

GROUP PROJECT-4

Project Report: Development of Comprehensive Dashboards Using Web Frameworks

Course: IE6600 Computation and Visualization

Spring Semester 2024

Group Number - 10

Sri Sai Prabhath Reddy Gudipalli (002207631)

Deepthi Umesha (002661626)

Ashwini Mahadevaswamy (002661627)

Yashwant Dontam (002844794)

Introduction

This project report explores and analyzes New York City property sales data from January 2016 to December 2022. Our goal is to reveal trends, distributions, and anomalies in real estate transactions across different neighborhoods to aid stakeholders in making informed decisions. We utilize advanced data manipulation and visualization techniques, including geocoding, data cleaning, and interactive visualizations through tools like Pandas and Streamlit. The insights derived from this analysis provide a comprehensive overview of the NYC real estate market, highlighting active areas and identifying sales patterns related to property characteristics.

Task: Data Acquisition and Inspection

```
1 import pandas as pd
2 df = pd.read_csv('NYC_Property_Sales_Data_Geocoded.csv', low_memory=False)
3
4 missing_percentages = (df.isnull().sum() / len(df)) * 100
5 columns_to_drop = missing_percentages[missing_percentages > 70].index
6 df.drop(columns=columns_to_drop, inplace=True)
7 print(columns_to_drop)
8
9 threshold = 0.5 * len(df.columns)
10 missing_values_per_row = df.isnull().sum(axis=1)
11 rows_to_drop = df[missing_values_per_row > threshold]
12 print("Number of rows with more than 50% missing values:", len(rows_to_drop))
13 df.drop(index=rows_to_drop.index, inplace=True)
14 print("Shape of the DataFrame after dropping rows:", df.shape)
15
16 df['YEAR BUILT'].replace('0', pd.NaT, inplace=True)
17 df['YEAR BUILT'] = pd.to_datetime(df['YEAR BUILT'], format='%Y', errors='coerce')
18 df.dropna(subset=['YEAR BUILT'], inplace=True)
19
20 def clean_square_feet(value):
21     try:
22         # Check if the value is not null and is a string
23         if pd.notnull(value) and isinstance(value, str):
24             # Remove commas and spaces
25             cleaned_value = value.replace(',', '').replace(' ', '')
26             # Convert to float
27             return float(cleaned_value)
28         else:
29             # Return None for non-string or NaN values
30             return None
31     except ValueError:
32         # Handle exception for strings like '- 0'
33         return None # or any other appropriate action
34
35 # Apply the cleaning function to 'LAND SQUARE FEET' column
36 df['LAND SQUARE FEET'] = df['LAND SQUARE FEET'].apply(clean_square_feet)
37 df['LAND SQUARE FEET'] = df['LAND SQUARE FEET'].astype(float)
38
39 def clean_square_feet(value):
40     try:
41         # Check if the value is not null and is a string
42         if pd.notnull(value) and isinstance(value, str):
43             # Remove commas and spaces
44             cleaned_value = value.replace(',', '').replace(' ', '')
45             # Convert to float
46             return float(cleaned_value)
47         else:
48             # Return None for non-string or NaN values
49             return None
50     except ValueError:
51         # Handle exception for strings like '- 0'
52         return None # or any other appropriate action
53
```

```

35 # Apply the cleaning function to 'LAND SQUARE FEET' column
36 df['LAND SQUARE FEET'] = df['LAND SQUARE FEET'].apply(clean_square_feet)
37 df['LAND SQUARE FEET'] = df['LAND SQUARE FEET'].astype(float)
38
39 def clean_square_feet(value):
40     try:
41         # Check if the value is not null and is a string
42         if pd.notnull(value) and isinstance(value, str):
43             # Remove commas and spaces
44             cleaned_value = value.replace(',', '').replace(' ', '')
45             # Convert to float
46             return float(cleaned_value)
47         else:
48             # Return None for non-string or NaN values
49             return None
50     except ValueError:
51         # Handle exception for strings like '- 0'
52         return None # or any other appropriate action
53
54 # Apply the cleaning function to 'LAND SQUARE FEET' column
55 df['GROSS SQUARE FEET'] = df['GROSS SQUARE FEET'].apply(clean_square_feet)
56 df['GROSS SQUARE FEET'] = df['GROSS SQUARE FEET'].astype(float)
57
58 df['LAND SQUARE FEET'] = df['LAND SQUARE FEET'].replace(-0.0, 0.0)
59 df['GROSS SQUARE FEET'] = df['GROSS SQUARE FEET'].replace(-0.0, 0.0)
60 df['LAND SQUARE FEET'].fillna(0.0, inplace=True)
61 df['GROSS SQUARE FEET'].fillna(0.0, inplace=True)
62
63 df['COMMERCIAL UNITS'] = df['COMMERCIAL UNITS'].fillna(0)
64 df['RESIDENTIAL UNITS'] = df['RESIDENTIAL UNITS'].fillna(0)
65
66 df['NTA'].fillna('Unknown', inplace=True)
67 df.drop(columns=['BIN', 'BBL', 'inplace=True)
68 df.drop(columns=['Census Tract', 'BOROUGH', 'Council District', 'Community Board'], inplace=True)
69 df['TAX CLASS AS OF FINAL ROLL'] = df.groupby('YEAR BUILT')['TAX CLASS AS OF FINAL ROLL'].transform(lambda x: x.fillna(x.mode()[0] if not x.mode().empty else "Unknown"))
70 df['BUILDING CLASS AS OF FINAL ROLL'] = df.groupby('YEAR BUILT')['BUILDING CLASS AS OF FINAL ROLL'].transform(lambda x: x.fillna(x.mode()[0] if not x.mode().empty else "Unknown"))
71
72 # Ensure 'Latitude' and 'Longitude' are in numeric form (if they're not already)
73 df['Latitude'] = pd.to_numeric(df['Latitude'], errors='coerce')
74 df['Longitude'] = pd.to_numeric(df['Longitude'], errors='coerce')
75
76 # Group by 'BLOCK' (or another geographical marker) and interpolate within each group
77 df['Latitude'] = df.groupby('BLOCK')['Latitude'].transform(lambda x: x.interpolate(method='linear', limit_direction='both'))
78 df['Longitude'] = df.groupby('BLOCK')['Longitude'].transform(lambda x: x.interpolate(method='linear', limit_direction='both'))
79 df['Latitude'] = df.groupby('NEIGHBORHOOD')['Latitude'].transform(lambda x: x.fillna(x.mean()))
80 df['Longitude'] = df.groupby('NEIGHBORHOOD')['Longitude'].transform(lambda x: x.fillna(x.mean()))
81 df['TOTAL UNITS'] = df['RESIDENTIAL UNITS'] + df['COMMERCIAL UNITS']
82 df['ZIP CODE'] = df.groupby('NEIGHBORHOOD')['ZIP CODE'].transform(lambda x: x.fillna(x.mode()[0] if not x.mode().empty else np.nan))
83 df['SALE DATE'] = pd.to_datetime(df['SALE DATE'])
84 df['BUILDING CLASS CATEGORY'] = df['BUILDING CLASS CATEGORY'].str.replace(r'^\d+s+', '', regex=True)
85 df['BUILDING CLASS CATEGORY'] = df['BUILDING CLASS CATEGORY'].str.title()
86 df['ZIP CODE'] = df['ZIP CODE'].astype(str).str.replace('\.0', '', regex=True)

```

In the preprocessing stage of our analysis of the New York City Property Sales dataset, we executed several essential steps to ensure the data's integrity and usability. These steps, implemented using the Python library Pandas. Initially, we loaded the dataset, optimizing for memory efficiency due to the data's extensive size. We then identified columns with excessive missing values—specifically, those with over 70% absence—and removed these columns as they could compromise the validity of any subsequent analysis. Similarly, rows with more than half of their values missing were also excluded.

