

Insurance Pricing Analysis

June 14, 2025

1 Insurance Pricing Adequacy & Risk Segmentation Analysis

Author: Dustin Corbett

Data Source: [Insurance Claims & Policy Data – Kaggle](#)

1.1 Project Summary

This actuarial pricing analysis uses synthetic insurance policyholder data to evaluate whether premiums appropriately reflect modeled risk. Using Generalized Linear Models (GLMs), we estimate claim frequency and severity, calculate the expected pure premium, and simulate repricing scenarios to correct underpricing.

Goals: 1. Model claim frequency using Poisson regression 2. Model claim severity using Gamma regression 3. Estimate pure premiums (frequency \times severity) 4. Identify segments with systemic over/underpricing 5. Simulate repricing adjustments, capped at 20% for retention realism 6. Deliver actuarially sound business recommendations

1.2 Load and Preview Data

We begin by loading the synthetic policyholder dataset and inspecting its structure.

```
[1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split
import statsmodels.api as sm
from statsmodels.genmod.families import Poisson, Gamma
from statsmodels.genmod.families.links import log
```

```
[2]: df = pd.read_csv("data/data_synthetic.csv")
df.head()
```

```
[2]:   Customer ID  Age  Gender  Marital Status  Occupation  Income Level  \
0         84966   23  Female         Married  Entrepreneur         70541
1         95568   26   Male         Widowed      Manager         54168
2         10544   29  Female         Single  Entrepreneur         73899
3         77033   20   Male         Divorced  Entrepreneur         63381
```

4	88160	25	Female	Separated	Manager	38794
---	-------	----	--------	-----------	---------	-------

	Education Level	Geographic Information	Location	Behavioral Data	...	\
0	Associate Degree	Mizoram	37534	policy5	...	
1	Doctorate	Goa	63304	policy5	...	
2	Associate Degree	Rajasthan	53174	policy5	...	
3	Bachelor's Degree	Sikkim	22803	policy5	...	
4	Bachelor's Degree	West Bengal	92858	policy1	...	

	Customer Preferences	Preferred Communication Channel	Preferred Contact Time	\
0	Email	In-Person Meeting	Afternoon	
1	Mail	In-Person Meeting	Morning	
2	Email	Mail	Evening	
3	Text	In-Person Meeting	Anytime	
4	Email	Text	Weekends	

	Preferred Language	Risk Profile	Previous Claims History	Credit Score	\
0	English	1	3	728	
1	French	1	2	792	
2	German	2	1	719	
3	French	3	0	639	
4	English	0	3	720	

	Driving Record	Life Events	Segmentation Group
0	DUI	Job Change	Segment5
1	Clean	Retirement	Segment5
2	Accident	Childbirth	Segment3
3	DUI	Job Change	Segment3
4	Major Violations	Childbirth	Segment2

[5 rows x 30 columns]

1.3 Clean and Engineer Features

We prepare data for actuarial modeling: - Parse dates - Rename for clarity - Create Tenure, Claim Flag, and Credit Category - Engineer affordability ratio: Premium to Income Ratio

```
[3]: # Parse dates
df['Policy Start Date'] = pd.to_datetime(df['Policy Start Date'])
df['Policy Renewal Date'] = pd.to_datetime(df['Policy Renewal Date'])

# Rename columns for modeling clarity
df.rename(columns={
    'Income Level': 'Income',
    'Education Level': 'Education',
    'Geographic Information': 'State',
    'Behavioral Data': 'Policy Category',
```

```

    'Interactions with Customer Service': 'Support Channel',
    'Insurance Products Owned': 'Policy Product',
    'Coverage Amount': 'Coverage',
    'Premium Amount': 'Premium',
    'Previous Claims History': 'Past Claims',
    'Credit Score': 'Credit',
    'Driving Record': 'Driving',
    'Life Events': 'Life Events',
    'Segmentation Group': 'Segment'
}, inplace=True)

# Feature Engineering
df['Tenure (Days)'] = (df['Policy Renewal Date'] - df['Policy Start Date']).dt.
    ↪days
df['Premium to Income Ratio'] = df['Premium'] / df['Income']
df['Claim Flag'] = (df['Claim History'] > 0).astype(int)

# Credit category binning
def categorize_credit(score):
    if score < 580: return 'Poor'
    elif score < 670: return 'Fair'
    elif score < 740: return 'Good'
    elif score < 800: return 'Very Good'
    else: return 'Excellent'
df['Credit Category'] = df['Credit'].apply(categorize_credit)

```

1.4 Exploratory Data Analysis (EDA)

Before modeling, we explore key distributions and relationships in the data to better understand how customer and policy characteristics may influence premiums and claims.

