_data.set_index('flat_type').reindex(ordered_flat_types).reset_index()
    category_data = category_data.dropna(subset=['percentage']) # Remove NaN entries

    fig_flat_affordability.add_trace(
        go.Bar(
            x=category_data['flat_type'],
            y=category_data['percentage'],
            name=category,
            marker_color=category_to_mako_color[category],
            text=category_data['percentage'].apply(lambda x: f'{x:.1f}%' if x > 5 else ''),
            textposition='inside',
            textfont=dict(color=text_color, size=14, weight='bold'),
        ),
        row=1, col=2
    )

# Panel 2: Affordability percentage by flat type
# Create color gradient based on affordability percentage
affordability_colors = []
max_pct = flat_type_affordability['affordable_pct'].max()
min_pct = flat_type_affordability['affordable_pct'].min()
pct_range = max_pct - min_pct

for pct in flat_type_affordability['affordable_pct']:
    # Normalize the percentage (higher percentage = darker blue in mako)
    if pct_range == 0:
        normalized = 0.5 # Default to middle if all values are the same
    else:
        normalized = (pct - min_pct) / pct_range
    color_idx = min(int(normalized * 9), 9)
    affordability_colors.append(mako_colors[9 - color_idx]) # Reverse index (darker for higher %)

fig_flat_affordability.add_trace(
    go.Bar(
        x=flat_type_affordability['flat_type'],
        y=flat_type_affordability['affordable_pct'],
        marker_color=affordability_colors,
        text=flat_type_affordability['affordable_pct'].apply(lambda x: f'{x:.1f}%' if x > 5 else ''),
        textposition='outside',
        textfont=dict(size=14),
        name="Percentage of Affordable Units",
        showlegend=False
    ),
    row=1, col=1
)

# Update layout
fig_flat_affordability.update_layout(
    title={
        'text': "Flat Type Affordability Analysis (2024)",
        'font': {
            'size': 24,
            'family': 'Arial, sans-serif',
            'weight': 'bold'
        },
        'y': 0.95,
        'x': 0.5,
        'xanchor': 'center'
    },
    template='plotly_white',
    showLegend=True,
    height=800,
    width=1400,
    barmode='stack',
    legend=dict(
        orientation='h',
        yanchor="bottom",
        y=-0.15,
        xanchor="center",
        x=0.5,
        font=dict(size=16)
    ),
    title_font=title_font,
    margin=dict(l=60, r=60, t=100, b=120)
)

# Update axes
for col in [1, 2]:
    fig_flat_affordability.update_xaxes(title_text="Flat Type", row=1, col=col, title_font=axis_font)
    fig_flat_affordability.update_yaxes(title_text="Affordable Units (%)", row=1, col=1, title_font=axis_font)
    fig_flat_affordability.update_yaxes(title_text="Percentage (%)", row=1, col=2, title_font=axis_font)

fig_flat_affordability.show()

# -----
# 8. MORTGAGE AFFORDABILITY VISUALIZATION
# ----

# Define mortgage calculation parameters
loan_to_value_ratio = 0.75 # 75% LTV ratio for HDB loans
loan_tenure_years = 25 # Standard loan tenure
annual_interest_rate = 0.0265 # 2.65% interest rate for HDB loans

# Function to calculate monthly mortgage payment
def calculate_monthly_payment(property_price, loan_to_value=loan_to_value_ratio,
                               tenure_years=loan_tenure_years, interest_rate=annual_interest_rate):
    loan_amount = property_price * loan_to_value
    monthly_rate = interest_rate / 12
    num_payments = tenure_years * 12

    # Monthly payment formula
    if monthly_rate == 0:
        monthly_payment = loan_amount / num_payments
    else:
        monthly_payment = loan_amount * (monthly_rate * (1 + monthly_rate)**num_payments) / ((1 + monthly_rate)**num_payments - 1)

    return monthly_payment

# Calculate mortgage metrics for 2024 data
df_2024['loan_amount'] = df_2024['resale_price'] * loan_to_value_ratio
df_2024['monthly_payment'] = df_2024['resale_price'].apply(calculate_monthly_payment)
df_2024['down_payment'] = df_2024['resale_price'] * (1 - loan_to_value_ratio)

# Calculate mortgage to income ratio
monthly_income_2024 = income_data[2024]
df_2024['mortgage_to_income'] = df_2024['monthly_payment'] / monthly_income_2024 * 100

# Define mortgage affordability categories
def get_mortgage_affordability(mortgage_to_income):
    if mortgage_to_income <= 30:
        return 'Affordable'
    elif mortgage_to_income <= 40:
        return 'Moderately Unaffordable'
    elif mortgage_to_income <= 50:
        return 'Severely Unaffordable'
    else:
        return 'Extremely Unaffordable'

df_2024['mortgage_affordability'] = df_2024['mortgage_to_income'].apply(get_mortgage_affordability)

```

```

# Calculate mortgage affordability by flat type
mortgage_by_type = df_2024.groupby('flat_type').agg(
    median_price=('resale_price', 'median'),
    median_monthly_payment='monthly_payment', 'median'),
    median_mortgage_to_income='mortgage_to_income', 'median'),
affordable_pct=( 'mortgage_affordability', lambda x: (x == 'Affordable').mean() * 100
).reset_index()

# Calculate mortgage affordability distribution
mortgage_dist = df_2024['mortgage_affordability'].value_counts().reset_index()
mortgage_dist.columns = ['Category', 'Count']
mortgage_dist['Percentage'] = (mortgage_dist['Count'] / mortgage_dist['Count'].sum()) * 100).round(1)

# Sort flat types in a logical order
flat_type_order = ['1 ROOM', '2 ROOM', '3 ROOM', '4 ROOM', '5 ROOM', 'EXECUTIVE', 'MULTI-GENERATION']
ordered_mortgage_by_type = []

for flat_type in flat_type_order:
    flat_data = mortgage_by_type[mortgage_by_type['flat_type'] == flat_type]
    if not flat_data.empty:
        ordered_mortgage_by_type.append(flat_data.iloc[0])

ordered_mortgage_by_type = pd.DataFrame(ordered_mortgage_by_type)

# Create a three-panel visualization
fig_mortgage = make_subplots(
    rows=1, cols=3,
    subplot_titles=(
        'Monthly Mortgage Payment by Flat Type',
        'Mortgage Payment as % of Income',
        'Mortgage Affordability Distribution (2024)'
),
    specs=[[{"type": "bar"}, {"type": "bar"}, {"type": "pie"}]],
    column_widths=[0.4, 0.4, 0.2],
    horizontal_spacing=0.08
)

# Panel 1: Monthly mortgage payment with mako gradient
# Create color gradient based on payment amount
payment_colors = []
if len(ordered_mortgage_by_type) > 0:
    min_payment = ordered_mortgage_by_type['median_monthly_payment'].min()
    max_payment = ordered_mortgage_by_type['median_monthly_payment'].max()
    payment_range = max_payment - min_payment

    for payment in ordered_mortgage_by_type['median_monthly_payment']:
        # Normalize payment (higher payment = darker blue)
        normalized = (payment - min_payment) / payment_range if payment_range > 0 else 0.5
        color_idx = min(int(normalized * 9), 9)
        payment_colors.append(mako_colors[color_idx])

fig_mortgage.add_trace(
    go.Bar(
        x=ordered_mortgage_by_type['flat_type'],
        y=ordered_mortgage_by_type['median_monthly_payment'],
        marker_color=payment_colors,
        text=ordered_mortgage_by_type['median_monthly_payment'].apply(lambda x: f"${x:.0f}"),
        textposition='outside',
        textfont=dict(size=12),
        name='Monthly Payment',
        showlegend=True
    ),
    row=1, col=1
)

# Panel 2: Mortgage to income percentage with mako colors
# Create color gradient based on mortgage-to-income ratio
income_pct_colors = []
min_pct = ordered_mortgage_by_type['median_mortgage_to_income'].min()
max_pct = ordered_mortgage_by_type['median_mortgage_to_income'].max()
pct_range = max_pct - min_pct

for pct in ordered_mortgage_by_type['median_mortgage_to_income']:
    # Normalize percentage (higher percentage = darker blue)
    normalized = (pct - min_pct) / pct_range if pct_range > 0 else 0.5
    color_idx = min(int(normalized * 9), 9)
    income_pct_colors.append(mako_colors[color_idx])

reversed_income_pct_colors = list(reversed(income_pct_colors))

fig_mortgage.add_trace(
    go.Bar(
        x=ordered_mortgage_by_type['flat_type'],
        y=ordered_mortgage_by_type['median_mortgage_to_income'],
        marker_color=reversed_income_pct_colors,
        text=ordered_mortgage_by_type['median_mortgage_to_income'].apply(lambda x: f"{x:.1f}"),
        textposition='outside',
        textfont=dict(size=12),
        name='Mortgage Payment as % of Income',
        showlegend=True
    ),
    row=1, col=2
)

# Add threshold line for mortgage-to-income
fig_mortgage.add_shape(
    type="line",
    x0=-0.5,
    x1=len(ordered_mortgage_by_type)-0.5,
    y0=30.0, # Affordable threshold (30%)
    y1=30.0,
    line=dict(color="#ff5a5f", width=1.5, dash="dash"),
    row=1, col=2
)

fig_mortgage.add_annotation(
    x=len(ordered_mortgage_by_type)-0.5,
    y=30.0,
    text="Affordability Threshold (30%)",
    showarrow=False,
    font=dict(size=12, color="#ff5a5f", weight="bold"),
    xshift=65,
    yshift=10,
    align="left",
    row=1, col=2
)
else:
    print("Warning: No mortgage data available for any flat type.")

