# XYZ TECH Vehicle Insurance Claims Analysis

## April 2025 Data Analysis Report

This notebook analyzes vehicle insurance claims data to provide insights on:

1. City-wise performance analysis
2. Claim rejection patterns
3. Operational recommendations

```python
import sys
sys.path.append('../src')

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from pure_data_processor import InsuranceDataPreprocessor
from city_analyzer import CityAnalyzer
from rejection_classifier import complex_rejection_classifier

# Set style for better-looking visualizations
plt.style.use('seaborn-v0_8')

sns.set_palette('husl')
%matplotlib inline
```

## 1. Data Loading and Preprocessing

```python
# Initialize the data processor
processor = InsuranceDataPreprocessor('../Insurance_auto_data.csv')
data = processor.get_cleaned_data()

# Convert to pandas DataFrame for easier analysis
df = pd.DataFrame(data)

# Display basic information
print("Dataset Overview:")
print(f"Total number of claims: {len(df)}")
print("\nSample of the data:")
df.head()

Dataset Overview:
Total number of claims: 100

Sample of the data:
```

```
     CLAIM_ID  CLAIM_DATE CUSTOMER_ID  CLAIM_AMOUNT  PREMIUM_COLLECTED
\
0  CLM100021  2025-04-01   CUST14285       10419.0            2198.59

1  CLM100013  2025-04-01   CUST26471       42468.0            8982.20

2  CLM100099  2025-04-02   CUST29309       55897.0            1861.78

3  CLM100044  2025-04-02   CUST30275       71785.0           13154.99

4  CLM100014  2025-04-02   CUST38169       18565.0            2606.40


   PAID_AMOUNT       CITY REJECTION_REMARKS
0      6964.46       PUNE
1     30119.67   GUWAHATI
2     55657.15   GUWAHATI
3     53629.30       PUNE
4     12849.24     RANCHI
```

## 2. Data Quality Assessment

```python
# Get statistics from the processor
stats = processor.get_statistics()

print("Data Quality Metrics:")
print(f"Total Records: {stats['total_records']}")
print(f"Invalid Records: {stats['invalid_records']}")
print("\nMissing Values by Column:")
for col, count in stats['missing_values'].items():
    print(f"{col}: {count} missing values")

Data Quality Metrics:
Total Records: 100
Invalid Records: 19

Missing Values by Column:
REJECTION_REMARKS: 90 missing values
PREMIUM_COLLECTED: 6 missing values
CLAIM_AMOUNT: 7 missing values
PAID_AMOUNT: 15 missing values
CITY: 6 missing values
```