Next, we transformed the 'YEAR BUILT' column from a string representation to a datetime format, handling any non-standard entries, such as '0', by converting them to a missing value indicator. We also cleaned numeric fields, particularly square footage measurements, by stripping non-numeric characters

and converting the strings to floats, ensuring that these values were correctly interpreted for quantitative analysis.

To refine the dataset further, we addressed missing geographical coordinates by interpolating 'Latitude' and 'Longitude' within groups of the same 'BLOCK' and 'NEIGHBORHOOD', which preserved spatial relationships. Additionally, we engineered a new 'TOTAL UNITS' feature by summing the counts of residential and commercial units, providing a more comprehensive metric for property analysis. Through these preprocessing measures, we established a robust foundation for detailed analysis, enabling us to draw more accurate and meaningful insights from the property sales data.

Task: Dashboard Development and Web Deployment

```
92 #1
93 import streamlit as st
94
95 # Set the title of your Streamlit app
96 st.subheader("NYC Property Sales Data Exploration")
97
98 # Create a sidebar for the date range filter
99 st.sidebar.header("Filter for Data Exploration")
100
101
102 # Get the minimum and maximum sale dates from your dataset for the date input range
103 min_date = df['SALE DATE'].min()
104 max_date = df['SALE DATE'].max()
105
106 # Use Streamlit's date_input widget to get a date range from the user
107 date_range = st.sidebar.date_input("Sale Date Range", value=(min_date, max_date), min_value=min_date, max_value=max_date)
108
109 # Filter the DataFrame based on the selected date range
110 filtered_df = df[(df['SALE DATE'] >= pd.to_datetime(date_range[0])) & (df['SALE DATE'] <= pd.to_datetime(date_range[1]))]
111
112 # Select only the specified columns for display
113 columns_to_display = ['NEIGHBORHOOD', 'BLOCK', 'ADDRESS', 'ZIP CODE', 'YEAR BUILT', 'SALE PRICE', 'SALE DATE']
114 filtered_df = filtered_df[columns_to_display]
115
116 # Display the filtered DataFrame
117 st.dataframe(filtered_df)
118
119 # Show the number of sales in the selected period
120 st.write(f"Total sales in the selected period: {len(filtered_df)}")
121
122
```

The code here utilizes Streamlit, an open-source app framework, for creating a simple web application for exploring property sales data in New York City.

The script starts by importing Streamlit and sets up a header for the web application titled "NYC Property Sales Data Exploration". This establishes the main purpose of the application for end-users. It introduces a sidebar with a date range filter, prompting users to select a specific period for analysis. Utilizing Streamlit's interactive `date_input` widget, the script fetches the minimum and maximum sale dates from

the dataset to set the boundaries for the date selection. Users can refine the data displayed by selecting a date range within these limits.

Once a date range is chosen, the script filters the DataFrame accordingly. It ensures that only the data entries with sale dates within the selected range are included in the subsequent analysis. The script then focuses on displaying the filtered data. Specific columns are selected for this purpose, such as 'NEIGHBORHOOD', 'BLOCK', 'ADDRESS', 'ZIP CODE', 'YEAR BUILT', 'SALE PRICE', and 'SALE DATE', providing a concise yet informative view. The filtered DataFrame is displayed in the main panel of the web application.

	NEIGHBORHOOD	BLOCK	ADDRESS	ZIP CODE	YEAR BUILT	SALE PRICE
0	CHELSEA	697	555 WEST 25TH STREET	10001	1926-01-01 00:00:00	43,300,000
1	CHELSEA	697	511 WEST 25TH STREET	10001	1917-01-01 00:00:00	148,254,147
2	CHELSEA	700	538 WEST 29TH STREET	10001	1910-01-01 00:00:00	11,000,000
3	CHELSEA	712	450 WEST 15TH	10011	1936-01-01 00:00:00	591,800,000
4	CHELSEA	746	340 WEST 23RD STREET	10011	1900-01-01 00:00:00	0
5	CHELSEA	802	158 WEST 27 STREET	10001	1913-01-01 00:00:00	99,350,000
6	CHELSEA	803	307 7 AVENUE	10001	1926-01-01 00:00:00	115,000,000
7	CHELSEA	697	521 WEST 25TH STREET	10001	1910-01-01 00:00:00	148,254,147
8	CHELSEA	772	250 WEST 23RD STREET	10011	1948-01-01 00:00:00	14,500,000
9	CHELSEA	772	254 WEST 23RD STREET	10011	1920-01-01 00:00:00	4,750,000

Total sales in the selected period: 558785

The image displayed from a web application, showcasing a table of property sales data filtered by a specific period. The table contains columns for 'NEIGHBORHOOD', 'BLOCK', 'ADDRESS', 'ZIP CODE', 'YEAR BUILT', and 'SALE PRICE'. Each row corresponds to a property sale in the Chelsea neighborhood, providing details such as the exact address, the block number, the postal code, the construction date of the property, and the sale price in USD. The sale prices vary significantly, highlighting the diversity in property values within the area. For instance, the property at '555 WEST 25TH STREET' sold for \$43,300,000, while the one at '340 WEST 23RD STREET' shows a sale price of \$0, which may indicate a non-monetary transaction or missing data.

Below the table, a summarizing statement reveals the total number of sales transactions in the selected period, amounting to 558,785. This reflects the application's capability to not only present detailed records but also provide aggregate information that can be vital for understanding overall market activity during a given timeframe.