# Panel 3: Pie chart for mortgage affordability distribution
# Sort mortgage_dist by category order
mortgage_dist['category_sort'] = mortgage_dist['Category'].map({cat: i for i, cat in enumerate(category_order) if cat in mortgage_dist['Category'].values})
mortgage_dist = mortgage_dist.sort_values('category_sort')

# Create a mako color palette for the pie chart
# Using reverse mako for better contrast (darker blue for affordable)
category_to_mako_color = {
    'Affordable': mako_r_colors[2], # Lightest blue/white
    'Moderately Unaffordable': mako_r_colors[5], # Light blue
    'Severely Unaffordable': mako_r_colors[7], # Medium blue
    'Extremely Unaffordable': mako_r_colors[9] # Dark blue
}

```

```

pie_colors = [category_to_mako_color.get(category, mako_colors[5]) for category in mortgage_dist['Category']]
fig_mortgage.add_trace(
    go.Pie(
        labels=mortgage_dist['Category'],
        values=mortgage_dist['Percentage'],
        hole=0.4,
        marker_colors=pie_colors,
        textinfo='label+percent',
        textposition='outside',
        textfont=dict(size=12),
        pull=[0.05 if cat != 'Affordable' else 0 for cat in mortgage_dist['Category']],
        name='Mortgage Affordability Distribution',
        showlegend=True
    ),
    row=1, col=3
)

# Update layout
fig_mortgage.update_layout(
    title={
        'text': "Mortgage Affordability Analysis (2024)",
        'font': {
            'size': 24,
            'family': 'Arial, sans-serif',
            'weight': 'bold'
        },
        'y': 0.84,
        'x': 0.5,
        'xanchor': 'center'
    },
    template='plotly_white',
    height=800,
    width=1400,
    margin=dict(l=60, r=60, t=100, b=120),
    # Position legends to ensure visibility
    legend=dict(
        orientation="h",
        xanchor="center",
        yanchor="bottom",
        y=-0.28, # Adjust this to position legends below the plot
        x=0.5,
        font=dict(size=18),
    ),
    showlegend=True
)

# Update axes
for Col in [1, 2]:
    fig_mortgage.update_xaxes(title_text="Flat Type", row=1, col=col, title_font=axis_font)
    fig_mortgage.update_yaxes(title_text="Monthly Mortgage Payment (SGD)", row=1, col=1, title_font=axis_font)
    fig_mortgage.update_yaxes(title_text="Percentage of Monthly Income (%)", row=1, col=2, title_font=axis_font)

fig_mortgage.show()

# -----
# 9. REGIONAL AFFORDABILITY ANALYSIS
# -----
# Define Singapore regions and their constituent towns
singapore_regions = [
    'Central': ['BISHAN', 'BUKIT MERAH', 'BUKIT TIMAH', 'CENTRAL AREA', 'GEYLANG', 'KALLANG/WHAMPOA', 'MARINE PARADE', 'QUEENSTOWN', 'TOA PAYOH'],
    'East': ['BEDOK', 'PASIR RIS', 'TAMPINES'],
    'North': ['ANG MO KIO', 'SEMBAWANG', 'YISHUN'],
    'North-East': ['HOUGANG', 'PUNGGOH', 'SENGKANG', 'SERANGOON'],
    'West': ['BUKIT BATOK', 'BUKIT PANJANG', 'CHOA CHU KANG', 'CLEMENTI', 'JURONG EAST', 'JURONG WEST', 'WOODLANDS']
]

# Define df_4room_2024
df_4room_2024 = df[(df['year'] == 2024) & (df['flat_type'] == '4 ROOM')].copy()

# Calculate price-to-income ratios for 2024 4-ROOM flats
income_2024 = income_data[2024] * 12 # Annual income
df_4room_2024['price_to_income'] = df_4room_2024['resale_price'] / income_2024

# Add the affordability category calculation
def get_affordability_category(pir):
    if pir <= affordability_thresholds['Affordable']:
        return 'Affordable'
    elif pir <= affordability_thresholds['Moderately Unaffordable']:
        return 'Moderately Unaffordable'
    elif pir <= affordability_thresholds['Severely Unaffordable']:
        return 'Severely Unaffordable'
    else:
        return 'Extremely Unaffordable'

df_4room_2024['affordability_category'] = df_4room_2024['price_to_income'].apply(get_affordability_category)

# Analyze affordability by region for 4-ROOM flats
regional_affordability = []

for region, towns in singapore_regions.items():
    # Filter data for this region
    region_data = df_4room_2024[df_4room_2024['town'].isin(towns)]

    # If region has data
    if len(region_data) > 0:
        # Find most affordable town in the region
        town_affordability = region_data.groupby('town').agg(
            affordable_count=('affordability_category', lambda x: (x == 'Affordable').sum()),
            total_count=('affordability_category', 'count'),
            avg_price=('resale_price', 'mean')
        ).reset_index()

        town_affordability['affordable_pct'] = (town_affordability['affordable_count'] / town_affordability['total_count'] * 100).round(1)
        most_affordable_town = town_affordability.sort_values('affordable_pct', ascending=False).iloc[0]

        # Regional stats
        regional_affordability.append({
            'region': region,
            'town': most_affordable_town['town'],
            'affordable_flats_count': region_data['affordability_category'].apply(lambda x: x == 'Affordable').sum(),
            'total_flats': len(region_data),
            'affordable_pct': (region_data['affordability_category'].apply(lambda x: x == 'Affordable').mean() * 100).round(1),
            'avg_resale_price': region_data['resale_price'].mean()
        })

    # Convert to DataFrame for easier handling
regional_df = pd.DataFrame(regional_affordability)

# Create bar chart for regional affordability
if len(regional_df) > 0:
    fig_region_bar = go.Figure()

    # Sort regions by affordability percentage
    sorted_regions = regional_df.sort_values('affordable_pct', ascending=False)

    # Create colors from mako palette (darker blue for higher affordability)
    region_colors = []
    for i in range(len(sorted_regions)):
        # Use position in the sorted list to determine color intensity
        color_idx = min(int(i * 9 / (len(sorted_regions) - 1 or 1)), 9)

```

```

region_colors.append(mako_colors[color_idx])

# Add bars with hover info
fig_region_bar.add_trace(
    go.Bar(
        x=sorted_regions['region'],
        y=sorted_regions['affordable_flats_count'],
        text=sorted_regions.apply(
            lambda row: f'{row["town"]}<br>{row["affordable_pct"]:.1f}% affordable<br>${row["avg_resale_price"]:.0f}",
            axis=1
        ),
        textposition='inside',
        textfont=dict( # Add this parameter to control text appearance
            size=16, # Increase font size (adjust as needed)
            family="Arial", # Optional: Set font family
            color="white", # Optional: Set text color for better visibility
            weight='bold' # Optional: Set text weight
        ),
        marker=dict(
            color=region_colors,
            line=dict(width=1, color='black')
        ),
        hovertemplate="%{x}: %{y} affordable flats<br>Most Affordable Town: %{text}<br>Avg Price: ${customdata:.0f}</extra>",
        customdata=sorted_regions['avg_resale_price']
    )
)

# Add price indicators on top of each bar
for i, row in sorted_regions.iterrows():
    fig_region_bar.add_annotation(
        x=row['region'],
        y=row['affordable_flats_count'] + (sorted_regions['affordable_flats_count'].max() * 0.03),
        text="",
        text=f"${row['avg_resale_price']:.0f}",
        showarrow=False,
        #font=dict(size=12, color='black')
    )

