# Airbnb Price Predictive model
# CIS 512
# Prof. Sumanlata Ghosh

Presented by: **Anand Vekariya**

# Data Cleaning

# Drop column

| | | |
|---|---|---|
| Listing Url,<br>Scrape ID,<br>Name,<br>Summary,<br>Space,<br>Description,<br>Experiences Offered,<br>Neighborhood Overview,<br>Notes,<br>Access,<br>Interaction,<br>House Rules,<br>Thumbnail Url,<br>Medium Url,<br>Picture Url,<br>XL Picture Url,<br>Host URL, | Host Name,<br>Host About,<br>Host Acceptance Rate,<br>Host Thumbnail Url,<br>Host Picture Url,<br>Host Neighbourhood,<br>Host Listings Count,<br>Host Total Listings Count,<br>Host Verifications,<br>Neighbourhood,<br>Neighbourhood Cleansed,<br>Neighbourhood Group Cleansed,<br>State,<br>Zipcode,<br>Smart Location,<br>Country Code,<br>Geolocation,<br>Cancellation Policy | Latitude,<br>Longitude,<br>Square Feet,<br>Has Availability,<br>Availability 30,<br>Availability 60,<br>Availability 90,<br>Availability 365,<br>Calendar last Scraped,<br>First Review,<br>Review Scores Accuracy,<br>Review Scores Cleanliness,<br>Review Scores Checkin,<br>Review Scores Communication,<br>Review Scores Location,<br>Review Scores Value,<br>Jurisdiction Names |

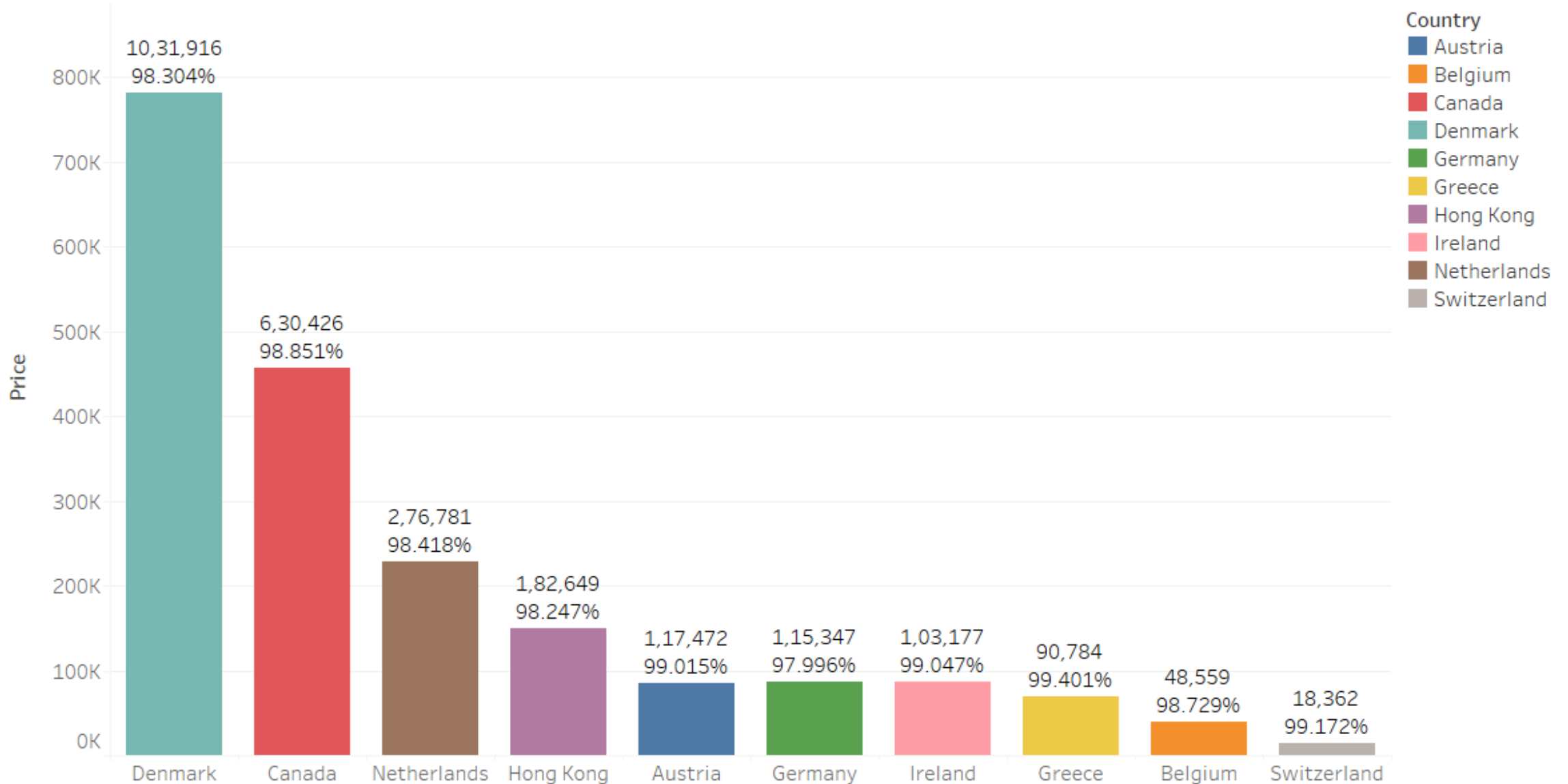| Columns | Data Description | Action Taken |
|---------|------------------|--------------|
| Amenities | Multiple facilities and requirements | Split into categories; new column = Yes/No. |
| Price | Contains NA values. | Check weekly/monthly values; calculate mean price based on country. |
| Bedrooms, Beds, Bathrooms | Contains NA values | Use 'Accommodates' column to find actual numbers; otherwise, set to 0. |
| Host Location Match | Host location matches property location | Host Location (split) = "Street" or "City." |
| Last Rented (in months) | Findout Last time rented. | 'Last Rented (in months)' = 'Last Scraped' - 'Last Review.' |
| Host Age | Determine host age. | 'Host Age' = 'Last Scraped' - 'Host Since.' |