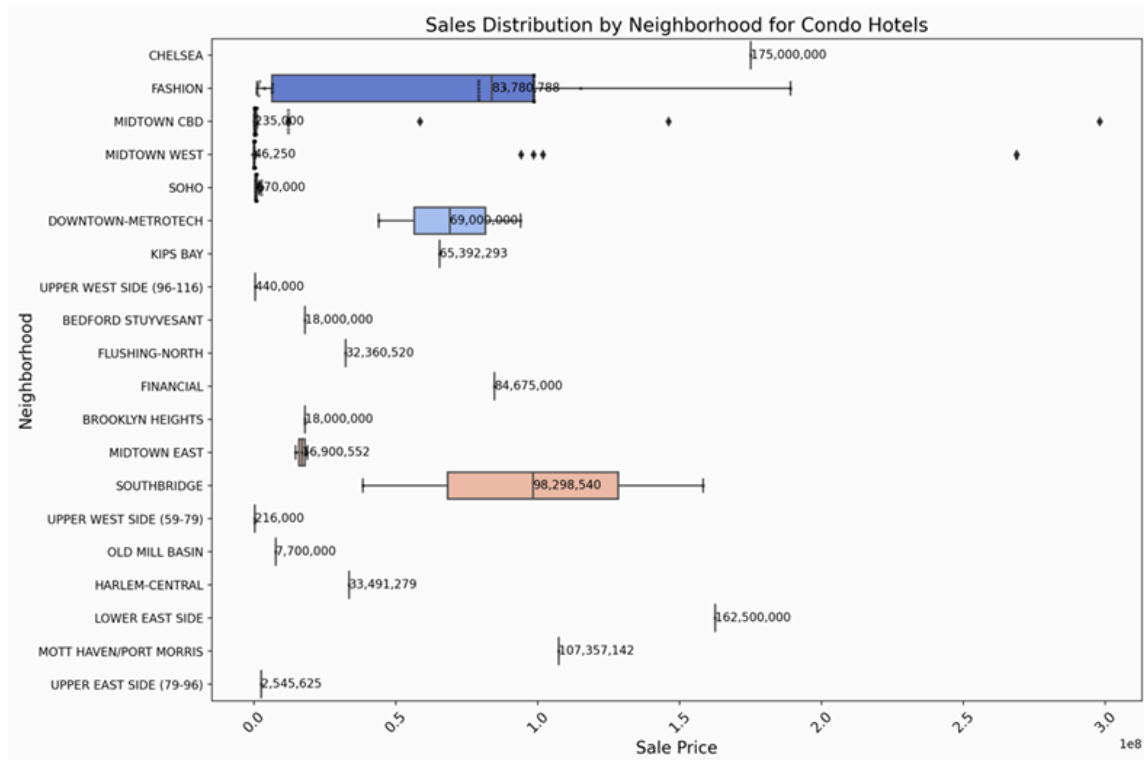
```
123 #2
124 import streamlit as st
125 import pandas as pd
126 import seaborn as sns
127 import matplotlib.pyplot as plt
128 from matplotlib.ticker import StrMethodFormatter
129
130 # Assuming your DataFrame is named 'df'
131 # Convert 'SALE PRICE' to numeric and filter
132 df = df[(df['SALE PRICE'] > 0) & (df['SALE PRICE'].notnull())]
133
134 # Set the title of your Streamlit app
135 st.subheader("Sales Distribution by Neighborhood")
136
137 # Create a selectbox for Building Class Category
138 selected_building_class = st.selectbox("Select Building Class Category:", df['BUILDING CLASS CATEGORY'].unique())
139
140 # Filter based on selected building class
141 filtered_df = df[df['BUILDING CLASS CATEGORY'] == selected_building_class]
142
143 # Start the plot
144 plt.figure(figsize=(14, 10)) # Increased figure size for clarity
145 sns.boxplot(
146     x='SALE PRICE',
147     y='NEIGHBORHOOD',
148     data=filtered_df,
149     palette='coolwarm', # A lighter palette for a fresh look
150 )
151
152 # Improve readability
153 plt.xticks(rotation=45, fontsize=12)
154 plt.yticks(fontsize=10)
155 plt.xlabel("Sale Price", fontsize=14)
156 plt.ylabel("Neighborhood", fontsize=14)
157 plt.title(f"Sales Distribution by Neighborhood for {selected_building_class}", fontsize=16)
158
159 # Overlay with swarmplot to show individual data points, commented out by default due to performance concerns
160 # sns.swarmplot(x='SALE PRICE', y='NEIGHBORHOOD', data=filtered_df, color='black', alpha=0.5, size=3)
161
162 # Annotate median prices
163 medians = filtered_df.groupby(['NEIGHBORHOOD'])['SALE PRICE'].median().sort_values(ascending=False)
164 for i, neighborhood in enumerate(filtered_df['NEIGHBORHOOD'].unique()):
165     median_val = medians[neighborhood]
166     plt.text(median_val, i, f'{median_val:,.0f}', va='center', fontsize=10, color='black')
167
168 # Display the plot
169 st.pyplot(plt)
170
```

The provided script is a Python-based tool that uses Streamlit for interactivity, along with Pandas for data handling, and Seaborn and Matplotlib for visualization, enabling users to interactively explore property sale prices by neighborhood. It includes a Streamlit select box for filtering data by building class and generates a boxplot to visualize sale price distributions, which are enhanced with median price annotations for immediate insights. This setup offers an intuitive interface for users to analyze the real estate market trends within various neighborhoods without delving into the underlying code.

Sales Distribution by Neighborhood

Select Building Class Category:

Condo Hotels



The above depict an interactive visualization of New York City's Condo Hotel sales, grouped by neighborhood. A boxplot graphically represents the range and outliers of sale prices, allowing for quick comparison across areas. Users can select different building classes from a dropdown to dynamically update the chart, making it a flexible tool for analyzing real estate market trends.

```

172 #3
173 import streamlit as st
174 import pandas as pd
175 import numpy as np
176 import altair as alt
177
178
179 # Extract year from 'SALE DATE' column
180 df['YEAR'] = df['SALE DATE'].dt.year
181
182 # Sidebar filters
183 st.sidebar.header('Filters for Sales Time Series Analysis')
184 min_date = min(df['SALE DATE']).date()
185 max_date = max(df['SALE DATE']).date()
186 start_date = st.sidebar.date_input('Start Date', min_value=min_date, max_value=max_date, value=min_date)
187 end_date = st.sidebar.date_input('End Date', min_value=min_date, max_value=max_date, value=max_date)
188
189 # Convert start_date and end_date to datetime
190 start_date = pd.to_datetime(start_date)
191 end_date = pd.to_datetime(end_date)
192
193 # Filter data based on selected dates
194 filtered_data = df[(df['SALE DATE'] >= start_date) & (df['SALE DATE'] <= end_date)]
195
196 # Group data by YEAR and calculate total sales count
197 time_series_data = filtered_data.groupby('YEAR').size().reset_index(name='Total Sales')
198
199 st.subheader("NYC Property Sales Time Series Analysis")
200
201 # Plot time series
202 chart = alt.Chart(time_series_data).mark_line().encode(
203     x='YEAR:O', # Using ordinal scale for discrete years
204     y='Total Sales'
205 ).properties(
206     width=800,
207     height=500
208 ).interactive()
209
210 st.altair_chart(chart, use_container_width=True)
211
212 # Summary statistics
213 total_sales = filtered_data.shape[0]
214 average_price = filtered_data['SALE PRICE'].mean()
215 median_price = filtered_data['SALE PRICE'].median()
216
217 st.subheader('Summary Statistics')
218 st.write(f'Total Sales: {total_sales}')
219 st.write(f'Average Sale Price: ${average_price:,.2f}')
220 st.write(f'Median Sale Price: ${median_price:,.2f}')
221
222

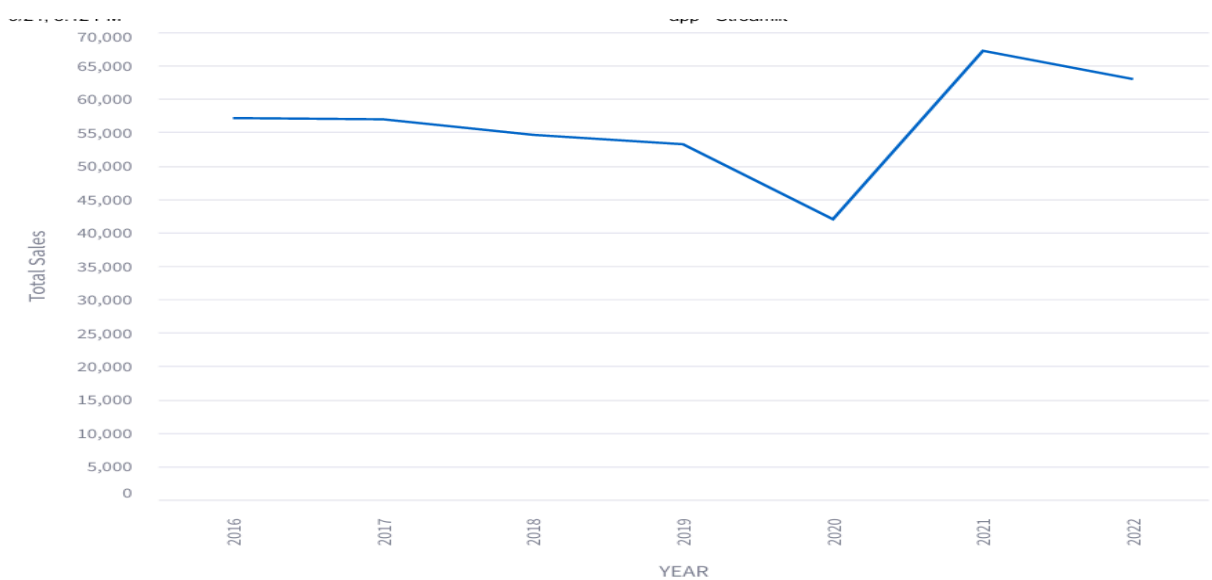
```