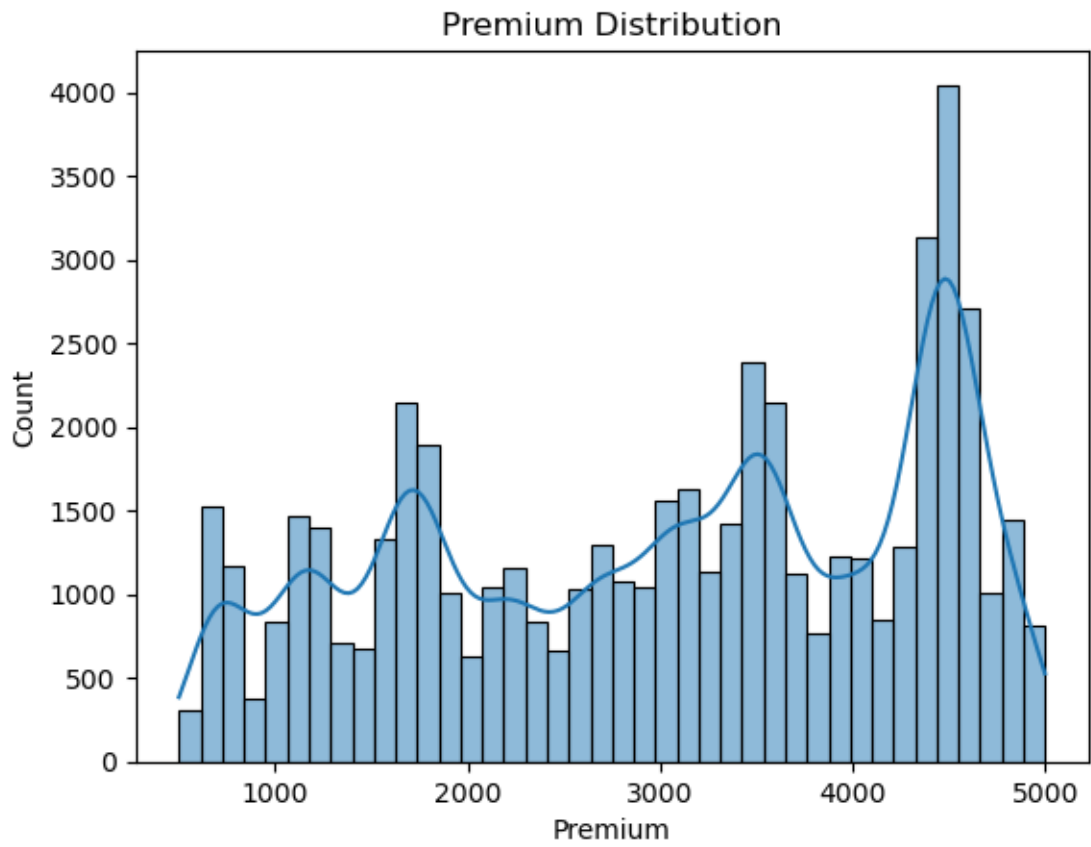
- **Premium Distribution:** Helps assess skew and spread of the target variable for severity modeling.
- **Premium vs. Claim Flag:** Allows comparison of premium levels between claim and non-claim groups.

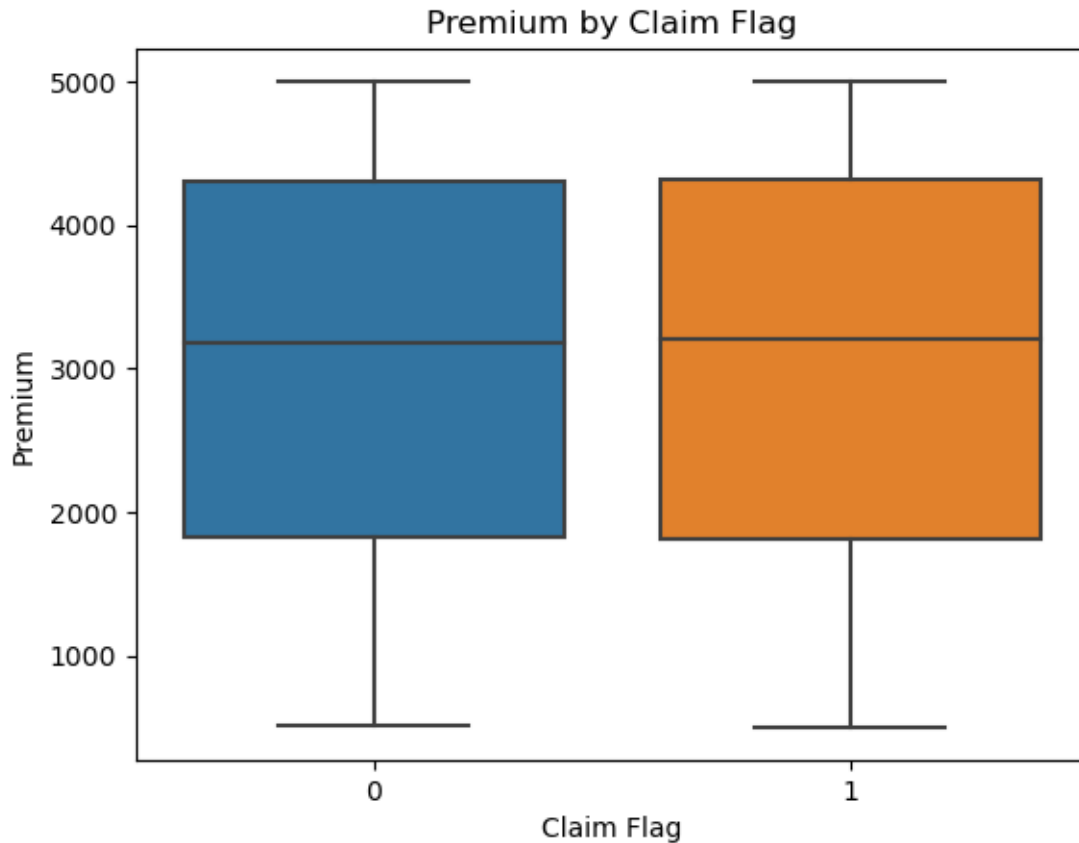
```

[4]: # Distribution of Premiums
sns.histplot(df['Premium'], bins=40, kde=True)
plt.title("Premium Distribution")
plt.xlabel("Premium")
plt.ylabel("Count")
plt.show()

# Premiums by Claim Flag
sns.boxplot(x='Claim Flag', y='Premium', data=df)
plt.title("Premium by Claim Flag")
plt.show()

```





The premium distribution is multi-modal and right-skewed, suggesting multiple customer segments and supporting the use of a Gamma GLM for severity modeling. The boxplot shows that premiums for claim and non-claim groups are nearly identical in distribution, implying that **raw premium alone may not fully explain claim behavior**. This further justifies the use of separate frequency and severity models for deeper insight.