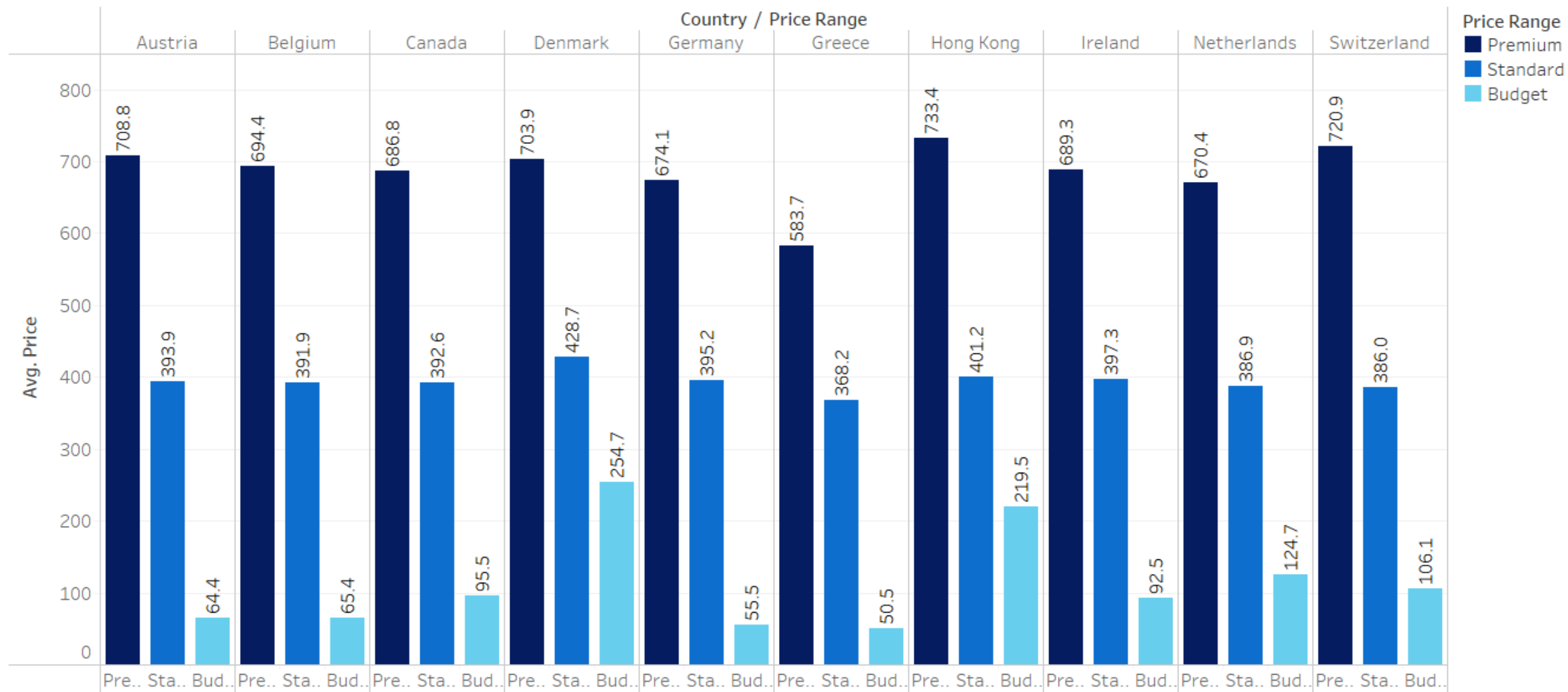| Columns | Data Description | Action Taken |
|---|---|---|
| Property Type | Unorganized property types. | Categorize based on house types. |
| Features | Long strings with essential details. | Split features; create new column with Yes/No. |
| Transit, Market, Weekly Price, Monthly Price, License | Long strings with mostly null values. | If empty, set to 'No'; else, set to 'Yes.' |
| Host Response Time, Beds, Security Deposit, Cleaning Fee, Review Scores Rating, Extra People, Reviews per Month | Contains NA values. | Fill with 0 for int/float; else, set to "Not Mention." |