The script is designed to analyze property sales over time within a Streamlit web application. It begins by extracting the year from a 'SALE DATE' column in the data and then uses Streamlit to create interactive sidebar filters that allow users to select a date range for analysis. With the specified date range, the script filters the data and groups it by year to calculate the total number of sales per year.

An Altair line chart is then created to display this time series data, showing how the number of sales has changed over time. The chart is configured to be interactive, providing users with an engaging way to explore the sales trends.

In addition, the script computes summary statistics, including the total number of sales, the average sale price, and the median sale price for the filtered data. These statistics are displayed under a 'Summary Statistics' subheader in the Streamlit app, offering users quick numerical insights into the property sales data within their selected time frame.

NYC Property Sales Time Series Analysis



Summary Statistics

Total Sales: 394112

Average Sale Price: \$1,785,120.51

Median Sale Price: \$695,000.00

The above visualization is a line chart from a time series analysis of property sales in New York City, spanning from 2016 to 2022. The chart indicates the number of property sales each year, revealing trends such as a decline leading up to 2020, a sharp dip in 2020, followed by a significant recovery in 2021, and a slight drop in 2022. This visual suggests fluctuating real estate activity, with a notable impact around 2020, possibly due to external factors such as economic conditions or events impacting the housing market.

```

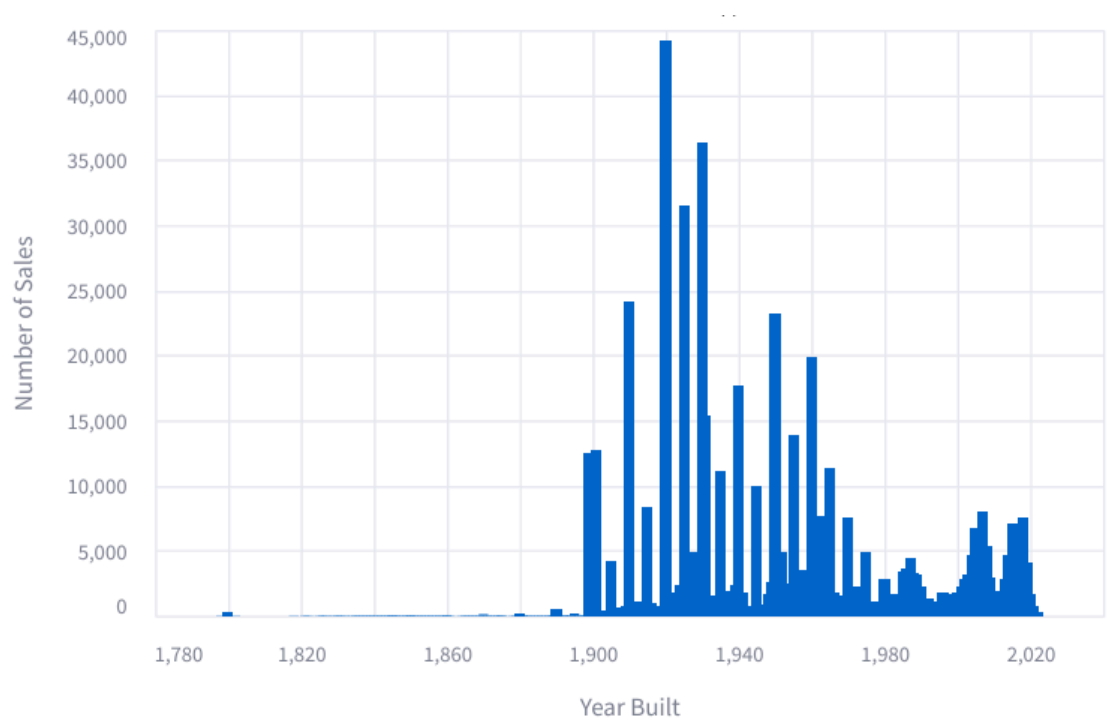
223 #4
224 import streamlit as st
225 import pandas as pd
226 import altair as alt
227
228 st.sidebar.header('Filter for Year Built')
229
230
231 # Filter by Year Built (Building Age)
232 year_built_range = st.sidebar.slider("Select Year Built Range:", 1798, 2022, (1798, 2022))
233
234 # Load your cleaned dataset
235 df = pd.read_csv("NYC_Property_Sales_Data.csv")
236
237 # Convert 'YEAR BUILT' column to integer type
238 df['YEAR BUILT'] = pd.to_numeric(df['YEAR BUILT'], errors='coerce')
239
240 # Apply filters to the DataFrame
241 filtered_df = df[(df['YEAR BUILT'] >= year_built_range[0]) &
242                 (df['YEAR BUILT'] <= year_built_range[1])]
243
244 # Display histogram for Building Age
245 histogram_chart = alt.Chart(filtered_df).mark_bar().encode(
246     x=alt.X('YEAR BUILT', title='Year Built'),
247     y=alt.Y('count()', title='Number of Sales'),
248     tooltip=['YEAR BUILT', 'count()']
249 ).properties(
250     width=600,
251     height=400
252 ).interactive()
253
254 st.subheader("Distribution of Property Sales by Year Built")
255 st.altair_chart(histogram_chart)

```

The script creates an interactive feature within a Streamlit application that allows users to filter property sales data based on the year a property was built. A slider in the application's sidebar enables users to define a range of years for which they want to analyze sales. The data is then loaded from a CSV file, and the 'YEAR BUILT' column is converted to a numeric type to ensure proper filtering.

Applying the selected year range as a filter to the dataset, the script generates a histogram using Altair that displays the distribution of the number of sales per year built within the chosen range. The histogram is designed to be interactive, allowing users to explore different facets of the data in a more engaging way. The visualization aids in identifying trends or patterns based on the age of properties sold within the specified timeframe.