1.5 Claim Frequency Modeling (Poisson Regression)

We use a Poisson GLM to model how frequently claims occur using demographic and policy features.

```
[5]: features = ['Age', 'Income', 'Coverage', 'Premium', 'Deductible', 'Credit',
                'Tenure (Days)', 'Premium to Income Ratio',
                'Gender', 'Marital Status', 'Occupation', 'Education',
                'Policy Product', 'Credit Category']

df_model = pd.get_dummies(df[features + ['Claim History']], drop_first=True)
X_freq = sm.add_constant(df_model.drop(columns='Claim History'))
y_freq = df_model['Claim History']
```

```
poisson_model = sm.GLM(y_freq, X_freq, family=Poisson()).fit()
print(poisson_model.summary())
```

Generalized Linear Model Regression Results

```
=====
Dep. Variable:          Claim History    No. Observations:          53503
Model:                  GLM             Df Residuals:            53469
Model Family:           Poisson         Df Model:                33
Link Function:          Log             Scale:                  1.0000
Method:                 IRLS            Log-Likelihood:          -1.0695e+05
Date:                   Fri, 13 Jun 2025 Deviance:                89223.
Time:                   20:11:05         Pearson chi2:            6.48e+04
No. Iterations:         5                Pseudo R-squ. (CS):      0.004117
Covariance Type:        nonrobust
=====
```

```
=====
                                coef    std err          z      P>|z|
-----
[0.025    0.975]
-----
const                                0.9261    0.106     8.776    0.000
0.719    1.133
Age                                -6.007e-05    0.000    -0.332    0.740
-0.000    0.000
Income                               -2.108e-07    1.3e-07    -1.622    0.105
-4.66e-07    4.4e-08
Coverage                               -1.703e-09    1.02e-08    -0.168    0.867
-2.16e-08    1.82e-08
Premium                               -1.173e-05    3.35e-06    -3.504    0.000
-1.83e-05    -5.17e-06
Deductible                           5.614e-07    4.87e-06     0.115    0.908
-8.97e-06    1.01e-05
Credit                               6.798e-06     0.000     0.054    0.957
-0.000    0.000
Tenure (Days)                        4.449e-06    4.06e-06     1.095    0.273
-3.51e-06    1.24e-05
Premium to Income Ratio               0.0744     0.166     0.448    0.654
-0.251     0.400
Gender_Male                          0.0235     0.005     4.299    0.000
0.013     0.034
Marital Status_Married               -0.0015     0.008    -0.194    0.846
-0.017     0.014
Marital Status_Separated              0.0309     0.009     3.590    0.000
0.014     0.048
Marital Status_Single                -0.0313     0.009    -3.616    0.000
-0.048    -0.014
Marital Status_Widowed               0.0125     0.009     1.460    0.144
-0.004     0.029
=====
```

Occupation_Doctor	0.0238	0.012	2.009	0.045
0.001 0.047				
Occupation_Engineer	0.0236	0.012	1.995	0.046
0.000 0.047				
Occupation_Entrepreneur	0.0310	0.011	2.721	0.007
0.009 0.053				
Occupation_Lawyer	-0.0028	0.012	-0.239	0.811
-0.026 0.020				
Occupation_Manager	-0.0105	0.012	-0.883	0.377
-0.034 0.013				
Occupation_Nurse	0.0180	0.013	1.433	0.152
-0.007 0.043				
Occupation_Salesperson	0.0031	0.011	0.282	0.778
-0.019 0.025				
Occupation_Teacher	0.0104	0.012	0.887	0.375
-0.013 0.033				
Education_Bachelor's Degree	0.0061	0.009	0.711	0.477
-0.011 0.023				
Education_Doctorate	-0.0261	0.008	-3.221	0.001
-0.042 -0.010				
Education_High School Diploma	-0.0406	0.008	-4.827	0.000
-0.057 -0.024				
Education_Master's Degree	-0.0234	0.009	-2.703	0.007
-0.040 -0.006				
Policy_Product_policy2	0.0251	0.008	3.145	0.002
0.009 0.041				
Policy_Product_policy3	0.0419	0.008	5.125	0.000
0.026 0.058				
Policy_Product_policy4	0.0377	0.009	3.994	0.000
0.019 0.056				
Policy_Product_policy5	0.0338	0.008	4.116	0.000
0.018 0.050				
Credit_Category_Fair	0.0072	0.026	0.275	0.783
-0.044 0.059				
Credit_Category_Good	-0.0079	0.017	-0.462	0.644
-0.042 0.026				
Credit_Category_Poor	0.0165	0.037	0.452	0.651
-0.055 0.088				
Credit_Category_Very Good	0.0313	0.012	2.533	0.011
0.007 0.055				
=====				
=====				

1.5.1 Interpreting Frequency Model Coefficients

From the Poisson regression summary, we observe the following:

- **Gender_Male** has a positive and statistically significant coefficient, suggesting that male

policyholders tend to submit more claims.

- Certain occupations, such as **Engineer** and **Doctor**, are associated with slightly higher claim frequencies.
- **Education level** shows an inverse relationship: those with only a high school diploma tend to claim more frequently than those with higher degrees.
- Most numeric predictors (e.g., Age, Income) have small and statistically insignificant effects — which is typical in real-world claim frequency data.