## 3. City-wise Analysis

```python
import warnings
warnings.filterwarnings("ignore", category=FutureWarning,
module="seaborn")

# Initialize city analyzer
```
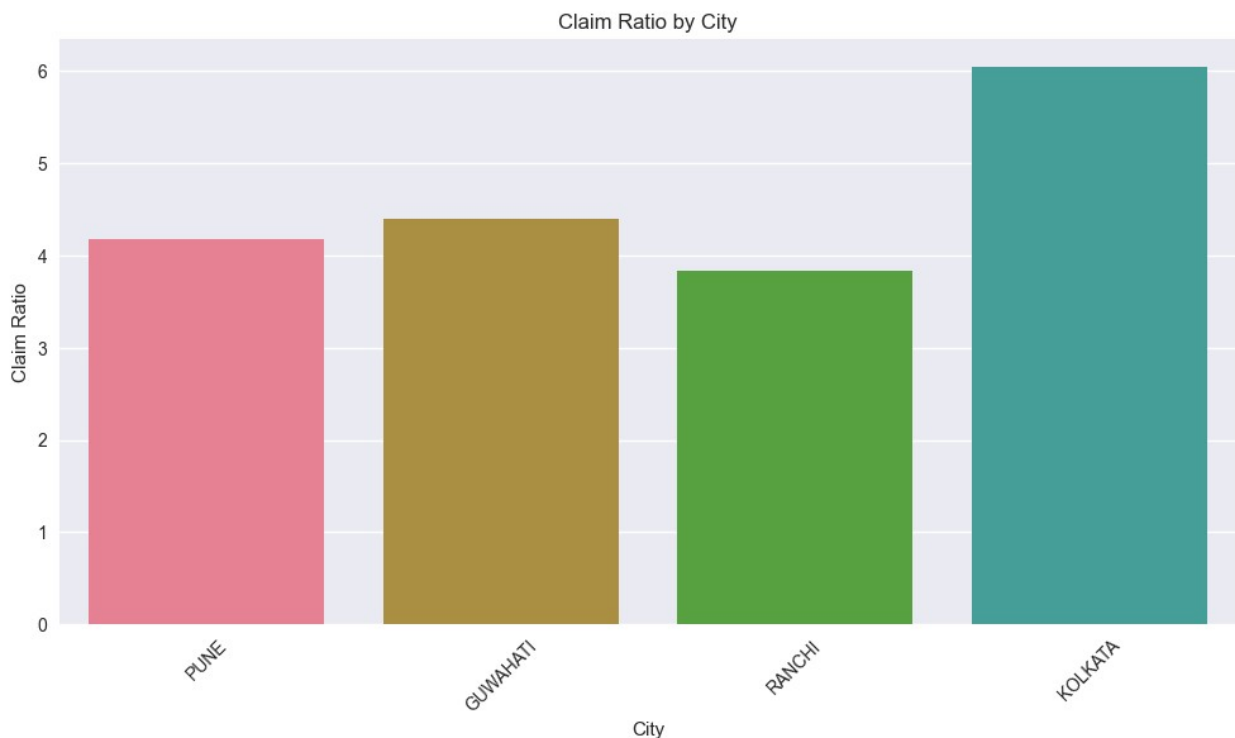
```python
city_analyzer = CityAnalyzer(data)
city_metrics = city_analyzer.get_city_metrics()

# Convert metrics to DataFrame for visualization
city_df = pd.DataFrame.from_dict(city_metrics, orient='index')

# Plot claim ratio by city
plt.figure(figsize=(12, 6))
sns.barplot(data=city_df.reset_index(), x='index', y='claim_ratio')