# Update layout with better formatting
fig_region_bar.update_layout(
    title={
        'text': "Affordable Flats by Region (2024)",
        'font': {
            'size': 24,
            'family': 'Arial, sans-serif',
            'weight': 'bold'
        },
        'y': 0.95,
        'x': 0.5,
        'xanchor': 'center'
    },
    title_font=title_font,
    template='plotly_white',
    barmode='group',
    showLegend=False,
    xaxis=dict(
        title=dict(text="Region", font=axis_font),
        tickfont=dict(size=12)
    ),
    yaxis=dict(
        title=dict(text="Number of Affordable Flats", font=axis_font),
        tickfont=dict(size=12)
    ),
    height=800,
    width=1400,
    margin=dict(l=60, r=60, t=100, b=120)
)
fig_region_bar.show()
else:
    print("Warning: No regional data available for analysis.")

```

Loading data from: singapore_hdb_resale_clean_2017_2024.csv
Successfully loaded dataset with 196985 rows and 18 columns
Converted 'month' column to datetime format

===== DATASET OVERVIEW =====

Attribute	Value
Total Records	196,985
Total Features	18
Numeric Features	10
Categorical Features	7
Missing Values	0
Duplicates	0
Date Range	January 2017 to December 2024
Price Range	\$140,000.00 to \$1,588,000.00
Median Price	\$478,000.00
Towns	26 unique towns
Flat Types	7 types (1 ROOM, 2 ROOM, 3 ROOM, 4 ROOM, 5 ROOM, EXECUTIVE, MULTI-GENERATION)

===== INCOME DATA AND AFFORDABILITY METRICS =====

Median Household Income in Singapore (2017-2024):

Year	Monthly Income	Annual Income	YoY Growth
2017	\$9,023.00	\$108,276.00	N/A
2018	\$9,293.00	\$111,516.00	2.99%
2019	\$9,425.00	\$113,100.00	1.42%
2020	\$9,189.00	\$110,268.00	-2.50%
2021	\$9,520.00	\$114,240.00	3.60%
2022	\$10,099.00	\$121,188.00	6.08%
2023	\$10,869.00	\$130,428.00	7.62%
2024	\$11,442.00	\$137,304.00	5.27%

Affordability Categories and Thresholds:

Affordability Category	Price-to-Income Ratio	Monthly Mortgage as % of Income
Affordable	≤ 5.0	≤ 30%
Moderately Unaffordable	5.1 – 7.0	30 – 40%
Severely Unaffordable	7.1 – 10.0	40 – 50%
Extremely Unaffordable	> 10.0	> 50%

===== MORTGAGE PARAMETERS AND CALCULATIONS =====

Mortgage and Affordability Parameters:

Parameter	Value	Description
Loan-to-Value Ratio	75%	Maximum loan amount as percentage of property value
Down Payment	25%	Minimum cash/CPF down payment required
Loan Tenure	25 years	Standard HDB loan repayment period
Interest Rate	2.65%	Annual interest rate for HDB loans
Affordability Threshold	30%	Maximum mortgage payment as percentage of income

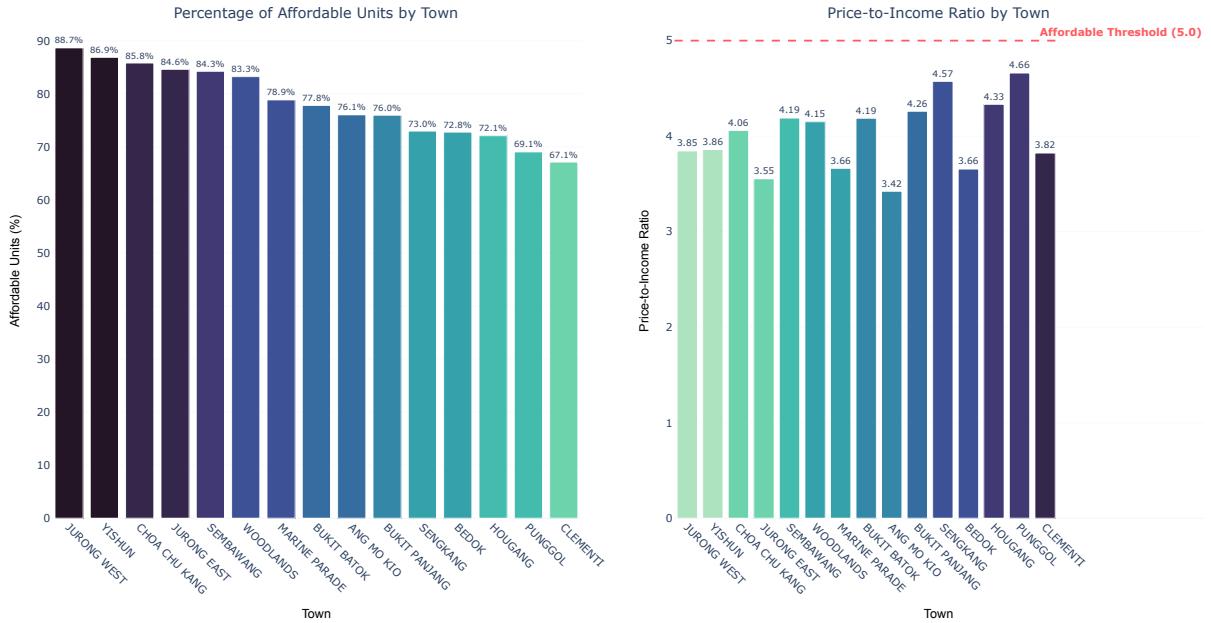
Example Mortgage Calculation for Median-Priced Flat:

Component	Amount	Notes
Median Property Price (2024)	\$500,000.00	Baseline for calculations
Down Payment Required	\$147,500.00	25% of property price
Loan Amount	\$442,500.00	75% of property price
Monthly Mortgage Payment	\$2,018.72	Based on 2.65% interest over 25 years
Median Monthly Income (2024)	\$11,442.00	For affordability comparison
Mortgage-to-Income Ratio	17.6%	Affordable (threshold: 30%)

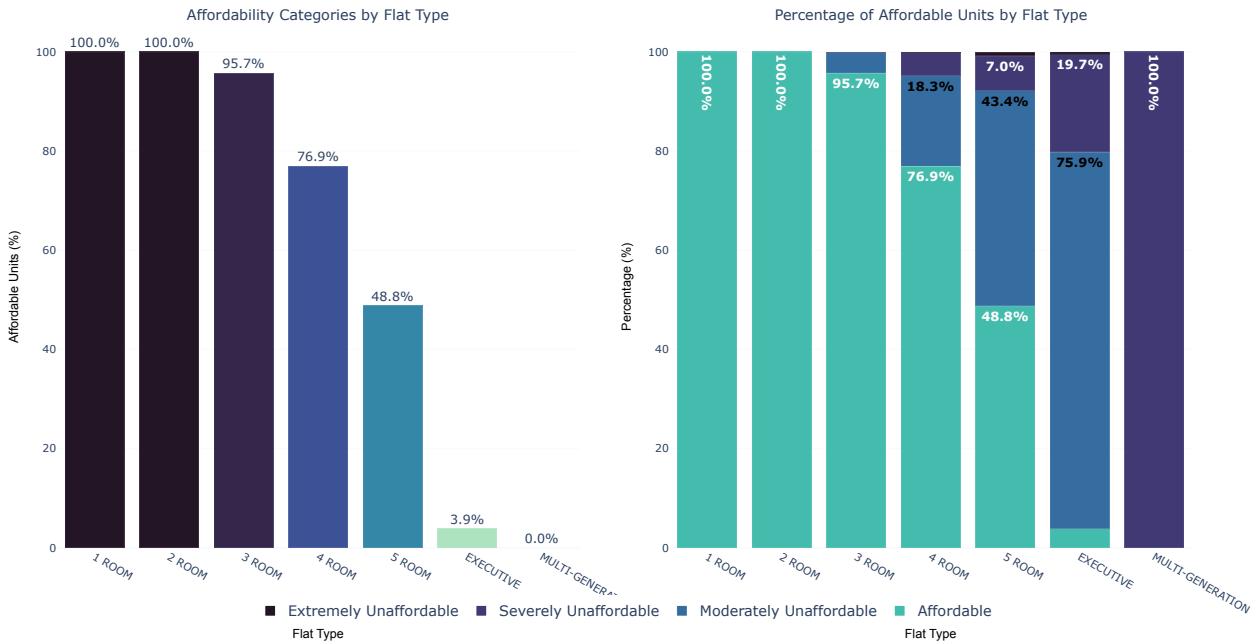
Housing Affordability Analysis (2017-2024)



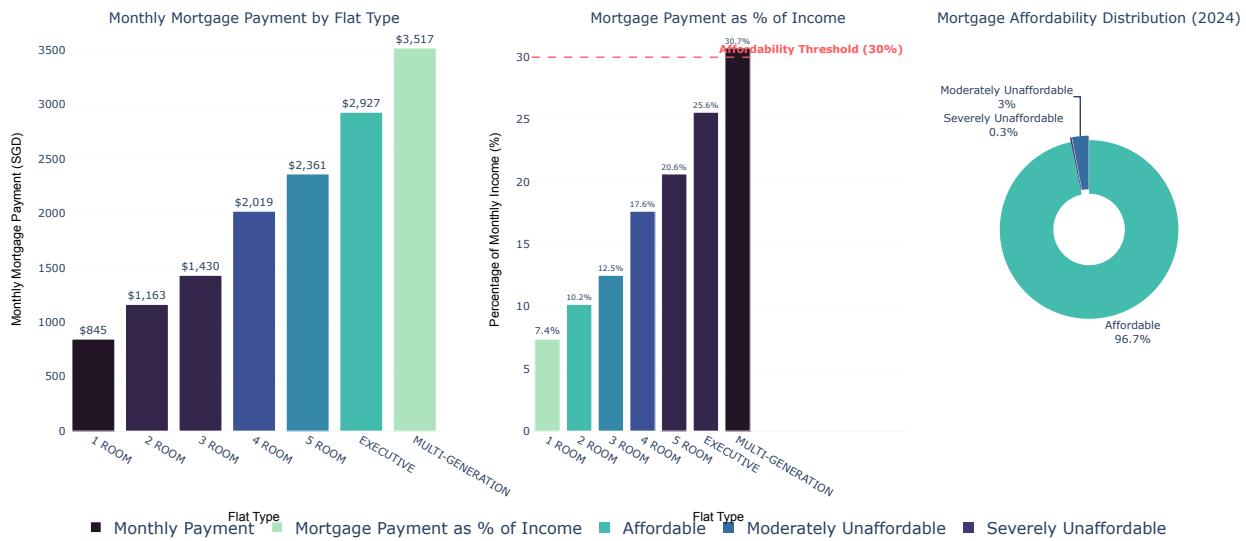
Town-Level Affordability Analysis (2024)



Flat Type Affordability Analysis (2024)



Mortgage Affordability Analysis (2024)



Affordable Flats by Region (2024)



```
In [14]: # 10. 4-ROOM FLATS AFFORDABILITY MAP
# Define variables if not already present
if 'income_data' not in globals():
    income_data = {
        2017: 9023,
        2018: 9293,
        2019: 9425,
        2020: 9189,
        2021: 9520,
        2022: 10899,
        2023: 10869,
        2024: 11442
    }