Distribution of Property Sales by Year Built



The above visualization displays a histogram illustrating the distribution of property sales in New York City by year built. The x-axis represents the construction year of properties, while the y-axis indicates the number of sales. There is a notable concentration of sales for properties built in the early 20th century, with peaks occurring around the 1920s. The frequency of sales declines for properties built in the mid-to-late 20th century, with a resurgence in sales for more recent constructions. This visualization aids in understanding popular times for property development in relation to sales activity.

```

257 #5
258 import streamlit as st
259 import pandas as pd
260 import altair as alt
261
262
263 # Set the title of your Streamlit app
264 st.subheader("Property Characteristics vs. Sale Price Analysis")
265
266 # Create filter widgets
267 selected_property_characteristic = st.selectbox("Select Property Characteristic:",
268                                                ['LAND SQUARE FEET', 'TOTAL UNITS', 'GROSS SQUARE FEET'])
269 selected_neighborhood = st.selectbox("Select Neighborhood:", df['NEIGHBORHOOD'].unique())
270 selected_tax_class = st.selectbox("Select Tax Class:", df['TAX CLASS AS OF FINAL ROLL'].unique())
271
272 # Apply filters to the DataFrame
273 filtered_df = df[(df['NEIGHBORHOOD'] == selected_neighborhood) &
274                 (df['TAX CLASS AS OF FINAL ROLL'] == selected_tax_class)]
275
276 # Create scatter plot
277 scatter_plot = alt.Chart(filtered_df).mark_circle().encode(
278     x=selected_property_characteristic,
279     y='SALE PRICE',
280     tooltip=['ADDRESS', 'SALE PRICE', 'BUILDING CLASS CATEGORY']
281 ).properties(
282     width=800,
283     height=500
284 ).interactive()
285
286 # Display scatter plot
287 st.subheader("Property Characteristics vs. Sale Price")
288 st.altair_chart(scatter_plot)
289
290 # Summary statistics
291 average_price = filtered_df['SALE PRICE'].mean()
292 median_price = filtered_df['SALE PRICE'].median()
293 total_sales = len(filtered_df)
294
295 st.subheader("Summary Statistics")
296 st.write(f"Average Sale Price: ${average_price:,.2f}")
297 st.write(f"Median Sale Price: ${median_price:,.2f}")
298 st.write(f"Total Sales: {total_sales}")
299

```

The above code sets up an interactive web application using Streamlit for visualizing the relationship between property characteristics and sale prices. It allows users to select a property characteristic, neighborhood, and tax class from dropdown menus. Applying these filters to the dataset, it then generates an Altair scatter plot to visually explore how the chosen characteristic impacts sale prices in the selected neighborhood and tax class. Additionally, it calculates and displays summary statistics, including the average and median sale prices, as well as the total number of sales, offering a quick numerical insight into the filtered subset of property sales data.

Property Characteristics vs. Sale Price Analysis

Select Property Characteristic:

LAND SQUARE FEET

▼

Select Neighborhood:

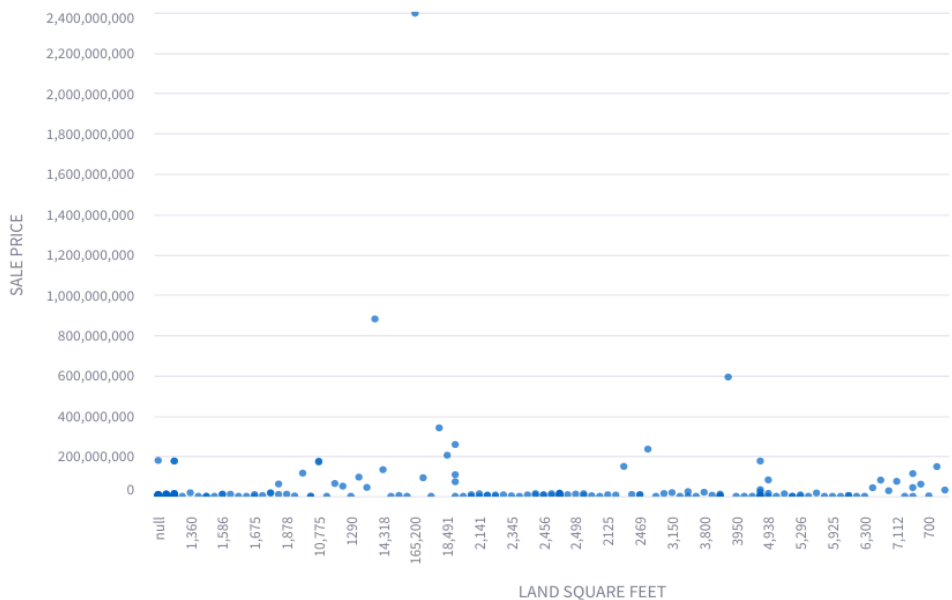
CHELSEA

▼

Select Tax Class:

4

▼



Summary Statistics

Average Sale Price: \$32,092,956.80

Median Sale Price: \$2,290,000.00

Total Sales: 288

The above visualization illustrates a scatter plot depicting the relationship between land square footage and property sale prices. Each point represents a property, with the x-axis indicating the land size in square feet and the y-axis showing the sale price in USD. The plot reveals a distribution of sale prices across various land sizes, with a cluster of data points at the lower end of both axes, suggesting that smaller land sizes are more commonly sold and generally at lower prices. There are also outliers indicating that some larger properties command significantly higher prices, although such sales are less frequent. This visualization is useful for understanding how land size influences property value in the real estate market.

```

301 #6
302 import streamlit as st
303 import pandas as pd
304 import geopandas as gpd
305 import plotly.express as px
306 import altair as alt
307 from shapely.geometry import Point
308
309 # Create Streamlit app
310 st.subheader('NYC Property Sales by Neighborhood')
311
312 # Filter for Date Range
313 start_date = df['SALE DATE'].min()
314 end_date = df['SALE DATE'].max()
315
316 # Filter for Property Type
317 property_type = st.selectbox('Property Type', df['BUILDING CLASS CATEGORY'].unique())
318
319 # Filter data based on selected filters
320 filtered_df = df[(df['SALE DATE'] >= start_date) & (df['SALE DATE'] <= end_date) &
321                 (df['BUILDING CLASS CATEGORY'] == property_type)]
322
323 # Plot map using Plotly
324 fig = px.scatter_mapbox(filtered_df, lat="Latitude", lon="Longitude", color="NEIGHBORHOOD",
325                         hover_data=["ADDRESS", "SALE PRICE"],
326                         mapbox_style="carto-positron", zoom=10)
327
328 fig.update_layout(margin={"r":0,"t":0,"l":0,"b":0})
329
330 # Display map with points
331 st.plotly_chart(fig)
332
333

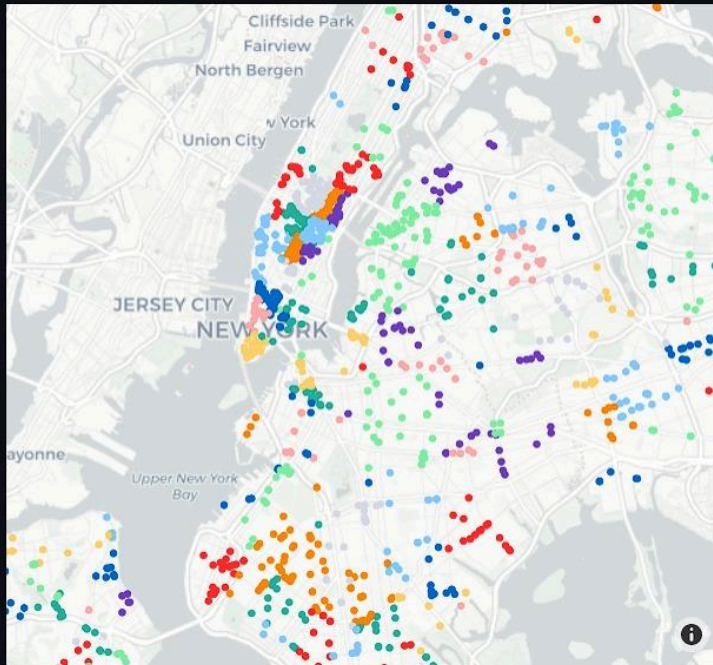
```