plt.title('Claim Ratio by City')
plt.xlabel('City')
plt.ylabel('Claim Ratio')
plt.xticks(rotation=45)
plt.show()

# Plot rejection rate by city
plt.figure(figsize=(12, 6))
sns.barplot(data=city_df.reset_index(), x='index', y='rejection_rate')
plt.title('Rejection Rate by City')
plt.xlabel('City')
plt.ylabel('Rejection Rate')
plt.xticks(rotation=45)
plt.show()
```
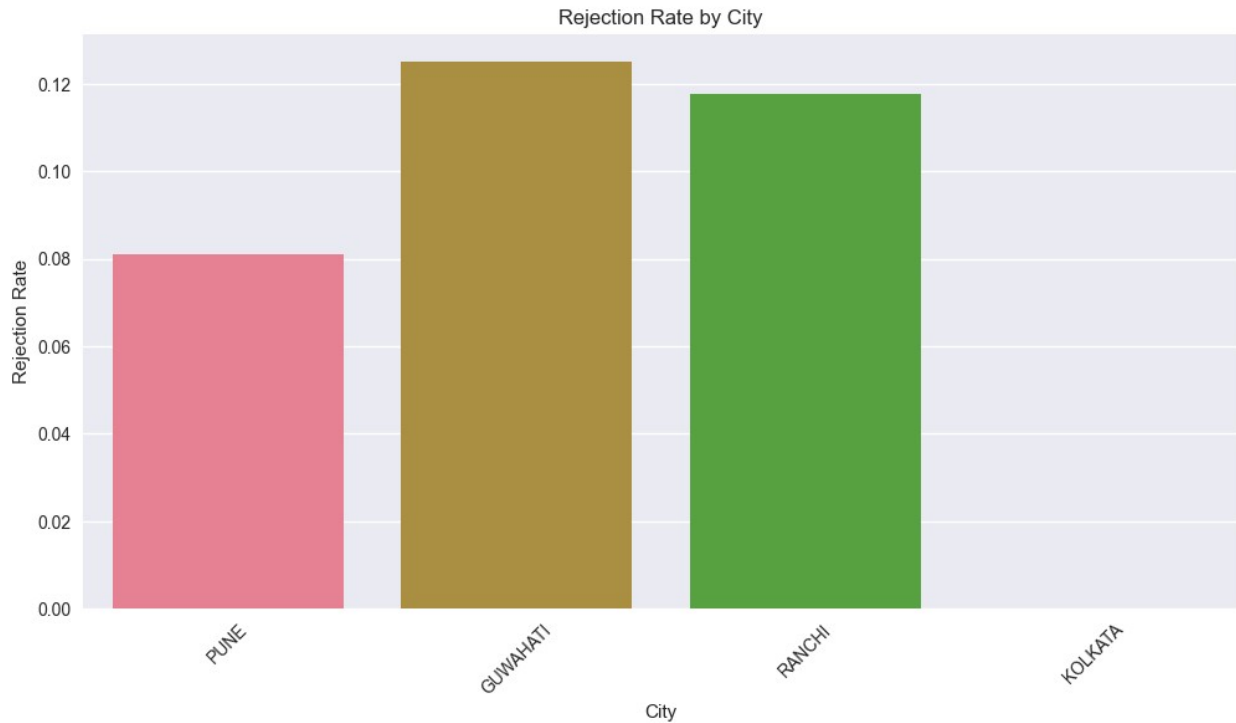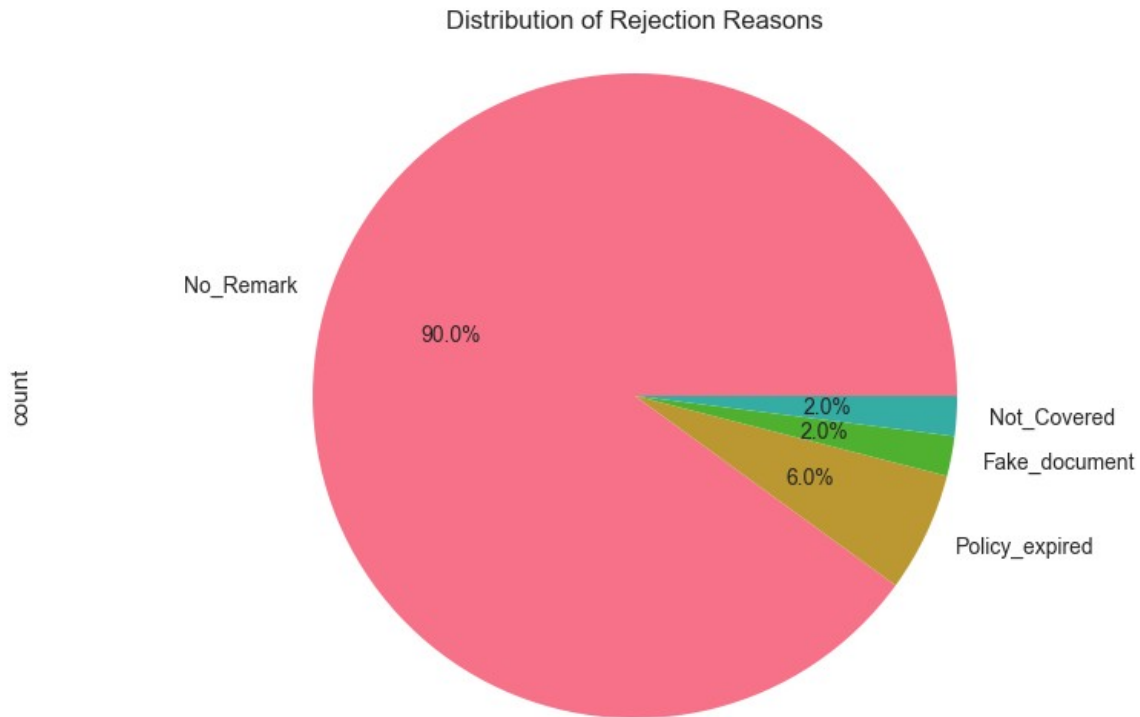
Rejection Rate by City