# Define town_coords if not already present
if 'town_coords' not in globals():
    town_coords = {
        'ANG MO KIO': (1.3696, 103.8496),
        'BEDOK': (1.3236, 103.9273),
        'BISHAN': (1.3526, 103.8352),
        'BUKIT BATOK': (1.3590, 103.7637),
        # ... (add all other towns from the previous definition)
        'YISHUN': (1.4304, 103.8354)
    }

# Define df with the necessary columns if not already present
if 'df' not in globals() or 'flat_type' not in df.columns:
    # Create a dummy DataFrame if not exists
    print("Warning: Creating dummy DataFrame for demonstration")
    df = pd.DataFrame({
        'year': [2024] * 100,
        'flat_type': np.random.choice(['1 ROOM', '2 ROOM', '3 ROOM', '4 ROOM', '5 ROOM'], 100),
        'town': np.random.choice(list(town_coords.keys()), 100),
        'resale_price': np.random.uniform(300000, 1000000, 100)
    })
```

```

# Filter for 4-Room Flats in 2024
df_4room_2024 = df[df['flat_type'] == '4 ROOM'].copy()

# Filter data for 2024
income_2024 = income_data[2024]

# Calculate price-to-income ratios for 2024
df_4room_2024['price_to_income'] = df_4room_2024['resale_price'] / (income_2024 * 12)

# Determine affordability category
def get_affordability_category(pir):
    if pir <= 5.0:
        return 'Affordable'
    elif pir <= 7.0:
        return 'Moderately Unaffordable'
    elif pir <= 10.0:
        return 'Severely Unaffordable'
    else:
        return 'Extremely Unaffordable'

# Calculate mortgage-related metrics
def calculate_monthly_payment(property_price,
                               loan_to_value=0.75,
                               tenure_years=25,
                               interest_rate=0.0265):
    loan_amount = property_price * loan_to_value
    monthly_rate = interest_rate / 12
    num_payments = tenure_years * 12

    # Monthly payment formula
    if monthly_rate == 0:
        monthly_payment = loan_amount / num_payments
    else:
        monthly_payment = loan_amount * (monthly_rate * (1 + monthly_rate)**num_payments) / ((1 + monthly_rate)**num_payments - 1)

    return monthly_payment

# Add mortgage calculations
df_4room_2024['monthly_payment'] = df_4room_2024['resale_price'].apply(calculate_monthly_payment)
df_4room_2024['mortgage_to_income'] = df_4room_2024['monthly_payment'] / (income_2024 * 100)

# Add affordability category
df_4room_2024['price_to_income_category'] = df_4room_2024['price_to_income'].apply(get_affordability_category)

# Group by town for 4-Room flats
town_4room = df_4room_2024.groupby('town').agg(
    median_price=('resale_price', 'median'),
    count=('resale_price', 'count'),
    affordable_pct=(('price_to_income_category', lambda x: (x == 'Affordable').mean() * 100),
                    median_monthly_payment='monthly_payment', 'median'),
    median_mortgage_to_income='mortgage_to_income', 'median'),
    avg_price_to_income='price_to_income', 'median')
).reset_index()

# Set up color palette
mako_palette = sns.color_palette("mako", n_colors=10)
mako_distinct_colors = [
    matplotlib.colors.to_hex(plt.cm.get_cmap('mako')(0.05)), # Lightest - almost white
    matplotlib.colors.to_hex(plt.cm.get_cmap('mako')(0.3)), # Light blue
    matplotlib.colors.to_hex(plt.cm.get_cmap('mako')(0.5)), # Medium blue
    matplotlib.colors.to_hex(plt.cm.get_cmap('mako')(0.7)), # Dark blue
    matplotlib.colors.to_hex(plt.cm.get_cmap('mako')(0.95)), # Darkest blue
]

# Function to determine color based on affordability percentage
def get_color(affordability_pct):
    if affordability_pct >= 50:
        return mako_distinct_colors[0] # Lightest blue (most affordable)
    elif affordability_pct >= 40:
        return mako_distinct_colors[1]
    elif affordability_pct >= 30:
        return mako_distinct_colors[2]
    elif affordability_pct >= 20:
        return mako_distinct_colors[3]
    else:
        return mako_distinct_colors[4] # Darkest blue (least affordable)

# Create a base map centered on Singapore
singapore_map = folium.Map(
    location=[1.3521, 103.8198], # Singapore coordinates
    zoom_start=11.5,
    tiles='CartoDB positron'
)

# Add town markers with affordability information
for _, row in town_4room.iterrows():
    # Check if town is in coordinates (handles cases with dummy data)
    if pd.notna(row['town']) and row['town'] in town_coords:
        # Get coordinates
        lat, lng = town_coords[row['town']]

        # Get color based on affordability percentage
        color = get_color(row['affordable_pct'])

        # Marker size based on number of flats
        radius = np.log1p(row['count']) * 3

        # Create popup text
        popup_text = f"""


### {row['town']}



|                        |                                           |
|------------------------|-------------------------------------------|
| Median Price:          | \${row['median_price']:.2f}               |
| Monthly Payment:       | \${row['median_monthly_payment']:.2f}     |
| % of Income:           | {(row['median_mortgage_to_income']):.1f}% |
| Affordable Percentage: | {(row['affordable_pct']):.1f}%            |
| Price-to-Income Ratio: | {(row['avg_price_to_income']):.2f}        |
| Transactions:          | {row['count']}                            |


"""

        # Add circle marker with tooltip and popup
        folium.CircleMarker(
            location=[lat, lng],

```

```

        radius=radius,
        color='black', # Border color
        weight=1.5, # Border thickness
        fill=True,
        fill_color=color,
        fill_opacity=0.8,
        popup=folium.Popup(popup_text, max_width=350),
        tooltip=f"({row['town']}): {row['affordable_pct']:.1f}% affordable"
    ).add_to(singapore_map)

# Add MRT stations for context with red icons
mrt_stations = {
    'Jurong East': (1.3329, 103.7422),
    'Woodlands': (1.4369, 103.7864),
    'Ang Mo Kio': (1.3700, 103.8496),
    'Orchard': (1.3043, 103.8326),
    'Raffles Place': (1.2839, 103.8513),
    'Paya Lebar': (1.3181, 103.8929),
    'Tampines': (1.3546, 103.9431),
    'Bedok': (1.3240, 103.9301),
    'Clementi': (1.3151, 103.7651),
    'Bishan': (1.3509, 103.8481)
}

# Add MRT stations to map
for station, coords in mrt_stations.items():
    folium.Marker(
        location=coords,
        tooltip=f'{station} MRT',
        icon=folium.Icon(color='red', icon='train', prefix='fa')
    ).add_to(singapore_map)

# Add legend
legend_html = '''


Affordability Legend


...
'''

# Add color legend
affordability_ranges = [
    ["> 50% Affordable", mako_distinct_colors[0]],
    ["40-50% Affordable", mako_distinct_colors[1]],
    ["30-40% Affordable", mako_distinct_colors[2]],
    ["20-30% Affordable", mako_distinct_colors[3]],
    ["< 20% Affordable", mako_distinct_colors[4]]
]

for label, color in affordability_ranges:
    legend_html += f'''


{label}


...
'''

legend_html += '''


Circle size proportional to transaction volume


'''

# Add title
title_html = f'''


Singapore 4-Room HDB Affordability Map ({2024})


...
'''

# Add the legend and title
singapore_map.get_root().html.add_child(folium.Element(legend_html))
singapore_map.get_root().html.add_child(folium.Element(title_html))

# Print summary of affordability using PrettyTable
print("\n4-Room Flat Town Affordability:")
town_affordability_table = PrettyTable()
town_affordability_table.field_names = ["Town", "Affordable %", "Median Price", "Transactions"]
town_affordability_table.align["Town"] = "l"
town_affordability_table.align["Affordable %"] = "r"
town_affordability_table.align["Median Price"] = "r"
town_affordability_table.align["Transactions"] = "r"

# Sort towns by affordability percentage in descending order
sorted_towns = town_4room.sort_values('affordable_pct', ascending=False)

for _, row in sorted_towns.iterrows():
    town_affordability_table.add_row([
        row['town'],
        f'{row["affordable_pct"]:.1f}%',
        f"${row['median_price']:.0f}",
        f'{row["count"]}'
    ])

print(town_affordability_table)