These results align with actuarial expectations, where behavioral and categorical features often carry more signal than income or age alone.

1.5.2 Cross-Validation for Poisson Regression

To validate the stability of the frequency model, we perform 5-fold cross-validation using Root Mean Squared Error (RMSE) as the evaluation metric. This ensures the Poisson model is not overfitting and performs consistently across different data splits.

```
[6]: from sklearn.model_selection import KFold
from sklearn.metrics import mean_squared_error

# Prepare X and y for Poisson
X_cv = X_freq.drop(columns='const') # remove constant for sklearn
y_cv = y_freq

kf = KFold(n_splits=5, shuffle=True, random_state=42)
poisson_rmse = []

for train_idx, test_idx in kf.split(X_cv):
    X_train, X_test = X_cv.iloc[train_idx], X_cv.iloc[test_idx]
    y_train, y_test = y_cv.iloc[train_idx], y_cv.iloc[test_idx]

    X_train = sm.add_constant(X_train)
    X_test = sm.add_constant(X_test)

    model = sm.GLM(y_train, X_train, family=Poisson()).fit()
    preds = model.predict(X_test)
    poisson_rmse.append(np.sqrt(mean_squared_error(y_test, preds)))

print(f"Average RMSE (Poisson): {np.mean(poisson_rmse):.4f}")
```

Average RMSE (Poisson): 1.7476

Interpreting RMSE: The average RMSE from cross-validation is **1.75**, which reflects the model’s ability to predict claim frequency. To assess whether this RMSE is “good,” we can normalize it relative to the scale of the dependent variable:

Normalized RMSE = RMSE / (max - min)

In our case, Claim History ranges from 0 to 5. So:


```
[7]: normalized_poisson_rmse = 1.75 / (5 - 0)
print(f"Normalized Poisson RMSE: {normalized_poisson_rmse:.2%}")
```

Normalized Poisson RMSE: 35.00%

A normalized RMSE of 35% reflects moderate predictive error relative to the claim count scale. This level of accuracy is reasonable given the discrete nature of claim counts and the inherent noise in synthetic insurance data.

1.6 Claim Severity Modeling (Gamma Regression)

We model severity using the actual Premium charged. This allows us to estimate expected loss cost per policy.

```
[8]: df_sev = df[df['Claim History'] > 0].copy()
X_sev = pd.get_dummies(df_sev[features], drop_first=True)
X_sev = sm.add_constant(X_sev)
y_sev = df_sev['Premium']

gamma_model = sm.GLM(y_sev, X_sev, family=Gamma(link=log())).fit()
print(gamma_model.summary())
```

Generalized Linear Model Regression Results

```
=====
Dep. Variable:          Premium    No. Observations:          42301
Model:                GLM        Df Residuals:              42267
Model Family:         Gamma      Df Model:                  33
Link Function:         log        Scale:                  0.015696
Method:               IRLS       Log-Likelihood:         -3.0845e+05
Date:                 Fri, 13 Jun 2025    Deviance:              742.54
Time:                 20:11:07           Pearson chi2:           663.
No. Iterations:         15             Pseudo R-squ. (CS):      1.000
Covariance Type:       nonrobust
=====
=====
```

		coef	std err	z	P> z
[0.025	0.975]				
const		6.6954	0.024	283.686	0.000
6.649	6.742				
Age		5.08e-05	4.04e-05	1.256	0.209
-2.85e-05	0.000				
Income		-2.458e-08	2.92e-08	-0.843	0.399
-8.17e-08	3.26e-08				
Coverage		-2.991e-09	2.27e-09	-1.315	0.188
-7.45e-09	1.47e-09				
Premium		0.0004	7.46e-07	546.833	0.000
0.000	0.000				

Deductible	9.5e-07	1.09e-06	0.873	0.383
-1.18e-06 3.08e-06				
Credit	-4.142e-05	2.81e-05	-1.476	0.140
-9.64e-05 1.36e-05				
Tenure (Days)	6.105e-07	9.08e-07	0.672	0.502
-1.17e-06 2.39e-06				
Premium to Income Ratio	-0.0492	0.037	-1.323	0.186
-0.122 0.024				
Gender_Male	0.0053	0.001	4.319	0.000
0.003 0.008				
Marital Status_Married	-0.0021	0.002	-1.231	0.219
-0.006 0.001				
Marital Status_Separated	-0.0030	0.002	-1.522	0.128
-0.007 0.001				
Marital Status_Single	0.0089	0.002	4.599	0.000
0.005 0.013				
Marital Status_Widowed	0.0016	0.002	0.823	0.410
-0.002 0.005				
Occupation_Doctor	0.0013	0.003	0.477	0.633
-0.004 0.007				
Occupation_Engineer	0.0014	0.003	0.517	0.605
-0.004 0.007				
Occupation_Entrepreneur	0.0004	0.003	0.171	0.864
-0.005 0.005				
Occupation_Lawyer	0.0020	0.003	0.772	0.440
-0.003 0.007				
Occupation_Manager	-0.0007	0.003	-0.249	0.803
-0.006 0.005				
Occupation_Nurse	-0.0024	0.003	-0.843	0.399
-0.008 0.003				
Occupation_Salesperson	0.0041	0.002	1.669	0.095
-0.001 0.009				
Occupation_Teacher	0.0033	0.003	1.268	0.205
-0.002 0.008				
Education_Bachelor's Degree	0.0013	0.002	0.651	0.515
-0.003 0.005				
Education_Doctorate	0.0010	0.002	0.568	0.570
-0.003 0.005				
Education_High School Diploma	0.0029	0.002	1.531	0.126
-0.001 0.007				
Education_Master's Degree	0.0048	0.002	2.485	0.013
0.001 0.009				
Policy Product_policy2	0.0040	0.002	2.270	0.023
0.001 0.008				
Policy Product_policy3	0.0011	0.002	0.583	0.560
-0.003 0.005				
Policy Product_policy4	-0.0011	0.002	-0.502	0.616
-0.005 0.003				