# Data Analysis
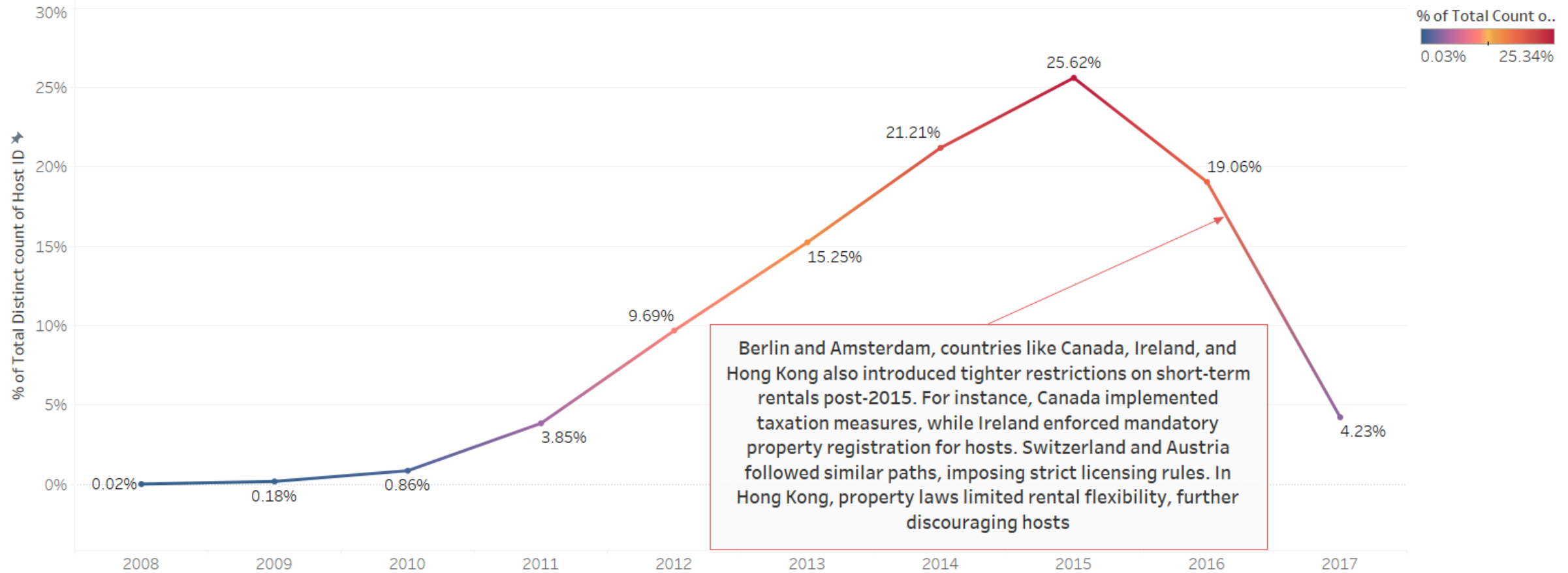
# Total Revenue by Country



Shows which countries generate the highest revenue from Airbnb rentals, emphasizing quality hosts (Superhosts).
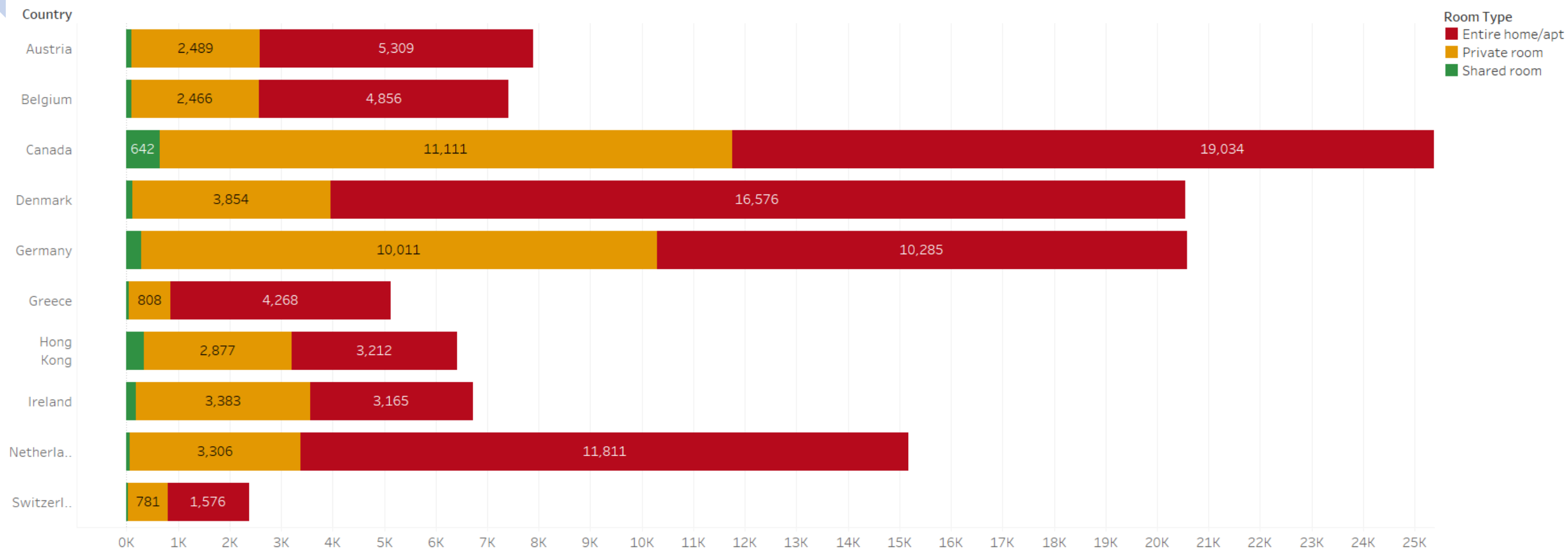
# Average Airbnb Prices Across Countries



It compares average Airbnb prices across countries by Budget, Standard, and Premium categories, highlighting Denmark and Hong Kong as the most expensive, while Greece and Germany offer more affordable options.

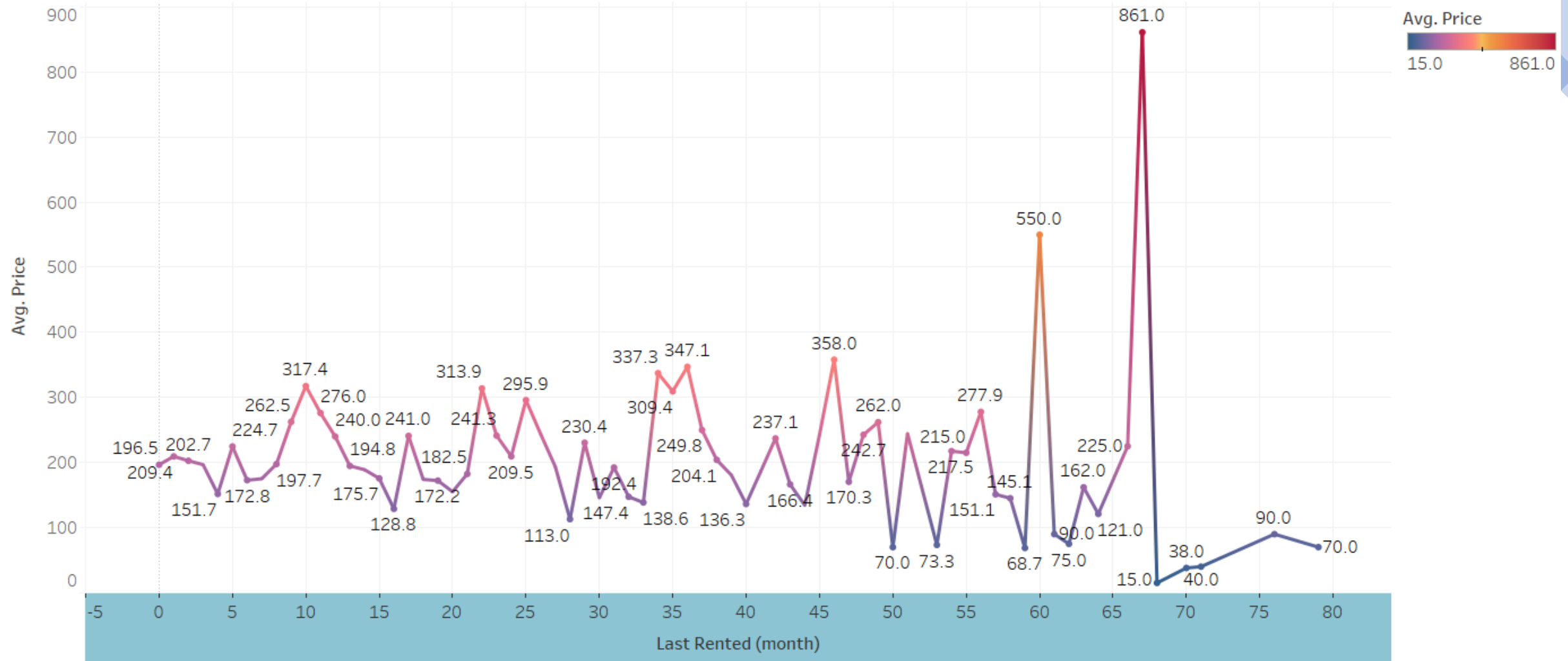# Yearly Host Registration Rates on Airbnb



Berlin and Amsterdam, countries like Canada, Ireland, and Hong Kong also introduced tighter restrictions on short-term rentals post-2015. For instance, Canada implemented taxation measures, while Ireland enforced mandatory property registration for hosts. Switzerland and Austria followed similar paths, imposing strict licensing rules. In Hong Kong, property laws limited rental flexibility, further discouraging hosts

Host registrations began at 0.02% in 2008, gradually increasing to 25.63% by 2015. Following this peak, a significant drop occurred, with registrations falling to 4.23% by 2017. The rapid growth between 2011 and 2015 contrasts sharply with the steep decline in the subsequent years.

# Room Type Distribution by Country



**Country**

**Room Type**
- Entire home/apt
- Private room
- Shared room

| Country | Private room | Entire home/apt |
|---------|-------------|-----------------|
| Austria | 2,489 | 5,309 |
| Belgium | 2,466 | 4,856 |
| Canada | 642 / 11,111 | 19,034 |
| Denmark | 3,854 | 16,576 |
| Germany | 10,011 | 10,285 |
| Greece | 808 | 4,268 |
| Hong Kong | 2,877 | 3,212 |
| Ireland | 3,383 | 3,165 |
| Netherla.. | 3,306 | 11,811 |
| Switzerl.. | 781 | 1,576 |

Shows the popularity of different room types in each country, helping investors understand which types of properties are more in demand.