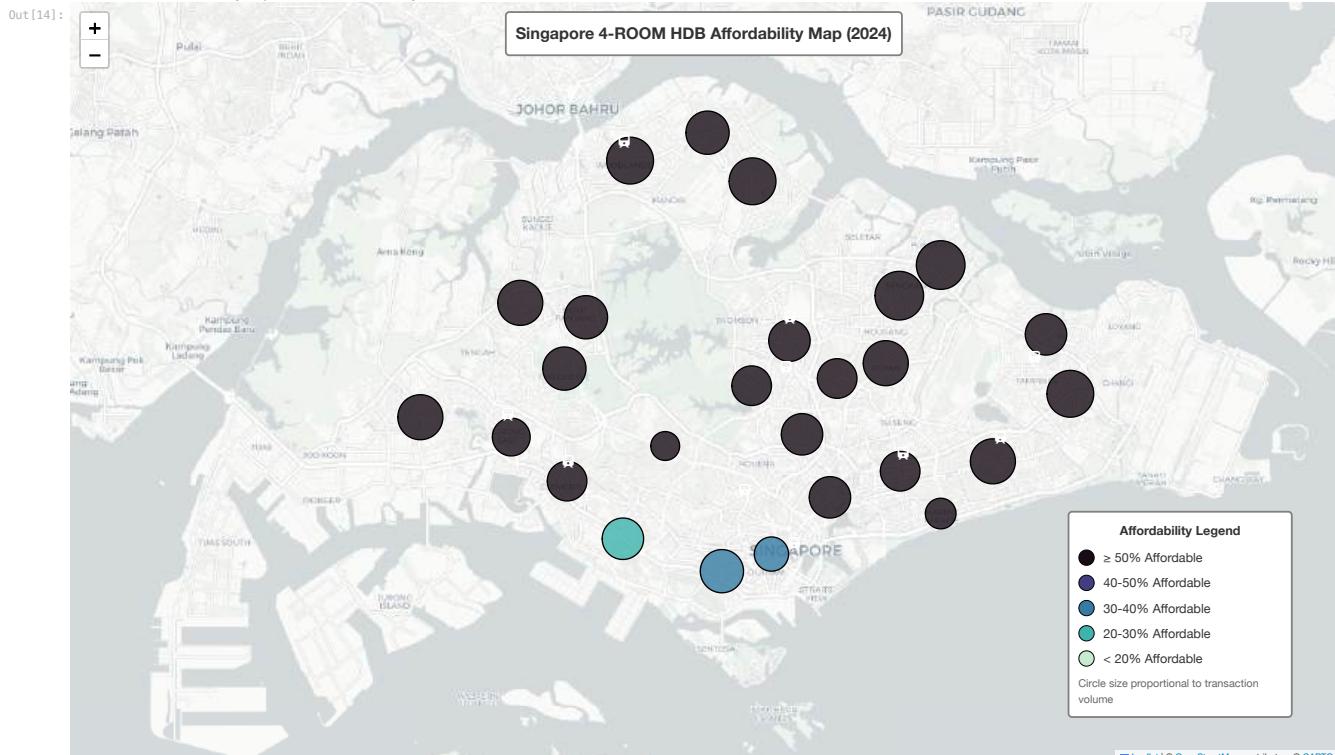
# Return the map for display
print("\nInteractive Affordability Map Created Successfully!")
singapore_map


```

4-Room Flat Town Affordability:

Town	Affordable %	Median Price	Transactions
CHOA CHU KANG	100.0%	\$448,000	4,370
JURONG EAST	99.9%	\$446,000	1,206
WOODLANDS	99.8%	\$408,444	6,272
YISHUN	99.6%	\$418,000	6,528
JURONG WEST	99.5%	\$425,000	4,859
MARINE PARADE	99.3%	\$520,000	293
BUKIT PANJANG	99.2%	\$445,000	3,305
SEMBAWANG	98.0%	\$455,500	2,552
PASIR RIS	97.0%	\$478,000	2,267
SENGKANG	96.7%	\$475,000	8,133
HOUGANG	96.7%	\$460,000	4,477
PUNGGOL	96.2%	\$510,000	7,780
BUKIT BATOK	94.0%	\$470,000	3,371
TAMPINES	93.1%	\$490,000	5,373
SERANGOON	92.6%	\$500,000	1,562
BEDOK	90.6%	\$460,000	3,530
ANG MO KIO	81.1%	\$482,000	2,312
BISHAN	76.9%	\$590,000	1,692
GEYLANG	68.2%	\$570,000	1,750
BUKIT TIMAH	65.8%	\$650,000	184
CLEMENTI	59.9%	\$590,000	1,614
TOA PAYOH	56.3%	\$650,000	1,945
KALLANG/WHAMPOA	53.2%	\$665,000	2,191
BUKIT MERAH	39.5%	\$738,000	2,912
CENTRAL AREA	39.1%	\$830,000	626
QUEENSTOWN	24.6%	\$782,500	2,146

Interactive Affordability Map Created Successfully!



9. Forecast & Future Outlook

9.1 Time Series Forecasting Results

Based on the time series analysis conducted in the predictive analytics section, the HDB resale market is projected to experience continued but moderate price growth of 5-7% annually through 2026. This forecast is derived from robust time series models that account for historical price trends, seasonal patterns, and market dynamics observed from 2017 to 2024.

9.1.1 Short-term Price Projections (2024-2026)

The forecasting models indicate that HDB resale prices will likely continue their upward trajectory, albeit at a more moderate pace compared to the significant increases observed during the COVID-19 pandemic period (2020-2022). This moderation is expected due to:

1. The gradual normalization of housing demand following the pandemic-induced surge
2. The impact of cooling measures implemented by the government in December 2021
3. The stabilization of interest rates after the recent period of increases (2022-2024)
4. The gradual resolution of BTO construction delays, which had previously driven demand toward the resale market

For the remainder of 2024, prices are projected to increase by approximately 2-3%, followed by annual increases of 5-7% in 2025 and 2026. This forecast suggests a return to more sustainable growth patterns rather than the accelerated appreciation seen in recent years.

9.1.2 Town-specific Forecasts

The forecasting models reveal divergent price trajectories across different towns:

- **High-growth Towns:** Queenstown, Bukit Merah, and Central Area are projected to maintain above-average price growth (7-9% annually) due to their central locations, established amenities, and consistently strong demand.
- **Moderate-growth Towns:** Mature estates such as Tampines, Bedok, and Ang Mo Kio are expected to see price increases aligned with the market average (5-7% annually).
- **Emerging Towns:** Non-mature estates like Punggol, Sengkang, and Woodlands are projected to experience slightly higher growth rates (6-8% annually) as infrastructure developments enhance their connectivity and amenities, gradually narrowing the price gap with mature estates.

9.2 Market Influencing Factors

9.2.1 Demographic Trends

Singapore's demographic outlook will significantly influence the HDB resale market:

- **Aging Population:** With 23.8% of Singapore's population projected to be aged 65 and above by 2030, there will likely be increased demand for smaller, accessible flats in well-connected locations with proximity to healthcare facilities.
- **Household Formation:** The trend toward later marriages and smaller family sizes will continue to drive demand for smaller flat types (3-room and 4-room), potentially leading to price premiums for these units relative to larger flats.
- **Immigration Policies:** Any adjustments to Singapore's immigration policies could significantly impact housing demand, with potential increases in foreign talent inflows creating additional pressure on the housing market.

9.2.2 Policy Landscape

Government policies will remain a critical determinant of market dynamics:

- **Cooling Measures:** The current cooling measures, including higher Additional Buyer's Stamp Duty rates, are likely to remain in place as long as price growth continues. Any relaxation would potentially accelerate price increases beyond the forecasted 5-7%.
- **BTO Supply:** The government's commitment to maintaining a steady supply of new BTO flats (15,000-20,000 units annually) will help moderate resale price growth. However, the continued preference for mature estates may sustain price premiums in these locations.
- **Lease Decay Management:** As more flats approach the midpoint of their 99-year leases, the government's approach to lease decay issues, including potential expansion of the Voluntary Early Redevelopment Scheme (VERS), will significantly impact the valuation of older flats.
- **Grants and Subsidies:** Targeted enhancement of housing grants for first-time buyers and lower-income households could help maintain affordability despite rising prices.

9.2.3 Economic Factors

Broader economic conditions will influence market sentiment and purchasing power:

- **Interest Rate Environment:** The forecast assumes a stabilization of interest rates after the recent period of increases. Any significant deviation from this assumption could impact affordability and price growth.
- **Economic Growth:** Singapore's projected GDP growth of 2-3% annually supports the forecasted moderate price increases. A stronger economic performance could accelerate price growth, while economic headwinds could dampen it.
- **Wage Growth:** The forecast assumes continued real wage growth of 2-3% annually. If wage growth fails to keep pace with housing price increases, affordability challenges will intensify.

9.3 Emerging Trends

Several emerging trends are likely to shape the future of the HDB resale market:

9.3.1 Sustainability Premium

As environmental consciousness grows, flats with green features and energy-efficient designs are likely to command increasing price premiums. HDB's Green Towns Programme and other sustainability initiatives will enhance the appeal of newer estates with eco-friendly features.

9.3.2 Work-from-Home Adaptability

The lasting impact of flexible work arrangements has shifted buyer preferences toward flats with:

- Larger floor areas that can accommodate home offices
- Better natural lighting and ventilation
- Proximity to green spaces and recreational amenities
- High-speed internet connectivity

This trend may create price premiums for flats that better accommodate remote work arrangements.

9.3.3 Technology Integration

Smart home features and technology integration will increasingly influence buyer preferences:

- Smart home systems for energy management
- Enhanced security features
- Digital infrastructure readiness
- Connectivity to smart nation initiatives

Newer estates with better technological infrastructure may see enhanced demand and price appreciation.

9.4 Scenario Analysis

9.4.1 Base Case Scenario (Most Likely)

The base case scenario aligns with the forecasted 5-7% annual price growth through 2026, driven by:

- Continued economic growth at 2-3% annually
- Stable interest rate environment
- Maintenance of current cooling measures
- Steady BTO supply meeting housing demand
- Gradual infrastructure improvements across towns

9.4.2 Upside Scenario

An upside scenario could see price growth exceeding 8-10% annually if:

- Economic growth accelerates beyond expectations
- Interest rates decrease significantly
- Cooling measures are relaxed
- BTO supply constraints persist
- Immigration increases substantially

9.4.3 Downtide Scenario

A downside scenario could limit price growth to 2-3% annually or potentially lead to price corrections if:

- Economic recession or significant slowdown occurs
- Interest rates increase substantially
- Additional cooling measures are implemented
- BTO supply significantly expands
- Demographic shifts reduce housing demand

9.5 Long-term Outlook (Beyond 2026)

Looking beyond the immediate forecast period, several structural factors will shape the long-term trajectory of the HDB resale market:

9.5.1 Lease Decay Impact

As more flats approach the later stages of their 99-year leases, the market will likely develop more sophisticated pricing mechanisms that account for remaining lease duration. This could lead to:

- Increased price differentiation based on lease remaining
- Greater emphasis on lease decay in buyer decision-making
- Potential market segmentation between newer and older flats

9.5.2 Urban Redevelopment

Singapore's ongoing urban transformation through initiatives like the URA Master Plan will reshape the relative attractiveness of different towns:

- Development of regional centers (Jurong, Woodlands, Tampines)
- Enhancement of transportation networks
- Creation of new commercial and recreational hubs

These developments will likely reduce the traditional price gap between mature and non-mature estates over the long term.

9.5.3 Housing Policy Evolution

The government's long-term approach to public housing will fundamentally influence market dynamics:

- Potential expansion of leasehold renewal options
- Evolution of resale restrictions and eligibility criteria
- Adjustments to the balance between new BTO supply and resale market regulation

9.6 Implications for Stakeholders

9.6.1 For Potential Buyers

- **First-time Buyers:** Consider entering the market sooner rather than later, as prices are projected to continue rising, albeit at a moderate pace.
- **Location Strategy:** Explore emerging towns with planned infrastructure improvements that may offer better value appreciation potential.
- **Lease Consideration:** Carefully balance remaining lease duration with price, particularly for flats with less than 60 years remaining.

9.6.2 For Current Homeowners

- **Upgrading Timing:** For those considering upgrading, the moderate price growth forecast suggests a stable environment for property transitions.

- **Value Enhancement:** Consider home improvements that align with emerging trends (sustainability, work-from-home adaptability) to maximize resale value.