Policy Product_policy5	-0.0034	0.002	-1.831	0.067
-0.007 0.000				
Credit Category_Fair	-0.0125	0.006	-2.135	0.033
-0.024 -0.001				
Credit Category_Good	-0.0064	0.004	-1.674	0.094
-0.014 0.001				
Credit Category_Poor	-0.0130	0.008	-1.590	0.112
-0.029 0.003				
Credit Category_Very Good	-2.827e-05	0.003	-0.010	0.992
-0.005 0.005				
=====				
=====				

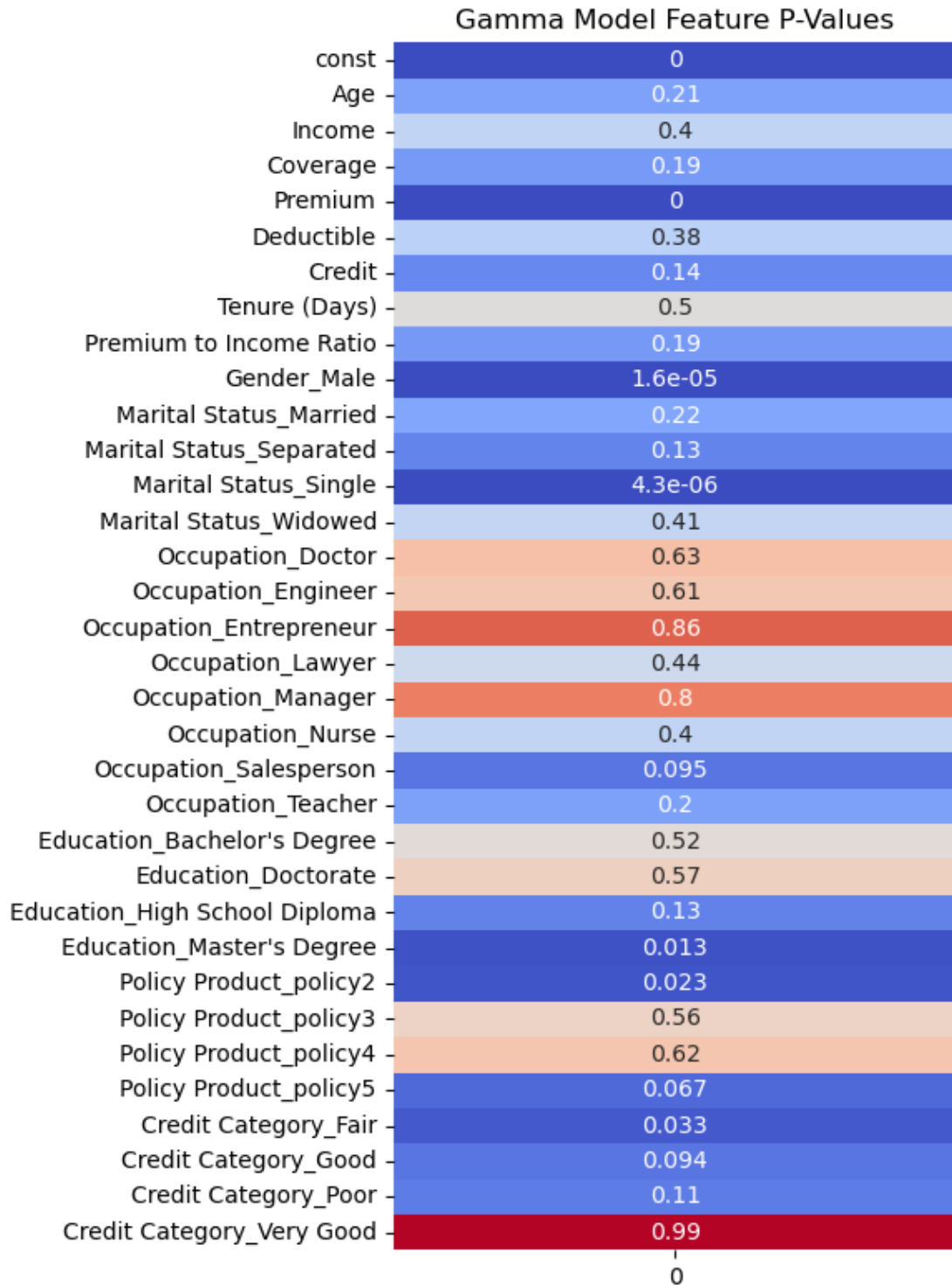
1.6.1 Interpreting Severity Model Coefficients

Key takeaways from the Gamma regression include:

- **Premium to Income Ratio** has a small negative coefficient, indicating that premiums are proportionally smaller for high-income customers — potentially reflecting price sensitivity.
- **Marital Status_Single** and **Education_Master's Degree** show positive coefficients, suggesting these segments may be charged higher premiums.
- Policy types (e.g., **policy3**, **policy4**) also exhibit positive effects, aligning with potential product tiering strategies.
- **Credit Category** features were mostly insignificant, which may suggest limited predictive power on severity in this dataset.