Current Average Price X Number of Months Since Last Rental

This chart displays the current average price of Airbnb Property X, indicating the number of months since its last rental. It provides insights into how rental availability may influence pricing.
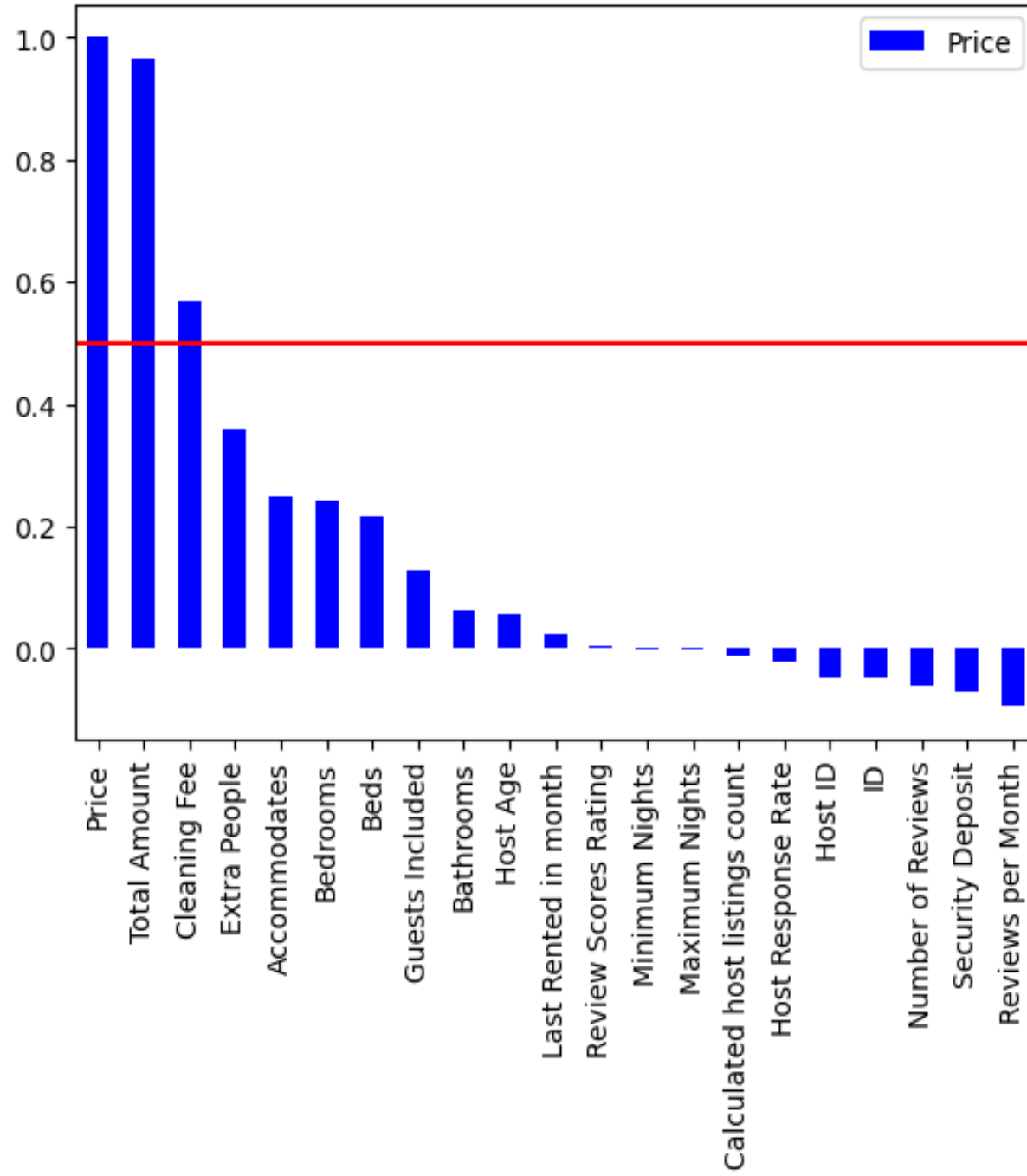
# Data Correlation

Correlation of Price with Other Features

| | ID | Host ID | Host Response Rate | Accommodates | Bathrooms | Bedrooms | Beds | Price | Security Deposit | Cleaning Fee | Guests Included | Extra People | Minimum Nights | Maximum Nights | Number of Reviews | Review Scores Rating | Calculated host listings count | Reviews per Month | Host Age | Last Rented in month | Total Amount |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Price | -0.05 | -0.05 | -0.02 | 0.25 | 0.06 | 0.24 | 0.22 | 1.00 | -0.07 | 0.57 | 0.13 | 0.36 | -0.00 | -0.00 | -0.06 | 0.00 | -0.01 | -0.09 | 0.06 | 0.02 | 0.96 |

Correlation with Price

# Feature Engineering

# Feature Engineering Process:

| Feature Engineering Process: | Categorical Variables Encoded: | Example: | Resulting Dataframe Shape: |
|---|---|---|---|
| • Converted categorical variables to numerical using one-hot encoding<br>• Created new columns for each category:<br>• 1 indicates presence 0 indicates absence<br>• Original categorical columns dropped | • **Host Response Time**<br>• **Country**<br>• **Property Type**<br>• **Room Type**<br>• **Bed Type**<br>• **Cancellation Policy** | • Before: Room Type = "Entire home/apt"<br>• After: Room Type_Entire home/apt = 1, others = 0 | • Rows: 123,061<br>• Columns: 91 |

# Preparing Data for Modeling

**Correlation Filtering:**

Removed columns with correlation < 0.02 to price

**Data Splitting:**

Training Data: 75% (~92,295 rows)

Testing/Validation Data: 25% (~30,766 rows)

**Data Standardization:**

Applied StandardScaler to features

Standardizes features by removing the mean and scaling to unit variance

**Dimensionality Reduction:**

Performed Principal Component Analysis (PCA)

Reduced feature set to 57 columns for modeling

Predictive Modeling

# Linear Regression



Histogram of Actual vs Predicted (Testing Data)

Histogram of Actual vs Predicted (Validation Data)

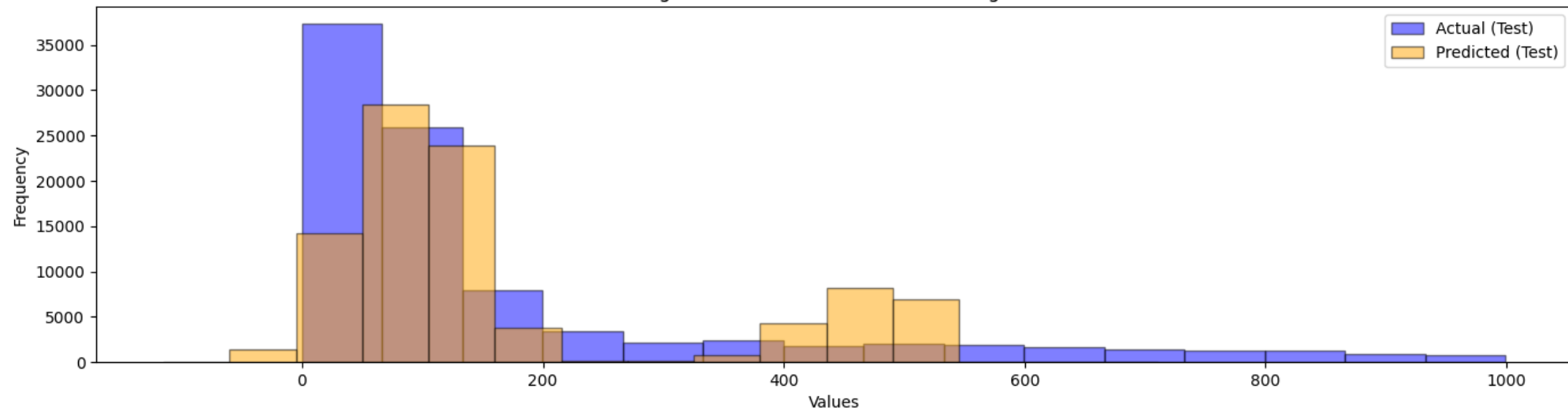# Random Forest Regression



Histogram of Actual vs Predicted (Testing Data)