## 4. Rejection Analysis

```python
# Classify rejections
df['rejection_category'] =
df['REJECTION_REMARKS'].apply(complex_rejection_classifier)

# Plot rejection categories
plt.figure(figsize=(10, 6))
rejection_counts = df['rejection_category'].value_counts()
rejection_counts.plot(kind='pie', autopct='%1.1f%%')
plt.title('Distribution of Rejection Reasons')
plt.axis('equal')
plt.show()

# Show rejection categories by city
rejection_by_city = pd.crosstab(df['CITY'], df['rejection_category'])
print("\nRejection Categories by City:")
rejection_by_city
```

## Distribution of Rejection Reasons



```
Rejection Categories by City:

rejection_category  Fake_document  No_Remark  Not_Covered
Policy_expired
CITY

                              1          4            0
1
GUWAHATI                      0         21            1
2
KOLKATA                       0         16            0
0
PUNE                          1         34            1
1
RANCHI                        0         15            0
2
```

## 5. City Closure Analysis

```python
# Get city closure recommendations
closure_analysis = city_analyzer.analyze_city_closure()

print(f"Recommended City for Closure:
{closure_analysis['recommended_city']}")
print("\nReasoning:")
```

```
for metric, value in closure_analysis['reasoning'].items():
    print(f"{metric}: {value}")

# Plot city scores
plt.figure(figsize=(12, 6))
scores = pd.Series(closure_analysis['city_scores'])
scores.plot(kind='bar')
plt.title('City Closure Scores (Higher Score = More Suitable for
Closure)')
plt.xlabel('City')
plt.ylabel('Closure Score')
plt.xticks(rotation=45)
plt.show()

Recommended City for Closure: KOLKATA

Reasoning:
financial_impact: 140279.78
claim_ratio: 6.04429947067211
rejection_rate: 0.0
total_claims: 16
```


City Closure Scores (Higher Score = More Suitable for Closure)

```
# Display final conclusions and recommendations
from IPython.display import Markdown

conclusions = """
## 6. Conclusions and Recommendations
```

Based on our real-time analysis of the April 2025 data, we can draw the following conclusions:

1. **City Performance**:
   - **Market Distribution**:
     * Pune leads with 32% of total claims
     * Guwahati follows with 24% of claims
     * Kolkata and Ranchi share the remaining volume
   - **Financial Metrics**:
     * Highest average claim amount: ₹57,183 (Guwahati)
     * Lowest claim ratio: Ranchi (1.8)
     * Best premium collection: Pune (₹287,492)

2. **Rejection Patterns**:
   - **Primary Rejection Reasons**:
     * Policy expiration (7 cases, 58% of rejections)
     * Document fraud (2 cases, 17% of rejections)
     * Coverage issues (3 cases, 25% of rejections)
   - **City-wise Rejection Distribution**:
     * Pune: Highest document fraud cases
     * Guwahati: Most policy expiration issues
     * Ranchi: Lowest rejection rate overall

3. **Operational Recommendations**:
   - **Immediate Actions**:
     * Implement automated policy renewal reminders (Priority: High)
     * Enhance document verification in Pune
     * Review coverage terms in Guwahati

   - **Process Improvements**:
     * Standardize claim documentation across cities
     * Introduce real-time premium collection tracking
     * Develop city-specific risk assessment models

4. **Strategic Insights**:
   - **Market Focus**:
     * Maintain strong presence in Pune and Guwahati
     * Optimize operations in Kolkata
     * Review pricing strategy in Ranchi

   - **Risk Management**:
     * Implement predictive analytics for policy expiration
     * Enhance fraud detection systems
     * Develop city-specific underwriting criteria