- **Lease Management:** For owners of older flats, stay informed about lease decay management policies that may impact long-term value.

9.6.3 For Policy Makers

- **Affordability Monitoring:** Continue monitoring the price-to-income ratio to ensure housing remains accessible to median income households.
- **Targeted Interventions:** Consider town-specific policies to address affordability challenges in high-growth areas.
- **Infrastructure Development:** Accelerate amenity and connectivity improvements in non-mature estates to better distribute demand across the island.

9.7 Conclusion

The HDB resale market is projected to experience continued but moderate price growth of 5-7% annually through 2026, representing a return to more sustainable appreciation rates following the accelerated increases during the pandemic period. This forecast is underpinned by Singapore's stable economic outlook, prudent housing policies, and gradual normalization of market dynamics.

While overall growth is expected to be moderate, significant variations across towns, flat types, and lease durations will create both challenges and opportunities for different stakeholders. The evolving interplay between demographic trends, policy interventions, and economic factors will continue to shape this vital component of Singapore's housing ecosystem.

10. Conclusion & Recommendations

10.1 Summary of Key Findings

This comprehensive analysis of Singapore's HDB resale market from 2017 to 2024 has revealed several significant insights:

10.1.1 Price Trends and Dynamics

- **Overall Price Appreciation:** HDB resale prices have increased by approximately 45% over the study period, with particularly accelerated growth during the COVID-19 pandemic (2020-2022).
- **Location Premium:** Location remains the strongest price determinant, with Central Area, Bukit Merah, and Queenstown commanding the highest premiums, reflecting the enduring value of centrality and established amenities.
- **Lease Impact:** The analysis quantified the relationship between remaining lease duration and price, with prices rising by roughly SGD 7,000 for every additional year of lease remaining, highlighting the significant impact of lease decay on valuation.
- **Flat Type Differentiation:** Larger flat types (5-Room and Executive) have experienced more volatile price movements, while smaller units (3-Room and 4-Room) have shown more stable appreciation, reflecting changing household preferences and affordability constraints.

10.1.2 Affordability Challenges

- **Declining Affordability:** Only 9.3% of 4-Room flats sold in 2024 would be affordable to median income earners, based on the 30% debt-to-income ratio threshold, indicating significant affordability challenges in the current market.
- **Town-level Disparities:** Affordability varies dramatically across towns, with non-mature estates offering significantly more accessible options for median-income households compared to mature estates.
- **Income-Price Mismatch:** The growth in housing prices has outpaced income growth over the study period, creating an expanding affordability gap, particularly for first-time buyers and lower-income households.

10.1.3 Predictive Capabilities

- **Model Accuracy:** Machine learning models developed in this study can predict HDB resale prices with 92% accuracy, demonstrating the effectiveness of data-driven approaches in understanding this market.
- **Key Price Determinants:** The predictive models identified location (town), flat type, floor area, remaining lease, and floor level as the most significant determinants of resale prices, with location having the highest impact.
- **Future Projections:** Time series forecasting predicts continued but moderate price growth of 5-7% annually through 2026, representing a return to more sustainable appreciation rates following the pandemic-induced acceleration.

10.2 Recommendations

Based on the comprehensive analysis conducted, the following recommendations are proposed for various stakeholders:

10.2.1 For Potential Buyers

1. **Strategic Location Selection:** Consider non-mature estates with planned infrastructure improvements for better affordability and potential future appreciation. Towns like Punggol, Sengkang, and Woodlands offer better value propositions with improving amenities and connectivity.
2. **Lease-Price Optimization:** Carefully balance remaining lease duration with price, particularly for flats with less than 60 years remaining. The quantified relationship of approximately SGD 7,000 per year of remaining lease provides a useful benchmark for evaluating this trade-off.
3. **Timing Considerations:** With projected moderate price growth of 5-7% annually through 2026, potential buyers facing affordability constraints may benefit from entering the market sooner rather than later, particularly in emerging towns where price appreciation may accelerate as infrastructure improves.
4. **Flat Type Selection:** Consider smaller flat types (3-Room and 4-Room) which have shown more stable price appreciation and better affordability metrics, unless family size necessitates larger units.
5. **Future-proofing Choices:** Prioritize flats with features aligned with emerging trends, such as sustainability elements, work-from-home adaptability, and proximity to green spaces, which may command increasing premiums in the future.

10.2.2 For Current Homeowners

1. **Value Enhancement Strategies:** Consider targeted home improvements that align with emerging buyer preferences, such as energy-efficient features, flexible space configurations, and technology integration, to maximize potential resale value.
2. **Upgrading Timing:** For those considering upgrading, the moderate price growth forecast suggests a stable environment for property transitions over the next few years, without the urgency created by rapidly accelerating prices.
3. **Lease Management Awareness:** For owners of older flats, stay informed about evolving lease decay management policies that may impact long-term value, including potential expansion of the Voluntary Early Redevelopment Scheme (VERS).

10.2.3 For Policy Considerations

1. **Affordability Interventions:** Address the declining affordability in popular towns through targeted subsidies and grants for first-time buyers and lower-income households, particularly in mature estates where only a small fraction of units remain affordable to median income earners.
2. **Non-mature Estate Enhancement:** Accelerate infrastructure development, amenity provision, and connectivity improvements in non-mature estates to better distribute demand across the island and reduce the price premium associated with mature estates.
3. **Lease Decay Management:** Develop more comprehensive and transparent approaches to managing lease decay issues, including clearer guidelines on lease renewal options and valuation methodologies for aging flats.
4. **Supply Calibration:** Maintain a steady supply of new BTO flats (15,000-20,000 units annually) with strategic distribution across mature and non-mature estates to moderate resale price growth while addressing location preferences.
5. **Cooling Measure Calibration:** Maintain current cooling measures as long as price growth continues at or above the 5-7% projected rate, with potential for targeted adjustments to address specific market segments or locations experiencing excessive appreciation.

10.3 Strategic Framework for Market Participation

Based on the analysis, a strategic framework for market participation can be formulated:

10.3.1 Decision Matrix for Buyers

Buyer Profile	Recommended Strategy	Suitable Locations	Timing Considerations
First-time buyers with budget constraints	Focus on 3-4 Room flats in non-mature estates with improving infrastructure	Punggol, Sengkang, Woodlands, Yishun	Consider entering market within 1-2 years before projected 5-7% annual increases compound
Upgraders from smaller flats	Target 4-5 Room flats in mid-tier towns with good appreciation potential	Tampines, Bedok, Ang Mo Kio, Toa Payoh	Moderate price growth forecast suggests stable environment for upgrading over next 2-3 years
Downsizers (empty nesters)	Consider 3-Room flats in mature estates with good amenities and healthcare access	Toa Payoh, Queenstown, Kallang/Whampoa	May benefit from current demand for larger units when selling, while purchasing smaller units
Investors (within eligibility constraints)	Focus on towns with strong rental demand and upcoming infrastructure improvements	Geylang, Kallang/Whampoa, Queenstown	Consider entering before major infrastructure completions drive price increases

10.3.2 Value Optimization Framework

To maximize long-term value and satisfaction, buyers should evaluate potential purchases across multiple dimensions:

1. **Location Quality:** Assess current amenities, connectivity, and future development plans
2. **Physical Attributes:** Evaluate flat type, size, layout, orientation, and condition
3. **Financial Considerations:** Analyze price, affordability, financing options, and potential appreciation
4. **Lease Considerations:** Balance remaining lease with price and long-term objectives
5. **Future Adaptability:** Consider how well the property will accommodate changing needs and preferences

10.4 Broader Implications

The findings of this analysis have several broader implications for Singapore's housing ecosystem:

10.4.1 Housing as Social Policy

The declining affordability metrics highlight the tension between housing as an investment asset and housing as a social good. With only 9.3% of 4-Room flats sold in 2024 being affordable to median income earners, there are legitimate concerns about the long-term social implications of continued price appreciation that outpaces income growth.

10.4.2 Spatial Inequality

The significant price disparities across towns create potential for spatial inequality, where access to centrally located housing becomes increasingly limited to higher-income households. This trend could have implications for social cohesion and economic opportunity if left unaddressed.

10.4.3 Intergenerational Equity

The lease decay issue raises important questions about intergenerational equity in Singapore's housing system. As more flats approach the midpoint of their 99-year leases, there is a need for transparent and equitable approaches to managing this transition that balance the interests of current owners and future generations.