The script configures a Streamlit web application to visualize New York City property sales on an interactive map. Users can filter sales by date and property type via Streamlit's sidebar widgets. The data is then plotted on a map using Plotly Express, with points representing sales, color-coded by neighborhood, and including details like address and sale price on hover. This visualization allows for an intuitive understanding of the geographical distribution of property sales within the selected parameters.

NYC Property Sales by Neighborhood

Property Type

21 OFFICE BUILDINGS



NEIGHBORHOOD

- CHELSEA
- CHINATOWN
- CIVIC CENTER
- CLINTON
- EAST VILLAGE
- FASHION
- FINANCIAL
- FLATIRON
- GRAMERCY
- GREENWICH VILLAGE-CENTRAL
- GREENWICH VILLAGE-WEST
- HARLEM-EAST
- HARLEM-WEST
- JAVITS CENTER
- LITTLE ITALY
- LOWER EAST SIDE
- MANHATTAN-UNKNOWN
- MIDTOWN CBD
- MIDTOWN EAST
- MIDTOWN WEST
- MURRAY HILL
- SOHO
- TIMES SQUARE

The above visualization showcases a map from a web application depicting the distribution of property sales in various neighborhoods across New York City, with a focus on office buildings. Color-coded dots on the map represent individual property sales, with the colors corresponding to different neighborhoods as indicated in the legend. This visual format provides a clear geographic distribution of office building sales, allowing users to quickly discern patterns and concentrations of sales activity across the city. The interactive dropdown menu enables users to select and filter the property type, in this case, '21 OFFICE BUILDINGS', for a tailored analytical view.

```

335 #7
336 import streamlit as st
337 import pandas as pd
338 import plotly.express as px
339
340 ## Create Streamlit app
341 st.subheader('Property Characteristics vs. Sale Price Analysis')
342
343 # Filter widgets
344 property_type = st.selectbox('Select Property Type', df['BUILDING CLASS CATEGORY'].unique(), key='property_type_selectbox')
345 neighborhood = st.selectbox('Select Neighborhood', df['NEIGHBORHOOD'].unique(), key='neighborhood_selectbox')
346 start_date = st.date_input('Start Date', min_value=pd.to_datetime(df['SALE DATE']).min(), max_value=pd.to_datetime(df['SALE DATE']).max(),
347 value=pd.to_datetime(df['SALE DATE']).min(), key='start_date_input')
348 end_date = st.date_input('End Date', min_value=pd.to_datetime(df['SALE DATE']).min(), max_value=pd.to_datetime(df['SALE DATE']).max(),
349 value=pd.to_datetime(df['SALE DATE']).max(), key='end_date_input')
350
351 # Convert start_date and end_date to datetime objects
352 start_date = pd.to_datetime(start_date)
353 end_date = pd.to_datetime(end_date)
354
355 # Filter data based on selected filters
356 filtered_df = df[(df['BUILDING CLASS CATEGORY'] == property_type) &
357 (df['NEIGHBORHOOD'] == neighborhood) &
358 (pd.to_datetime(df['SALE DATE']) >= start_date) &
359 (pd.to_datetime(df['SALE DATE']) <= end_date)]
360
361 # Plotly scatter plot
362 fig = px.scatter(filtered_df, x='TOTAL UNITS', y='SALE PRICE',
363 hover_data=['ADDRESS', 'LAND SQUARE FEET', 'GROSS SQUARE FEET'],
364 trendline='ols', title='Total Units vs. Sale Price',
365 labels={'TOTAL UNITS': 'Total Units', 'SALE PRICE': 'Sale Price'})
366
367 # Customize Layout
368 fig.update_layout(showlegend=True)
369
370 # Display plot
371 st.plotly_chart(fig)
372
373 # Summary statistics
374 average_price = filtered_df['SALE PRICE'].mean()
375 median_price = filtered_df['SALE PRICE'].median()
376 total_sales = len(filtered_df)
377
378 st.subheader("Summary Statistics")
379 st.write(f"Average Sale Price: ${average_price:,.2f}")
380 st.write(f"Median Sale Price: ${median_price:,.2f}")
381 st.write(f"Total Sales: {total_sales}")

```

The above code generates an interactive plot for a real estate web application, correlating property characteristics with sale prices. Utilizing Streamlit, users can filter properties by type and neighborhood and select a sale date range. After data filtering based on these selections, a Plotly scatter plot is created, highlighting the relationship between 'TOTAL UNITS' and 'SALE PRICE', enriched with a trend line. Additional hover details like 'ADDRESS' and 'GROSS SQUARE FEET' provide deeper insights. The plot is customizable and designed for clarity in presentation. Summary statistics such as average and median sale prices, as well as the total number of sales, are computed and displayed, giving users immediate access to key market indicators.

Property Characteristics vs. Sale Price Analysis

Select Property Type

21 OFFICE BUILDINGS

▼

Select Neighborhood

CHELSEA

▼

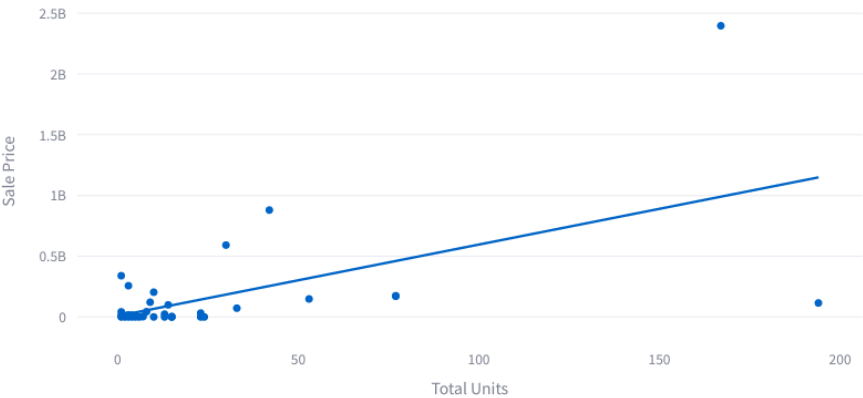
Start Date

2016/01/01

End Date

2022/12/31

Total Units vs. Sale Price



Summary Statistics

Average Sale Price: \$138,743,985.24

Median Sale Price: \$11,666,666.50

Total Sales: 42

The above visualization portrays a scatter plot visualizing the relationship between the total units in properties and their sale prices, along with a line indicating the trend of increasing prices with more units. Notably, there appears to be a positive correlation between the size of a property, in terms of total units, and its sale price. Additionally, summary statistics are provided below the graph, including an average sale price of approximately \$138.7 million, a median sale price of about \$11.7 million, and a total of 42 sales captured in the dataset. This information can help stakeholders gauge market value trends related to property size.