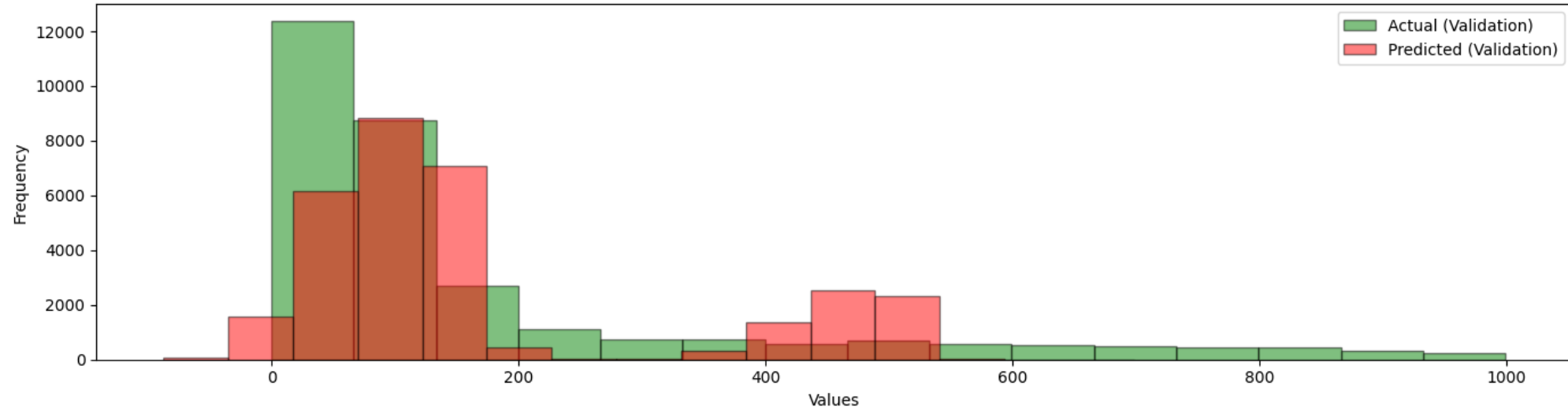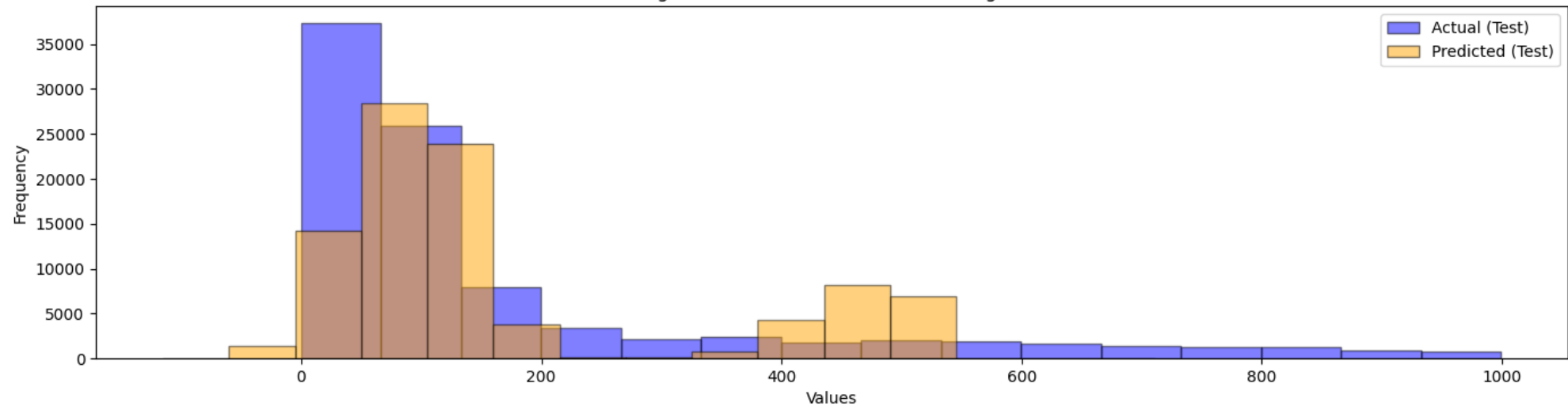Histogram of Actual vs Predicted (Validation Data)

# Lasso Regression



Histogram of Actual vs Predicted (Testing Data)
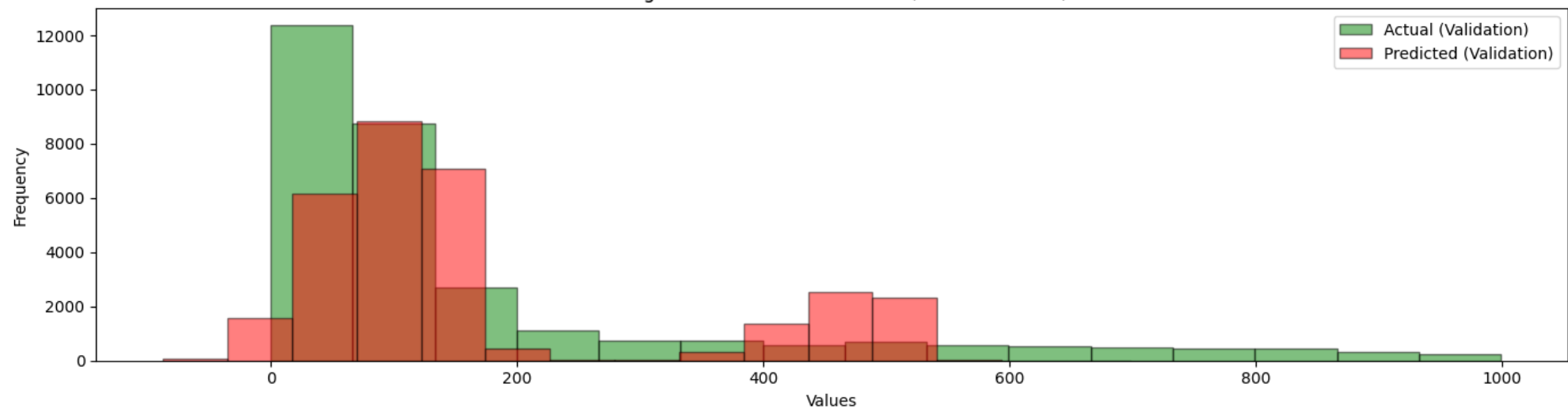
Histogram of Actual vs Predicted (Validation Data)

# Ridge Regression



Histogram of Actual vs Predicted (Testing Data)

Histogram of Actual vs Predicted (Validation Data)