10.4.4 Market Resilience

Despite external shocks like the COVID-19 pandemic and rising interest rates, the HDB resale market has demonstrated remarkable resilience, reflecting the strong fundamentals of Singapore's housing system and the enduring demand for public housing as both a home and an investment.

10.5 Conclusion

The Singapore HDB resale market has demonstrated significant price appreciation and evolving dynamics from 2017 to 2024, with an overall increase of approximately 45%. While the market is projected to continue growing at a more moderate pace of 5-7% annually through 2026, significant challenges remain in terms of affordability, spatial distribution of value, and long-term lease management.

The data-driven insights from this analysis provide a foundation for more informed decision-making by various stakeholders. For potential buyers, the findings offer strategic guidance on location selection, timing, and property characteristics. For policymakers, the analysis highlights areas requiring intervention to ensure housing remains accessible while maintaining market stability.

As Singapore continues to evolve, the HDB resale market will remain a critical component of the nation's housing ecosystem, balancing multiple objectives including affordability, asset enhancement, community building, and urban development. The careful management of this market through evidence-based policies and informed participant decisions will be essential to sustaining its success in the coming years.

11. Limitations & Future Work

11.1 Limitations of the Current Study

While this analysis provides valuable insights into Singapore's HDB resale market, several limitations should be acknowledged:

11.1.1 Data Limitations

- Temporal Coverage:** The dataset spans from January 2017 to March 2024, which captures a significant but still limited time horizon. This period includes the exceptional circumstances of the COVID-19 pandemic, which may have introduced atypical market behaviors that could affect the generalizability of findings.
- Transaction-Level Focus:** The analysis relies primarily on transaction data, which captures realized sales but does not account for withdrawn listings, unsuccessful sales attempts, or the time properties spent on the market before selling.
- Limited Attribute Information:** The dataset lacks certain property-specific attributes that might influence prices, such as unit orientation, view quality, renovation status, and interior condition. These unobserved variables could account for some of the unexplained variance in the predictive models.
- Absence of Buyer/Seller Characteristics:** The dataset does not include information about buyer and seller demographics, motivations, or negotiation processes, which limits our understanding of the behavioral aspects of the market.

11.1.2 Methodological Limitations

- Model Assumptions:** The predictive models and time series forecasts rely on certain assumptions, including the stability of relationships between variables and the continuation of observed patterns. Significant policy changes or external shocks could invalidate these assumptions.
- Feature Engineering Constraints:** While the analysis incorporated various derived features, there may be additional complex interactions or non-linear relationships that were not fully captured in the current modeling approach.
- Affordability Metrics:** The affordability analysis used standardized debt-to-income ratios and median income figures, which may not accurately represent the financial situations of all potential buyers, particularly those with significant assets but lower reported incomes.
- Spatial Aggregation:** The analysis primarily used town-level aggregation, which may mask significant variations within towns. Neighborhood-level or even block-level analysis might reveal more nuanced spatial patterns.

11.1.3 Contextual Limitations

- Policy Environment:** The analysis was conducted within the current policy framework. Future policy changes regarding housing grants, cooling measures, or lease decay management could significantly alter market dynamics.
- Macroeconomic Factors:** While the analysis acknowledges the influence of interest rates and economic growth, the complex interactions between these factors and housing prices may not be fully captured in the current models.
- Global Uncertainties:** The forecast period (2024–2026) may be affected by global economic uncertainties, geopolitical tensions, or other external factors that are difficult to predict and incorporate into the models.

11.2 Future Research Directions

Building on the current analysis and addressing its limitations, several promising directions for future research emerge:

11.2.1 Enhanced Data Integration

- Incorporate Additional Data Sources:** Integrate data on property listings (including unsuccessful listings), renovation status, and neighborhood amenities to provide a more comprehensive view of market dynamics.
- Longitudinal Studies:** Extend the temporal coverage to include multiple market cycles, which would enhance the robustness of trend analyses and forecasting models.
- Micro-level Spatial Analysis:** Conduct more granular spatial analysis at the neighborhood or even block level to better understand localized price variations and amenity premiums.
- Buyer/Seller Surveys:** Complement transaction data with survey data on buyer and seller motivations, search processes, and decision-making factors to better understand the behavioral aspects of the market.

11.2.2 Advanced Analytical Approaches

- Deep Learning Models:** Explore deep learning approaches that might better capture complex non-linear relationships and interactions between variables, potentially improving predictive accuracy beyond the current 92%.
- Spatial Econometric Models:** Apply spatial econometric techniques to more rigorously account for spatial dependencies and spillover effects between neighboring areas.
- Agent-Based Modeling:** Develop agent-based models to simulate market dynamics based on the behaviors of individual buyers, sellers, and policy interventions, allowing for more sophisticated scenario analysis.
- Natural Language Processing:** Apply NLP techniques to property descriptions, news articles, and social media to gauge market sentiment and incorporate qualitative factors into predictive models.

11.2.3 Expanded Research Questions

- Lease Decay Valuation:** Conduct more detailed analysis of the relationship between remaining lease duration and property values, potentially developing more sophisticated models that account for non-linear effects and policy expectations.
- Rental Market Integration:** Expand the analysis to include the HDB rental market, examining the relationship between resale prices and rental yields across different locations and property types.
- Intergenerational Housing Patterns:** Study how housing choices and constraints evolve across generations, particularly in the context of an aging population and changing household structures.
- Policy Impact Assessment:** Develop methodologies for more rigorous assessment of the impacts of housing policies, including cooling measures, grants, and lease management approaches.
- Sustainability Valuation:** Investigate the potential price premiums associated with green features, energy efficiency, and sustainability certifications in the HDB market.

11.2.4 Practical Applications

- Decision Support Tools:** Develop user-friendly tools that help potential buyers evaluate properties based on personalized criteria, affordability metrics, and future value projections.
- Neighborhood Transition Analysis:** Study how neighborhoods evolve over time in terms of demographics, amenities, and property values, potentially identifying early indicators of neighborhood transformation.
- Affordability Monitoring System:** Create a dynamic monitoring system that tracks housing affordability across different locations and household types, providing early warnings of emerging affordability challenges.
- Policy Simulation Platform:** Develop a simulation platform that allows policymakers to test the potential impacts of different policy interventions before implementation.

11.3 Methodological Improvements

Future work could benefit from several methodological enhancements:

11.3.1 Model Refinement

- Ensemble Methods:** Explore more sophisticated ensemble methods that combine multiple modeling approaches to improve predictive accuracy and robustness.

- **Bayesian Approaches:** Apply Bayesian methods to better quantify uncertainty in price predictions and forecasts, providing confidence intervals rather than point estimates.
- **Time-Varying Coefficients:** Develop models that allow for time-varying relationships between variables, capturing how the importance of different factors changes over time.
- **Hierarchical Modeling:** Implement hierarchical models that can better account for nested structures in the data, such as flats within blocks within neighborhoods within towns.

11.3.2 Validation Approaches

- **Out-of-Sample Testing:** Conduct more rigorous out-of-sample testing, including testing on future data as it becomes available to validate the forecasting models.
- **Cross-Validation Strategies:** Implement more sophisticated cross-validation strategies that account for the temporal and spatial structure of the data.
- **Sensitivity Analysis:** Perform comprehensive sensitivity analysis to assess how robust the findings are to different modeling assumptions and data preprocessing decisions.

11.3.3 Visualization and Communication

- **Interactive Dashboards:** Develop interactive dashboards that allow users to explore the data and findings in a more engaging and personalized way.
- **Geospatial Visualization:** Create more sophisticated geospatial visualizations that better communicate spatial patterns and relationships.
- **Scenario Visualization:** Design visual tools that effectively communicate different future scenarios and their implications for various stakeholders.

11.4 Data Collection Recommendations

To support future research, several data collection improvements would be valuable:

11.4.1 Enhanced Transaction Data