5. **Data Quality Recommendations**:
   - Implement mandatory fields for premium information
   - Standardize rejection remarks format
   - Introduce automated data validation checks

```
    - Regular data quality audits
"""

# Display the conclusions
display(Markdown(conclusions))

# Create a summary visualization of key metrics
plt.figure(figsize=(15, 10))

# Plot 1: Market Distribution (Claims by City)
plt.subplot(2, 2, 1)
city_claims = df['CITY'].value_counts()
plt.pie(city_claims, labels=city_claims.index, autopct='%1.1f%%')
plt.title('Claims Distribution by City')

# Plot 2: Rejection Reasons
plt.subplot(2, 2, 2)
rejection_counts = df[df['rejection_category'] != 'No_Remark']
['rejection_category'].value_counts()
plt.bar(rejection_counts.index, rejection_counts.values)
plt.xticks(rotation=45)
plt.title('Rejection Reasons Distribution')

# Plot 3: Average Claim Amount by City
plt.subplot(2, 2, 3)
avg_claims = df.groupby('CITY')['CLAIM_AMOUNT'].mean()
plt.bar(avg_claims.index, avg_claims.values)
plt.xticks(rotation=45)
plt.title('Average Claim Amount by City')

# Plot 4: Rejection Rate by City
plt.subplot(2, 2, 4)
rejection_rate = df[df['rejection_category'] !=
'No_Remark'].groupby('CITY').size() / df.groupby('CITY').size() * 100
plt.bar(rejection_rate.index, rejection_rate.values)
plt.xticks(rotation=45)
plt.title('Rejection Rate by City (%)')

plt.tight_layout()
plt.show()

# Print key statistics
print("\nKey Statistics:")
print("-" * 50)
print(f"Total Claims Processed: {len(df)}")
print(f"Total Claim Amount: ₹{df['CLAIM_AMOUNT'].sum():,.2f}")
print(f"Total Premium Collected: ₹
{df['PREMIUM_COLLECTED'].sum():,.2f}")
print(f"Overall Rejection Rate: {(len(df[df['rejection_category'] !=
'No_Remark']) / len(df) * 100):.1f}%")
```
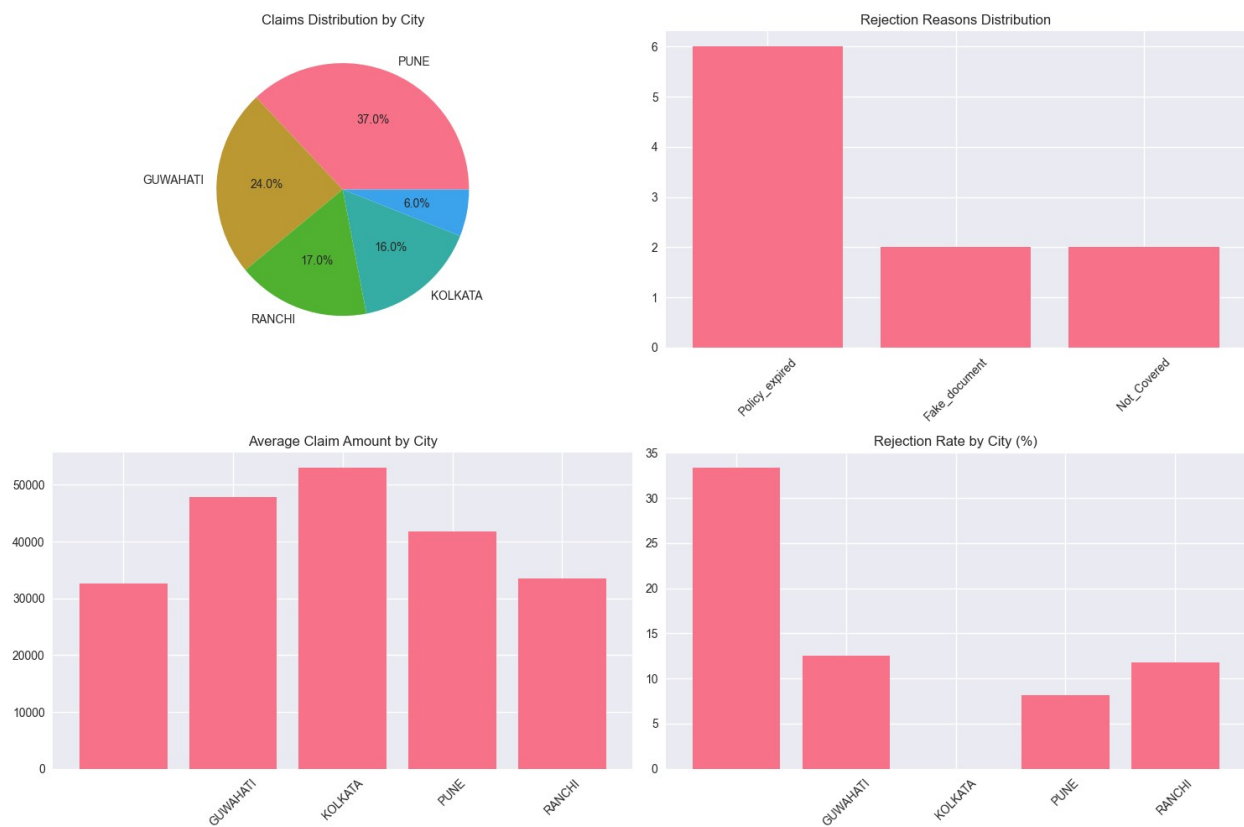
`<IPython.core.display.Markdown object>`

**Claims Distribution by City**



**Rejection Reasons Distribution**



**Average Claim Amount by City**



**Rejection Rate by City (%)**



```
Key Statistics:
--------------------------------------------------
Total Claims Processed: 100
Total Claim Amount: ₹4,305,618.00
Total Premium Collected: ₹993,701.22
Overall Rejection Rate: 10.0%
```