# Model Performance Comparison

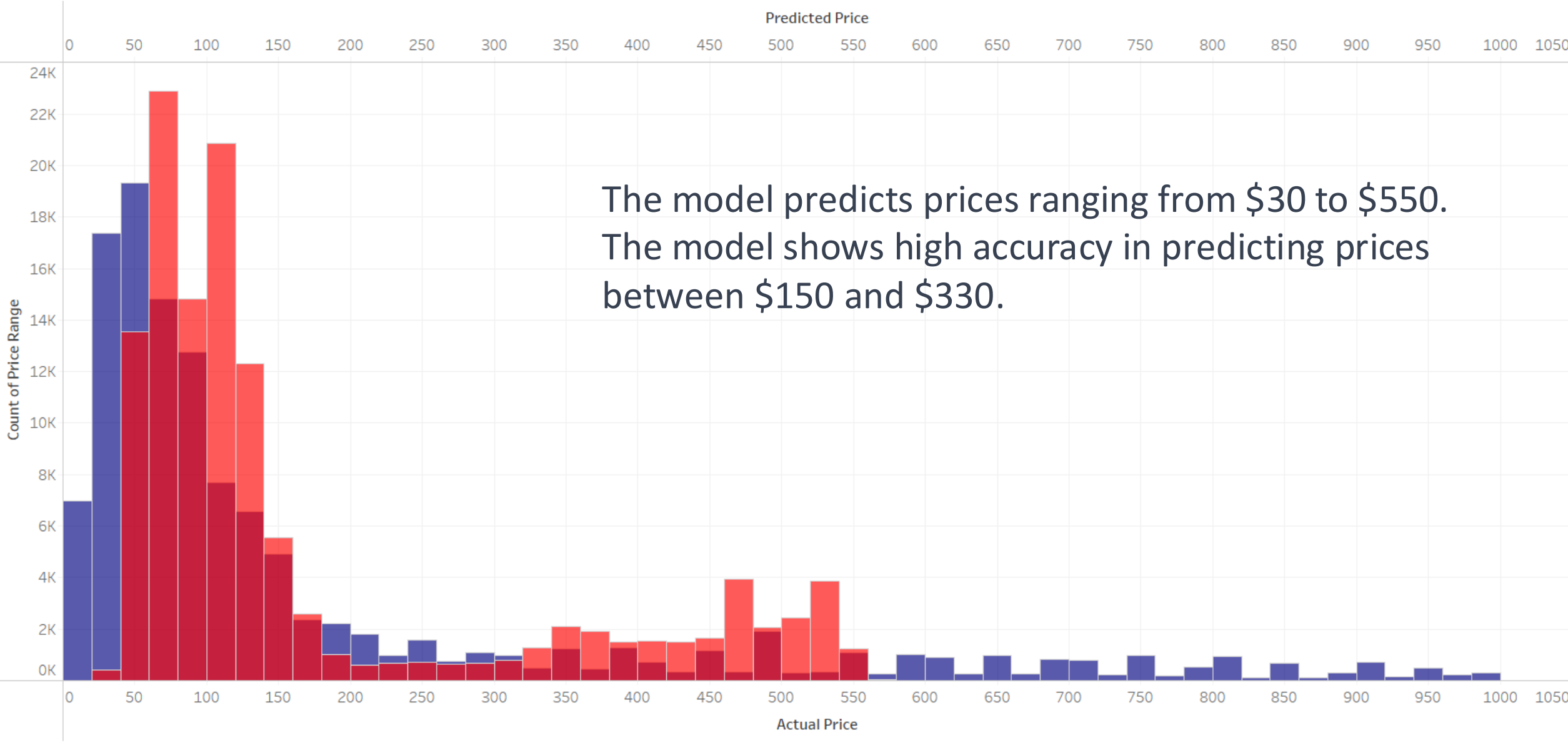| | Model | $R^2$ | Adjusted $R^2$ | Mean Squared Error (MSE) | Root Mean Square Error (RMSE) |
|---|---|---|---|---|---|
| **Train** | Linear Regression | 0.55 | 0.55 | 21416.28 | 146.34 |
| | Random Forest Regression | 0.61 | 0.61 | 18478.35 | 135.94 |
| | Lasso Regression | 0.55 | 0.55 | 21416.22 | 146.34 |
| | Ridge Ridge Regression | 0.55 | 0.55 | 21416.2 | 146.34 |
| | Model | $R^2$ | Adjusted $R^2$ | Mean Squared Error (MSE) | Root Mean Square Error (RMSE) |
| **Validation** | Linear Regression | 0.56 | 0.56 | 20997.67 | 144.91 |
| | Random Forest Regression | 0.6 | 0.6 | 18809.03 | 137.15 |
| | Lasso Regression | 0.56 | 0.56 | 20997.88 | 144.91 |
| | Ridge Ridge Regression | 0.56 | 0.56 | 20997.76 | 144.91 |

# Winning Model: Random Forest Regression

- Key Insights:
    - **Superior Accuracy**: Outperforms other models in both training and validation sets
    - **Higher Explanatory Power**: $R^2$ of **0.60-0.61 vs 0.55-0.56** for other models
    - **Lower Prediction Errors: ~12% lower MSE and ~6% lower RMSE** compared to other models
    - **Consistent Performance:** Minimal difference between training and validation metrics, indicating good generalization

- Conclusion:
    - Random Forest Regression demonstrates superior performance in:
        - Capturing complex data patterns
        - Making accurate predictions
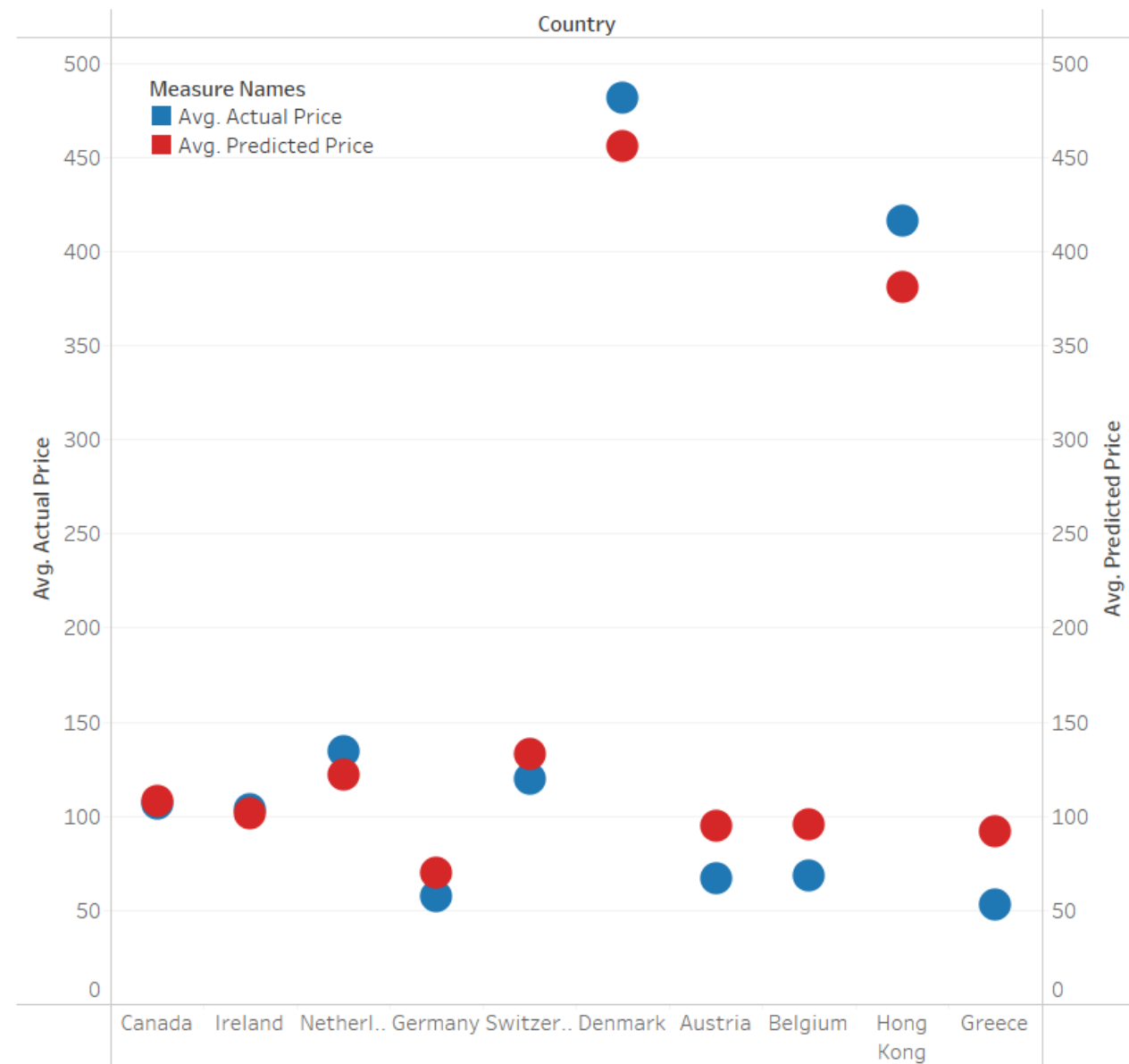        - Generalizing well to unseen data