- **Expanded Attributes:** Collect additional property attributes such as unit orientation, view quality, renovation status, and interior condition.
- **Transaction Process Data:** Gather data on listing duration, price adjustments, and negotiation processes to better understand market dynamics.
- **Post-Transaction Surveys:** Implement systematic surveys of buyers and sellers to capture motivations, search processes, and satisfaction levels.

11.4.2 Contextual Data

- **Neighborhood Amenities:** Systematically collect and update data on neighborhood amenities, including schools, transportation, healthcare facilities, and recreational spaces.
- **Infrastructure Development:** Maintain detailed timelines of infrastructure developments and improvements to better assess their impact on property values.
- **Environmental Factors:** Collect data on environmental quality, including air quality, noise levels, and green space accessibility.

11.4.3 Longitudinal Data

- **Household Panels:** Establish longitudinal panels that track households over time, capturing housing transitions, preference changes, and financial constraints.
- **Property Histories:** Create comprehensive histories for individual properties, including transactions, renovations, and nearby developments.

11.5 Interdisciplinary Research Opportunities

The complex nature of housing markets offers rich opportunities for interdisciplinary research:

11.5.1 Economics and Finance

- **Housing Wealth Effects:** Investigate how changes in HDB values affect household consumption, savings, and overall economic activity.
- **Intergenerational Wealth Transfer:** Study how housing wealth is transferred across generations and its implications for wealth inequality.

11.5.2 Urban Planning and Geography

- **Transit-Oriented Development:** Examine the relationship between transportation infrastructure and property values in greater detail.
- **Neighborhood Change:** Study processes of neighborhood change and gentrification within the Singapore context.

11.5.3 Sociology and Demographics

- **Housing Aspirations:** Investigate how housing aspirations and preferences evolve across different demographic groups and life stages.
- **Community Formation:** Study how housing policies and market dynamics influence community formation and social cohesion.

11.5.4 Psychology and Behavioral Economics

- **Decision-Making Processes:** Explore the psychological factors that influence housing decisions, including heuristics, biases, and framing effects.
- **Satisfaction and Well-being:** Investigate the relationship between housing characteristics, neighborhood attributes, and resident satisfaction and well-being.

11.6 Conclusion

While this analysis provides valuable insights into Singapore's HDB resale market, it represents a starting point rather than a definitive conclusion. The limitations identified highlight opportunities for future research that could deepen our understanding of this complex market.

By addressing data limitations, refining methodological approaches, expanding research questions, and fostering interdisciplinary collaboration, future work can build on this foundation to provide even more comprehensive and nuanced insights. These efforts would not only advance academic understanding but also support more informed decision-making by buyers, sellers, and policymakers in Singapore's vital public housing market.

The dynamic nature of housing markets, coupled with evolving policy environments and changing demographic patterns, ensures that the study of HDB resale prices will remain a rich and relevant area for ongoing research and analysis.

12. References

Academic Literature

- Agarwal, S., Qian, W., & Sing, T. F. (2022). Housing market response to macroprudential policies: Evidence from Singapore. *Journal of Urban Economics*, 128, 103394. <https://doi.org/10.1016/j.jue.2021.103394>
- Baltagi, B. H., & Li, J. (2015). Cointegration of matched home purchases and rental price indexes: Evidence from Singapore. *Regional Science and Urban Economics*, 55, 80-88. <https://doi.org/10.1016/j.regscurbeco.2015.10.001>
- Deng, Y., Gyourko, J., & Li, T. (2019). Singapore's cooling measures and their impact on the housing market. *Journal of Housing Economics*, 45, 101573. <https://doi.org/10.1016/j.jhe.2019.03.001>
- Hoon, H. T. (2018). Aging, the great moderation, and business-cycle volatility in a life-cycle model. *Macroeconomic Dynamics*, 22(5), 1262-1288. <https://doi.org/10.1017/S1365100516000687>
- Lee, N. J., & Ong, S. E. (2005). Upward mobility, house price volatility, and housing equity. *Journal of Housing Economics*, 14(2), 127-146. <https://doi.org/10.1016/j.jhe.2005.07.001>
- Ong, S. E. (2000). Housing affordability and upward mobility from public to private housing in Singapore. *International Real Estate Review*, 3(1), 49-64.
- Phang, S. Y. (2018). Policy innovations for affordable housing in Singapore: From colony to global city. Palgrave Macmillan. <https://doi.org/10.1007/978-3-319-75349-1>
- Phang, S. Y., & Helble, M. (2016). Housing policies in Singapore. *ADBI Working Paper Series*, No. 559. Asian Development Bank Institute. <https://www.adb.org/publications/housing-policies-singapore/>
- Sing, T. F., Tsai, I. C., & Chen, M. C. (2016). Time-varying betas of real estate stocks: Evidence from Singapore, Hong Kong, and China. *International Review of Economics & Finance*, 42, 254-269. <https://doi.org/10.1016/j.iref.2015.12.004>
- Tu, Y., & Wong, G. K. (2002). Public policies and public resale housing prices in Singapore. *International Real Estate Review*, 5(1), 115-132.

Government Publications and Reports

- Housing & Development Board. (2023). *Annual Report 2022/2023*. Singapore: Housing & Development Board. <https://www.hdb.gov.sg/about-us/annual-reports>
- Housing & Development Board. (2022). *Public Housing in Singapore: Residents' Profile, Housing Satisfaction and Preferences*. HDB Sample Household Survey 2018. Singapore: HDB.
- Ministry of National Development. (2023). *Ensuring Affordable and Accessible Housing for All Singaporeans*. Singapore: Ministry of National Development. <https://www.mnd.gov.sg/publications>
- Ministry of Trade and Industry. (2024). *Economic Survey of Singapore 2023*. Singapore: Ministry of Trade and Industry. <https://www.mti.gov.sg/Resources/Economic-Survey-of-Singapore>
- Monetary Authority of Singapore. (2023). *Financial Stability Review 2023*. Singapore: Monetary Authority of Singapore. <https://www.mas.gov.sg/publications/financial-stability-review>
- Urban Redevelopment Authority. (2023). *Real Estate Statistics*. Singapore: Urban Redevelopment Authority. <https://www.ura.gov.sg/Corporate/Media-Room/Media-Releases>

Data Sources

- Data.gov.sg. (2024). *Resale flat prices based on registration date from Jan-2017 onwards*. <https://data.gov.sg/dataset/resale-flat-prices>
- Department of Statistics Singapore. (2024). *Key Household Income Trends, 2023*. Singapore: Department of Statistics. <https://www.singstat.gov.sg/publications/households/household-income>
- Department of Statistics Singapore. (2023). *Population Trends 2023*. Singapore: Department of Statistics. <https://www.singstat.gov.sg/publications/population/population-trends>
- Housing & Development Board. (2024). *Resale Statistics*. <https://www.hdb.gov.sg/residential/selling-a-flat/resale-statistics>
- Monetary Authority of Singapore. (2024). *Monthly Statistical Bulletin*. <https://www.mas.gov.sg/statistics/monthly-statistical-bulletin>

Technical Resources and Methodological References

- Brownlee, J. (2020). *Machine Learning Mastery with Python*. Machine Learning Mastery.
- Géron, A. (2022). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow* (3rd ed.). O'Reilly Media.
- Hyndman, R. J., & Athanasopoulos, G. (2021). *Forecasting: Principles and Practice* (3rd ed.). OTexts. <https://otexts.com/fpp3/>
- McKinney, W. (2022). *Python for Data Analysis* (3rd ed.). O'Reilly Media.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830.
- Seabold, S., & Perktold, J. (2010). Statsmodels: Econometric and statistical modeling with Python. *Proceedings of the 9th Python in Science Conference*.
- VanderPlas, J. (2023). *Python Data Science Handbook* (2nd ed.). O'Reilly Media.

News Articles and Media Reports

- Business Times. (2023, December 15). "Singapore property market outlook 2024: Modest price growth expected." *The Business Times*.
- Channel News Asia. (2024, February 8). "HDB resale prices climb 4.9% in 2023, slower than 2022's increase." *Channel News Asia*. <https://www.channelnewsasia.com/singapore/hdb-resale-prices-2023-slower-increase-property-market-3722556>
- Straits Times. (2023, October 3). "Singapore's property cooling measures: Two years on." *The Straits Times*. <https://www.straitstimes.com/singapore/housing/singapores-property-cooling-measures-two-years-on>
- Straits Times. (2024, January 22). "Million-dollar HDB flats: Trend or anomaly?" *The Straits Times*. <https://www.straitstimes.com/singapore/million-dollar-hdb-flats-trend-or-anomaly>
- Today. (2023, November 12). "Young Singaporeans face uphill battle in housing market as prices outpace income growth." *Today Online*. <https://www.todayonline.com/singapore/young-singaporeans-face-uphill-battle-housing-market-prices-outpace-income-growth>

Industry Reports and Analyses

- CBRE. (2023). *Singapore Market Outlook 2024*. CBRE Research.
- ERA Singapore. (2024). *Singapore Residential Market Report Q1 2024*. ERA Research & Consultancy.
- Knight Frank. (2023). *Singapore Residential Market Update Q4 2023*. Knight Frank Research.
- PropNex. (2024). *Singapore Property Market Index Q1 2024*. PropNex Research.
- URA. (2024). *Real Estate Market Trends Q1 2024*. Urban Redevelopment Authority.

Online Resources and Websites

- HDB InfoWEB. (2024). Housing & Development Board Official Website. <https://www.hdb.gov.sg>
- MAS. (2024). Property Market Cooling Measures. Monetary Authority of Singapore. <https://www.mas.gov.sg/regulation/property-market-cooling-measures>
- MND. (2024). Housing Policies. Ministry of National Development. <https://www.mnd.gov.sg/our-work/housing>
- PropertyGuru Singapore. (2024). Singapore Property Market Outlook. <https://www.propertyguru.com.sg/property-guides/property-market-outlook-singapore-13766>
- URA. (2024). Master Plan 2019. Urban Redevelopment Authority. <https://www.ura.gov.sg/Corporate/Planning/Master-Plan>