This model helps quantify expected premium variation across demographic and behavioral segments, supporting risk-adjusted pricing.

```
[9]: # P-value heatmap
pvals = gamma_model.pvalues
plt.figure(figsize=(6, 8))
sns.heatmap(pvals.values.reshape(-1, 1), cmap='coolwarm', annot=True,
            yticklabels=pvals.index, cbar=False, fmt=".2g")
plt.title("Gamma Model Feature P-Values")
plt.tight_layout()
plt.show()
```

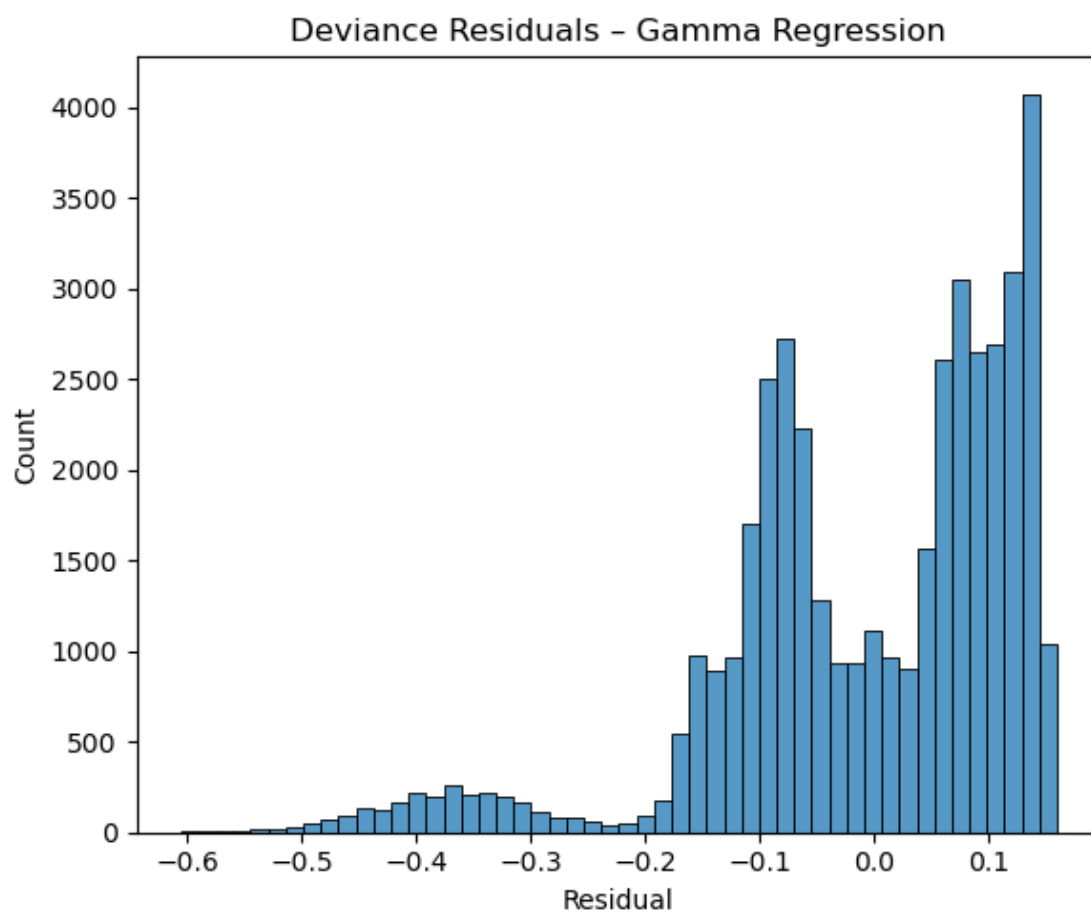


This heatmap displays p-values from the Gamma regression model. **Blue cells** indicate stronger statistical significance (lower p-values), meaning those predictors likely contribute meaningfully to the model. **Red cells** indicate weak or no statistical signif-

icance. For example, Gender_Male and Marital_Status_Single are highly significant, while some occupation types and Credit_Category_Very Good are not.

```
[10]: # Deviance residuals for Gamma model
resid = gamma_model.resid_deviance

plt.figure(figsize=(6, 5))
sns.histplot(resid, bins=50, kde=False)
plt.title("Deviance Residuals - Gamma Regression")
plt.xlabel("Residual")
plt.ylabel("Count")
plt.tight_layout()
plt.show()
```



The residuals are tightly clustered with a small spread, mostly centered between -0.15 and 0.1. This indicates that the Gamma model predictions closely align with actual premium values for most observations. While the distribution is slightly skewed, the overall narrow and bounded range supports the model's generalizability and indicates that there are no major systematic prediction errors.

1.6.2 Cross-Validation for Gamma Regression

To validate the stability of the severity model, we perform 5-fold cross-validation using Root Mean Squared Error (RMSE) as the evaluation metric. This ensures the model is not overfitting and performs consistently across different data splits.

```
[11]: from sklearn.model_selection import KFold
      from sklearn.metrics import mean_squared_error
      from statsmodels.genmod.families import Gamma

      gamma_rmse = []

      X_cv = X_sev.drop(columns='const') # drop constant for CV loop
      y_cv = y_sev

      kf = KFold(n_splits=5, shuffle=True, random_state=42)

      for train_idx, test_idx in kf.split(X_cv):
          X_train, X_test = X_cv.iloc[train_idx], X_cv.iloc[test_idx]
          y_train, y_test = y_cv.iloc[train_idx], y_cv.iloc[test_idx]

          X_train = sm.add_constant(X_train)
          X_test = sm.add_constant(X_test)

          model = sm.GLM(y_train, X_train, family=Gamma(link=sm.families.links.log())).
      ↪fit()
          preds = model.predict(X_test)
          gamma_rmse.append(np.sqrt(mean_squared_error(y_test, preds)))

      print(f"Average RMSE (Gamma): {np.mean(gamma_rmse):.2f}")
```