Actual V/S Pridicated Price Histogram

The model predicts prices ranging from $30 to $550. The model shows high accuracy in predicting prices between $150 and $330.

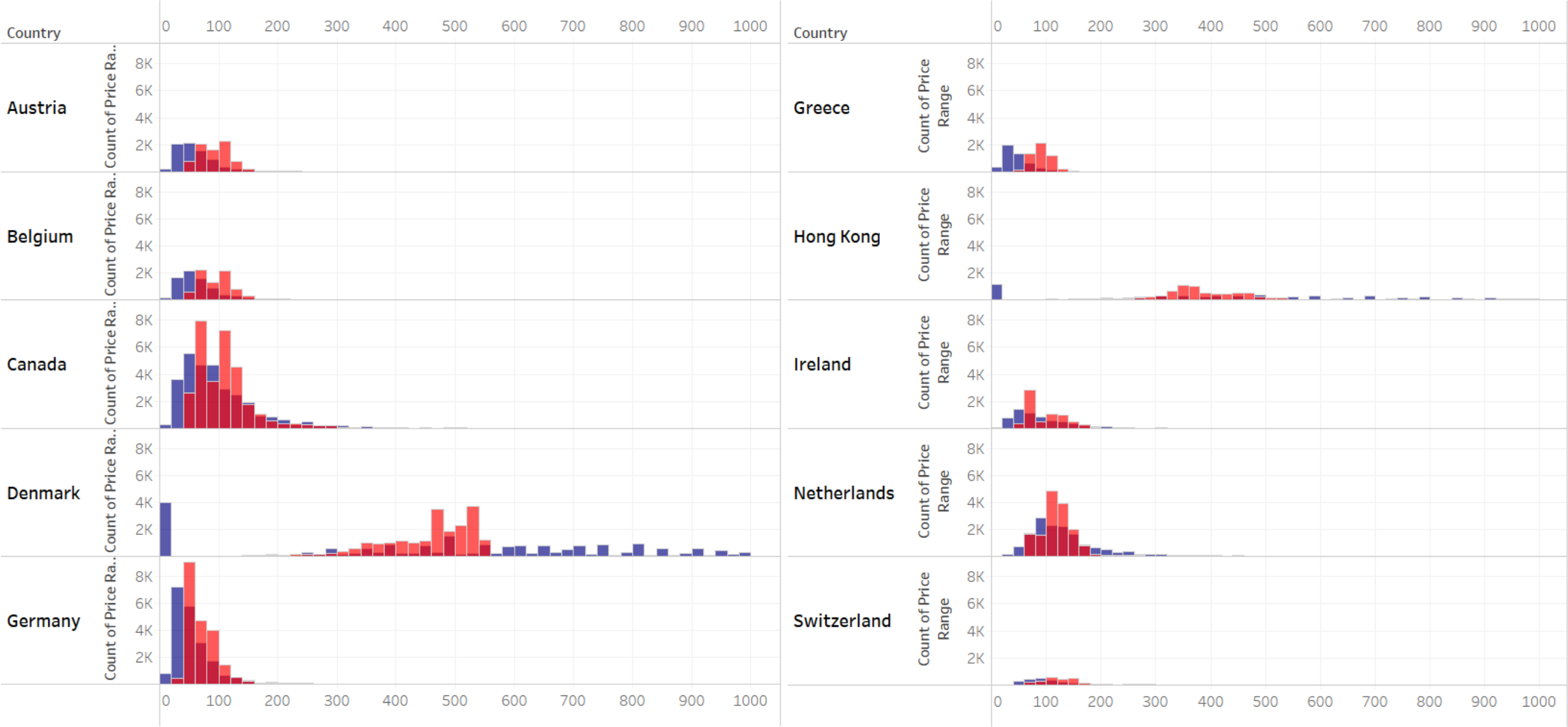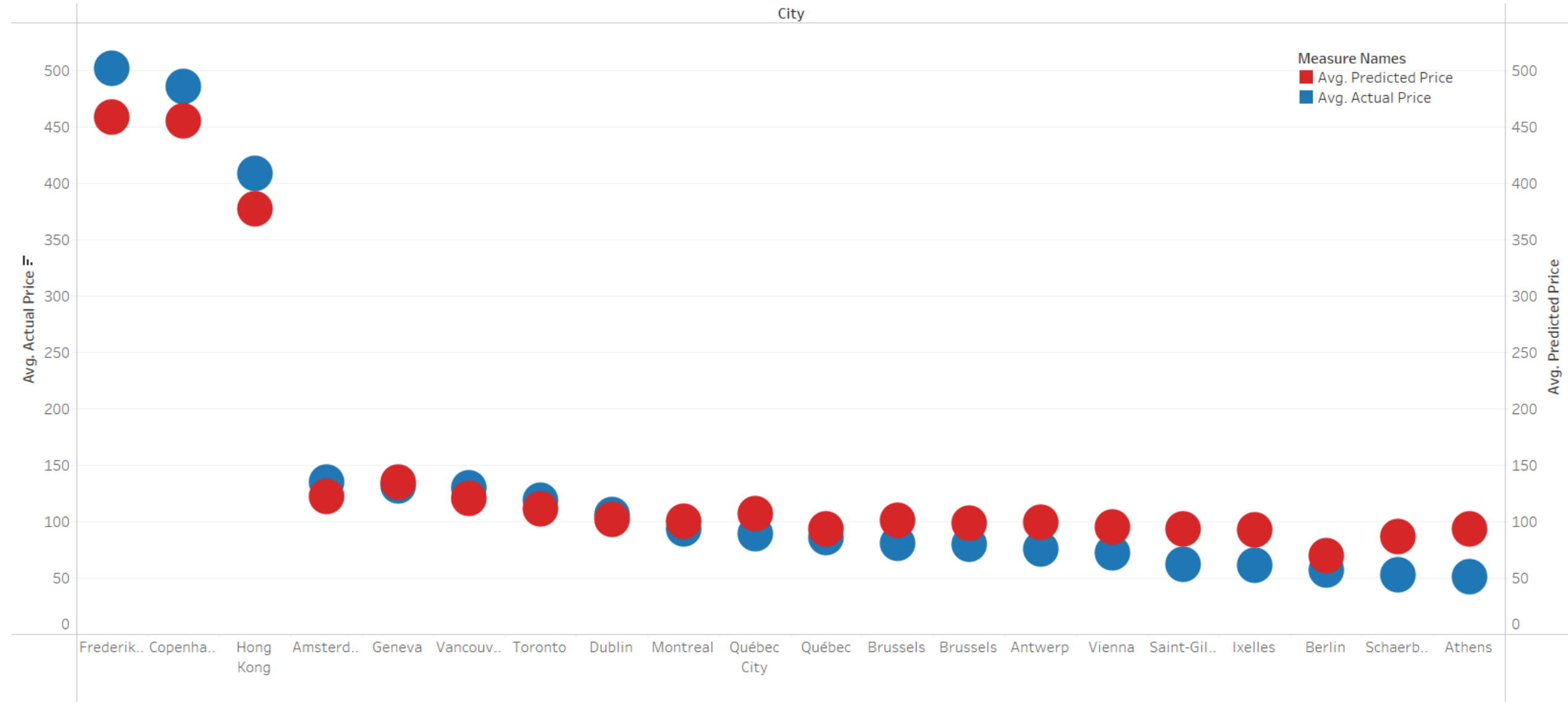## Actual V/S Pridicated Price by Contry