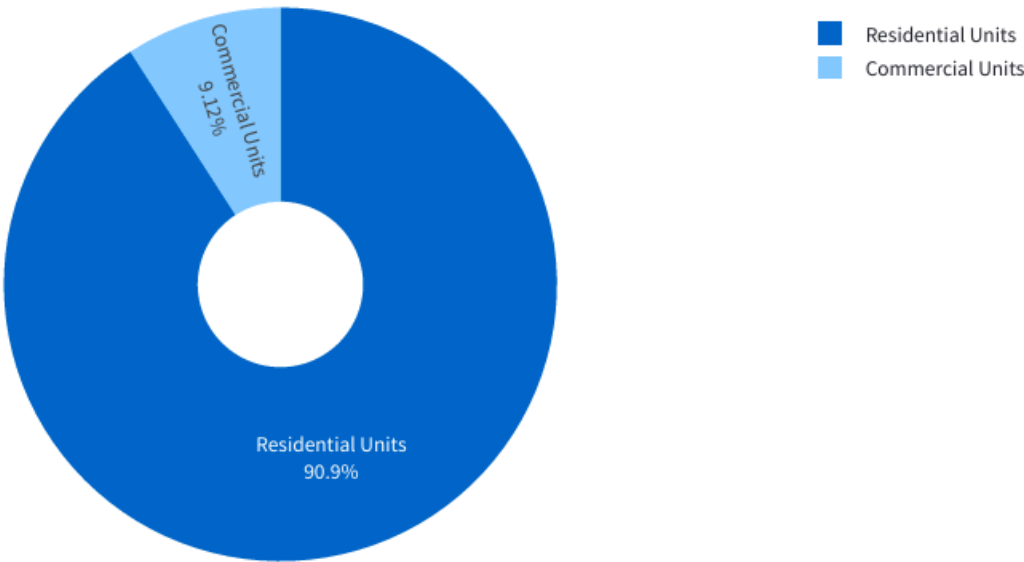
```

382 #8
383 import streamlit as st
384 import pandas as pd
385 import plotly.express as px
386
387 st.subheader('Unit Type Distribution')
388
389 # Calculate the total number of residential and commercial units
390 total_residential_units = df['RESIDENTIAL UNITS'].sum()
391 total_commercial_units = df['COMMERCIAL UNITS'].sum()
392
393 # Create a DataFrame for the pie chart
394 data = pd.DataFrame({
395     'Unit Type': ['Residential Units', 'Commercial Units'],
396     'Total Units': [total_residential_units, total_commercial_units]
397 })
398
399 # Create an interactive pie chart using Plotly Express
400 fig = px.pie(data, values='Total Units', names='Unit Type',
401             hover_name='Unit Type',
402             labels={'Unit Type': 'Unit Type'},
403             hole=0.3)
404
405 # Add labels to the pie chart sectors
406 fig.update_traces(textinfo='percent+label')
407
408 # Display the pie chart
409 st.plotly_chart(fig)
410

```

The above code utilizes Streamlit alongside Pandas and Plotly Express to create a pie chart for displaying the distribution of unit types in a property dataset. The Streamlit header indicates the focus of the visualization on 'Unit Type Distribution'. The script calculates the total counts of residential and commercial units within the dataset. A DataFrame is constructed with these totals, which is then used to generate a pie chart. This chart is designed to be interactive, with hover labels and a donut style (indicated by the 'hole' parameter set to 0.3), providing a visual breakdown of the proportion of residential to commercial units. Percentage labels are added to each segment for clarity. The resulting plot is displayed within the Streamlit app, offering a quick and easy understanding of the composition of property types.

Unit Type Distribution



The image displays a donut chart illustrating the distribution of unit types within a property dataset. The chart shows a predominant share of residential units, accounting for 90.9% of the total, while commercial units make up a smaller fraction at 9.1%. This visual representation emphasizes the significantly higher proportion of residential units in comparison to commercial spaces, offering an at-a-glance understanding of the composition of the properties analyzed. The chart serves as an analytical tool to quickly assess the balance of residential versus commercial property types within the market.

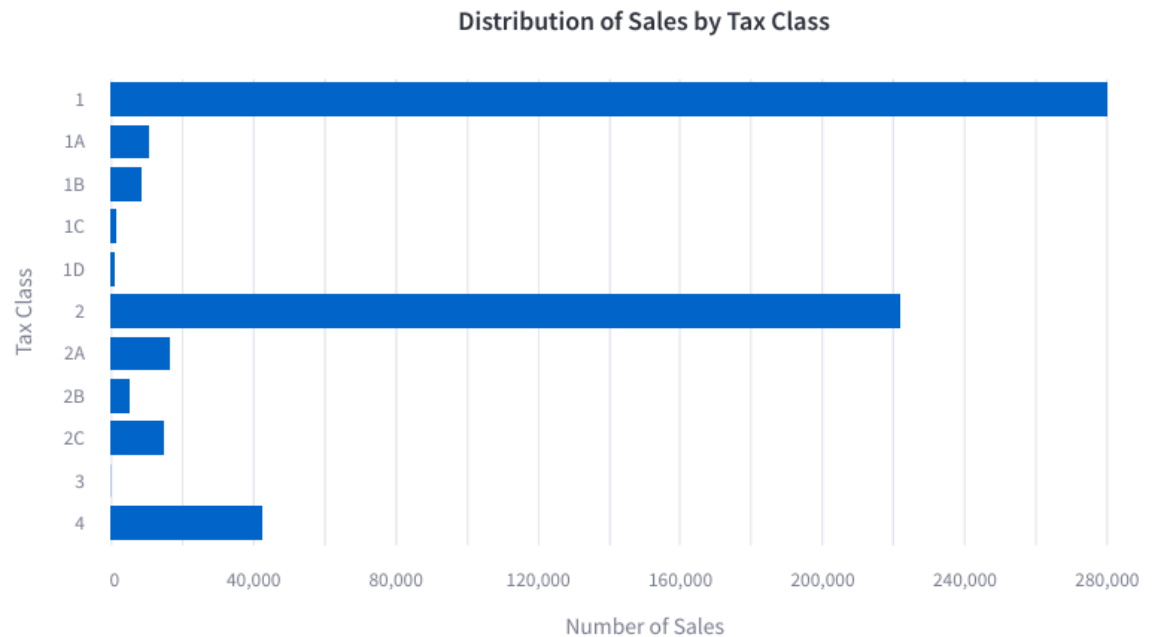
```

412 #9
413 import streamlit as st
414 import pandas as pd
415 import altair as alt
416
417
418 st.subheader('Distribution of Property Sales by Tax Class')
419 # Filter for Date Range
420 start_date = df['SALE DATE'].min()
421 end_date = df['SALE DATE'].max()
422
423 # Filter data based on selected filters
424 filtered_df = df[(df['SALE DATE'] >= start_date) & (df['SALE DATE'] <= end_date)]
425
426 # Create a bar chart for distribution of sales by tax class
427 tax_class_counts = filtered_df['TAX CLASS AS OF FINAL ROLL'].value_counts().reset_index()
428 tax_class_counts.columns = ['Tax Class', 'Number of Sales']
429
430 # Plotting with Altair
431 bar_chart = alt.Chart(tax_class_counts).mark_bar().encode(
432     y=alt.Y('Tax Class:O', title='Tax Class'),
433     x=alt.X('Number of Sales:Q', title='Number of Sales'),
434     tooltip=['Tax Class', 'Number of Sales']
435 ).properties(
436     width=600,
437     height=400
438 ).interactive()
439
440 # Add chart title and Labels
441 bar_chart = bar_chart.properties(
442     title="Distribution of Sales by Tax Class"
443 ).configure_axis(
444     labelFontSize=12,
445     titleFontSize=14
446 ).configure_title(
447     fontSize=16,
448     anchor='middle'
449 )
450
451 # Display the chart
452 st.altair_chart(bar_chart, use_container_width=True)
453