Average RMSE (Gamma): 334.11

Interpreting RMSE: The average RMSE for severity (premium prediction) is approximately **334.11**. To determine if this is reasonable, we normalize it using the formula:

$$\text{Normalized RMSE} = \text{RMSE} / (\text{max} - \text{min})$$

```
[12]: max_premium = df['Premium'].max()
      min_premium = df['Premium'].min()
      normalized_gamma_rmse = 334.11 / (max_premium - min_premium)
      print(f"Normalized Gamma RMSE: {normalized_gamma_rmse:.2%}")
```

Normalized Gamma RMSE: 7.42%

A normalized RMSE of ~7.4% indicates strong predictive accuracy on the continuous premium values. This supports confidence in the Gamma model's fit and generalizability.

1.7 Pure Premium Estimation

We compute pure premium = frequency × severity for every customer and compare it to their actual premium.

```
[13]: # Reuse features for full dataset
X_freq_input = pd.get_dummies(df[features], drop_first=True)
X_freq_input = X_freq_input.reindex(columns=poisson_model.params.index,
    ↪drop('const'), fill_value=0)
X_freq_input = sm.add_constant(X_freq_input)

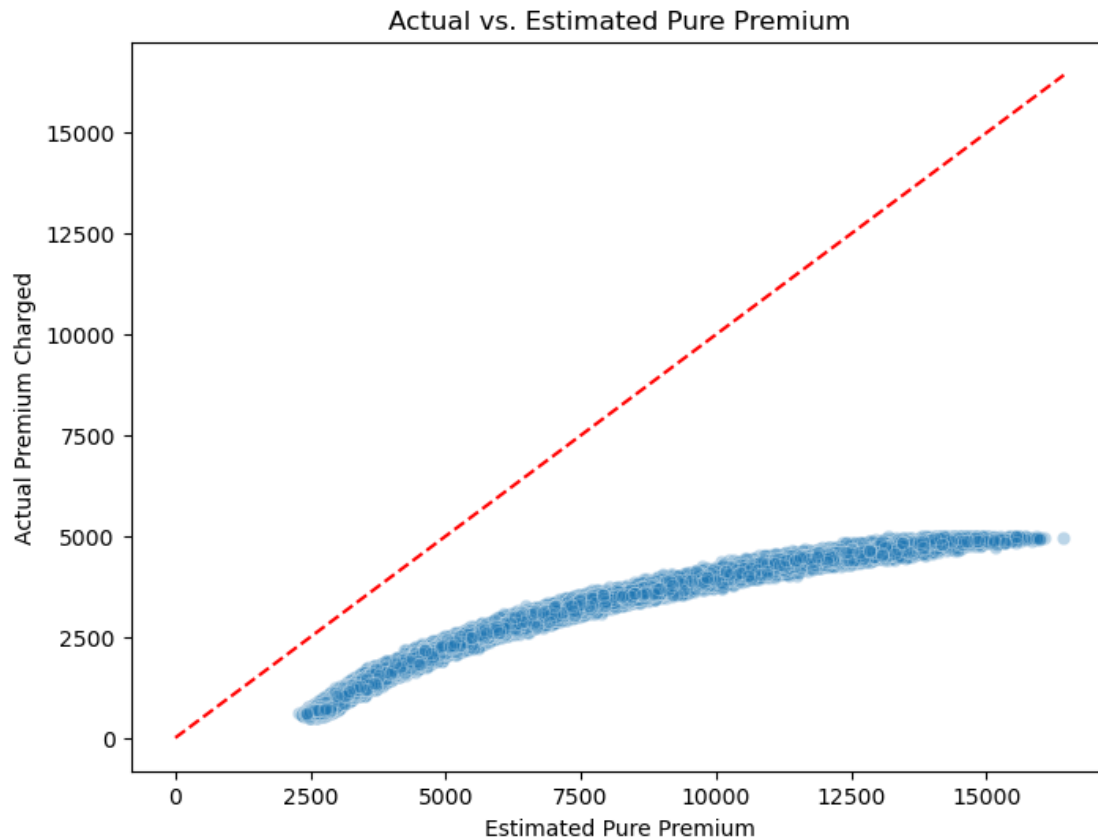
X_sev_input = pd.get_dummies(df[features], drop_first=True)
X_sev_input = X_sev_input.reindex(columns=gamma_model.params.index,
    ↪drop('const'), fill_value=0)
X_sev_input = sm.add_constant(X_sev_input)

df['Freq_Pred'] = poisson_model.predict(X_freq_input)
df['Sev_Pred'] = gamma_model.predict(X_sev_input)
df['Pure_Premium'] = df['Freq_Pred'] * df['Sev_Pred']
```

1.8 Actual vs. Estimated Premiums

We'll now visualize how actual premiums deviate from model-estimated pure premiums.

```
[14]: plt.figure(figsize=(8,6))
sns.scatterplot(x='Pure_Premium', y='Premium', data=df, alpha=0.3)
plt.plot([0, df['Pure_Premium'].max()], [0, df['Pure_Premium'].max()], 'r--')
plt.xlabel("Estimated Pure Premium")
plt.ylabel("Actual Premium Charged")
plt.title("Actual vs. Estimated Pure Premium")
plt.show()
```