## Number Of Listing by Contry

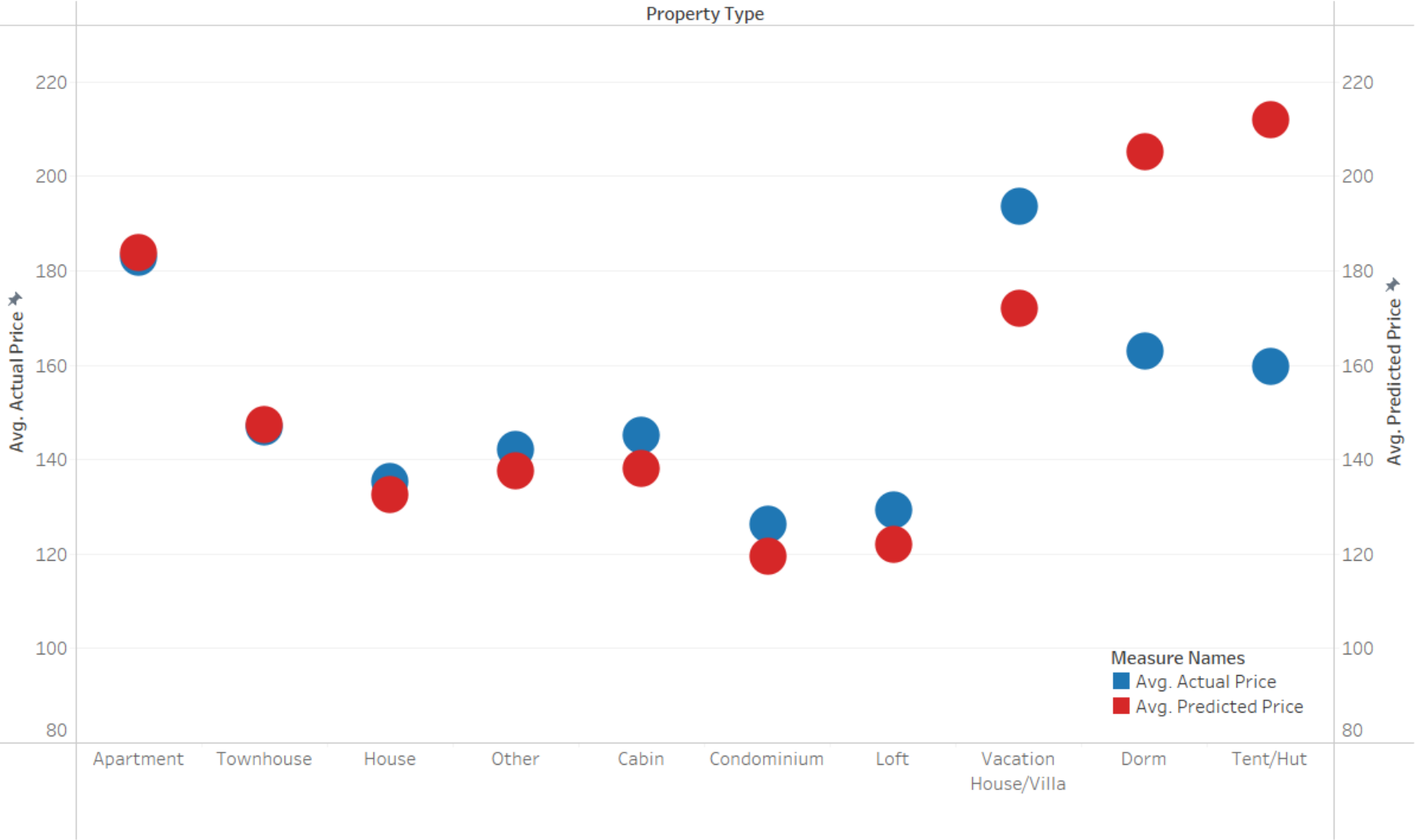Actual V/S Pridicated Price Histogram by Contry in Detail

# Actual V/S Pridicated Price by Top 20 City



In the graph, cities like **Geneva, Toronto, Dublin, Montreal, Québec** and **Berlin** show minimal differences between the actual and predicted average prices, indicating accurate predictions for these locations.

# Actual V/S Pridicated Price by Property Type

## Number Of Listing by Property Type



**Measure Names**
- ■ Avg. Actual Price
- ■ Avg. Predicted Price

| Property Type | |
|---|---|
| Apartment | 97,592 |
| House | 16,210 |
| Condominium | 3,879 |
| Loft | 1,603 |
| Townhouse | 1,219 |
| Vacation House/Villa | 954 |
| Other | 649 |
| Dorm | 544 |
| Cabin | 388 |
| Tent/Hut | 23 |

In the chart, the **predicted prices** for **Apartment** and **Townhouse** closely match their **actual prices**, indicating accurate predictions for these property types. For **House**, the predicted price is also quite close to the actual price, showing a high level of prediction accuracy.

# Conclusion

- **Canada**, **Netherlands**, and **Germany** are the top choices for investment, offering reliable price predictions and stable markets, particularly in **Toronto**, **Amsterdam**, and **Berlin**.

- Focus on major cities like **Toronto**, **Vancouver**, **Amsterdam**, **Rotterdam**, **Berlin**, and **Munich** where demand and price accuracy are high.

- Invest in **apartments** and **townhouses** for reliable returns. **Houses** are viable, but may have slight discrepancies in price predictions.

# Any Questions?