```

This above code outlines the process of creating an interactive bar chart that visualizes the distribution of property sales by tax class in a Streamlit web application. The chart is generated using Altair and is based on data filtered by a user-defined sale date range. After filtering, the data is grouped by the 'TAX CLASS AS OF FINAL ROLL' column, and the count of sales for each tax class is computed. The resulting bar chart presents the number of sales per tax class, complete with tooltips for additional detail. It's designed for user interaction and enhanced readability with customized fonts and labels, providing an informative visual summary of sales distribution which can easily be incorporated into a report for stakeholders.

Distribution of Property Sales by Tax Class



The above visualization presents a horizontal bar chart titled "Distribution of Property Sales by Tax Class", which provides a visual breakdown of the number of property sales grouped by different tax classes. The chart displays a series of tax classes on the y-axis and the corresponding number of sales on the x-axis. Tax class '2' appears to have the highest number of sales, suggesting it represents a significant portion of the property market. The visualization simplifies complex data into an easily interpretable format, enabling quick insights into which tax classes are most common among sold properties, information that can be valuable for market analysis and reporting.

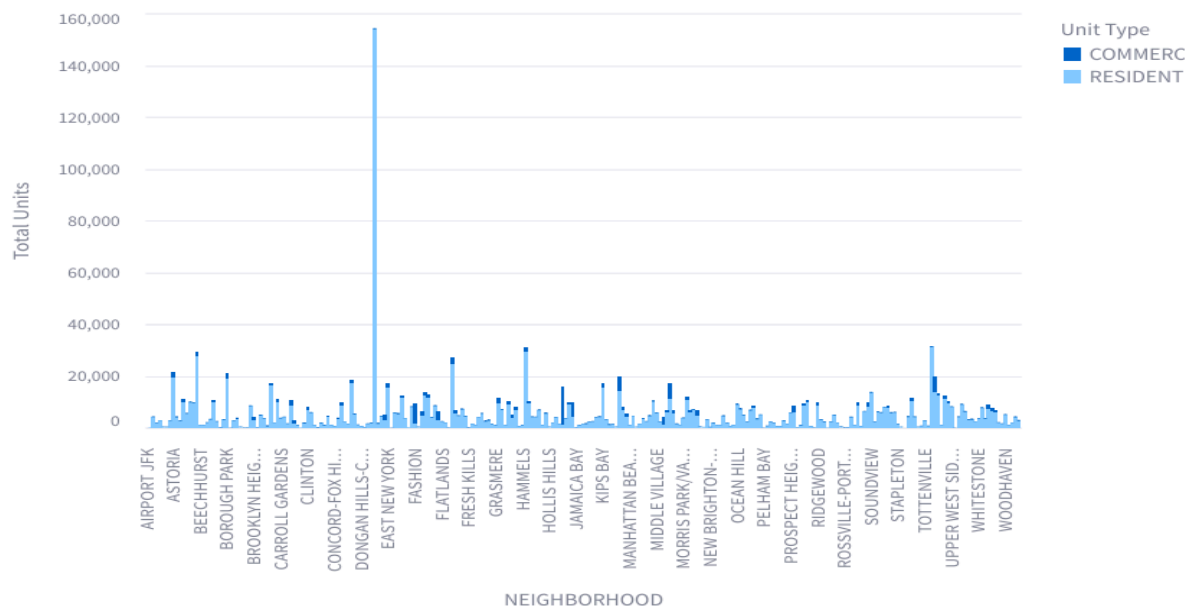
```

455 #10
456
457 st.subheader('Number of Residential and Commercial Units by Neighborhood')
458
459 # Group data by neighborhood and sum the residential and commercial units
460 units_by_neighborhood = df.groupby('NEIGHBORHOOD')[['RESIDENTIAL UNITS', 'COMMERCIAL UNITS']].sum().reset_index()
461
462 # Melt the DataFrame to Long format for easier plotting
463 units_by_neighborhood_melted = units_by_neighborhood.melt(id_vars='NEIGHBORHOOD', var_name='Unit Type', value_name='Total Units')
464
465 # Plot stacked bar chart
466 bar_chart = alt.Chart(units_by_neighborhood_melted).mark_bar().encode(
467     x='NEIGHBORHOOD:N',
468     y='Total Units:Q',
469     color='Unit Type:N',
470     tooltip=['NEIGHBORHOOD', 'Total Units', 'Unit Type']
471 ).properties(
472     width=800,
473     height=500
474 ).interactive()
475
476 # Display the chart
477 st.write(bar_chart)

```

The above code is set up to visualize the distribution of residential and commercial units across neighborhoods within a Streamlit application. It starts by grouping the data by neighborhood and summing up the residential and commercial units within each one. To prepare for visualization, it reshapes the grouped data into a long format suitable for plotting. An Altair stacked bar chart is then created, showing the total number of units for each neighborhood differentiated by unit type. The chart is configured to be interactive and is finally displayed in the Streamlit app, allowing for an intuitive comparison of residential and commercial properties across different neighborhoods.

Number of Residential and Commercial Units by Neighborhood



The above visualization depicts a bar chart visualizing the number of residential and commercial units within various neighborhoods. The bars are segmented by unit type, displaying a comparison between the two. This visualization clearly shows that residential units are more prevalent than commercial units across most neighborhoods. A particular neighborhood stands out with a significantly higher number of residential units compared to others. Such a chart is useful for identifying residential and commercial property distribution patterns, which can inform real estate investment strategies, urban planning, and market analysis.

Conclusion

In conclusion, the project report presented a detailed exploration of New York City's property sales data from 2016 to 2022, employing a variety of data processing and visualization techniques. Geocoding enriched the dataset with precise location coordinates, facilitating a geospatial analysis that led to interactive maps showcasing sales distributions across neighborhoods. Further, time-series analyses highlighted market dynamics over time, while data cleaning and preparation ensured the reliability of the subsequent analyses. Visualizations, including box plots and pie charts, illustrated the disparities in property features such as unit types, square footage, and the association of these features with sale prices.

This comprehensive analysis provided insights into the NYC real estate market's complexity, underlined by residential dominance over commercial units and a clear preference for certain neighborhoods and property types. The Streamlit web application's interactive nature allowed users to engage directly with the data, ensuring the project's findings were both accessible and actionable. The data-driven approach laid a strong foundation for future market predictions and investment strategies within the ever-changing landscape of New York City real estate.