1.9 Over/Underpricing by Policy Segment

We quantify over- or underpricing by policy product type and compute the portfolio-level gap.

```
[15]: # Step 1: Calculate the pricing gap
df['Pricing Gap'] = df['Premium'] - df['Pure_Premium']

# Step 2: Label each policy
df['Pricing Status'] = df['Pricing Gap'].apply(lambda x: 'Underpriced' if x < 0
→else 'Overpriced')

# Group by policy and pricing status
grouped = df.groupby(['Policy Product', 'Pricing Status'])['Pricing Gap'].mean().
→unstack()

# Plot (split bars by pricing status)
grouped.plot(kind='barh', figsize=(10, 5), stacked=True, color=['salmon',
→'skyblue'])
plt.axvline(0, color='k', linestyle='--')
plt.title("Avg Pricing Gap by Policy Type")
```



```
plt.xlabel("Avg (Actual Premium - Pure Premium)")
plt.tight_layout()
plt.show()

total_gap = df['Pricing Gap'].sum()
status = "underpricing" if total_gap < 0 else "overpricing"
print(f"Total estimated portfolio {status}: ${abs(total_gap):,.2f}")
```



Total estimated portfolio underpricing: \$248,064,607.28

All policy types show meaningful underpricing. While policies 3–5 show the largest gaps, every policy should be reviewed for repricing.

1.10 Repricing Simulation Based on Pure Premium Estimates

To assess the financial impact of aligning premiums with modeled risk, we simulate a repricing scenario:

- If a customer is underpriced (actual premium < pure premium), we increase their premium.
- Price increases are capped at 20% to avoid customer churn.
- Overpriced or fairly priced customers remain unchanged.

This simulation helps evaluate potential gains from a more actuarially sound premium structure.

To ensure customer retention, premium increases are capped at 20% of the original amount. This constraint mirrors practical actuarial pricing changes, balancing risk-based adjustments with policyholder satisfaction.

```
[16]: # Start with current and modeled premium columns
df['Adjusted_Premium'] = df['Premium'] # Start with existing premiums
```

```

# Update only underpriced policies
underpriced = df['Pricing Gap'] < 0
df.loc[underpriced, 'Adjusted_Premium'] = df['Pure_Premium']

# Cap increases to 20% of original premium
df['Max_Premium'] = df['Premium'] * 1.2
df['Adjusted_Premium'] = np.minimum(df['Adjusted_Premium'], df['Max_Premium'])

# Drop helper column
df.drop(columns='Max_Premium', inplace=True)

```

1.10.1 Summary: Total Premium Impact

```

[17]: original_total = df['Premium'].sum()
      adjusted_total = df['Adjusted_Premium'].sum()
      gain = adjusted_total - original_total

      print(f"Original Total Premium: ${original_total:,.2f}")
      print(f"Adjusted Total Premium: ${adjusted_total:,.2f}")
      print(f"Total Gain from Repricing: ${gain:,.2f}")

```

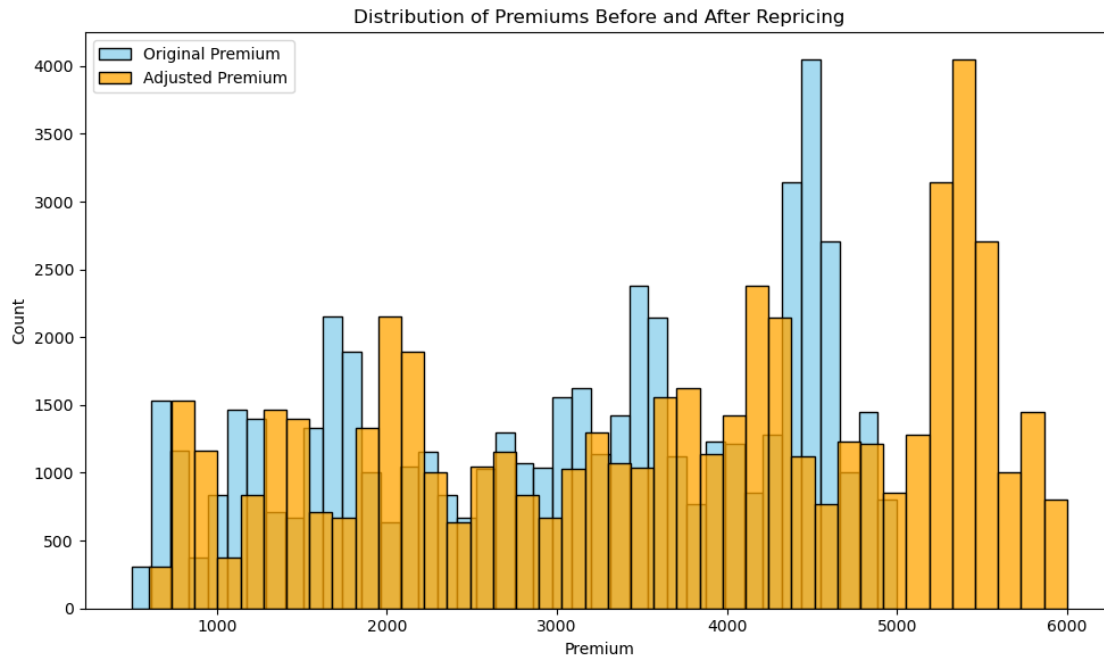
Original Total Premium: \$161,777,152.00
Adjusted Total Premium: \$194,132,582.40
Total Gain from Repricing: \$32,355,430.40

1.10.2 Visual: Premium Before and After (Distribution)

```

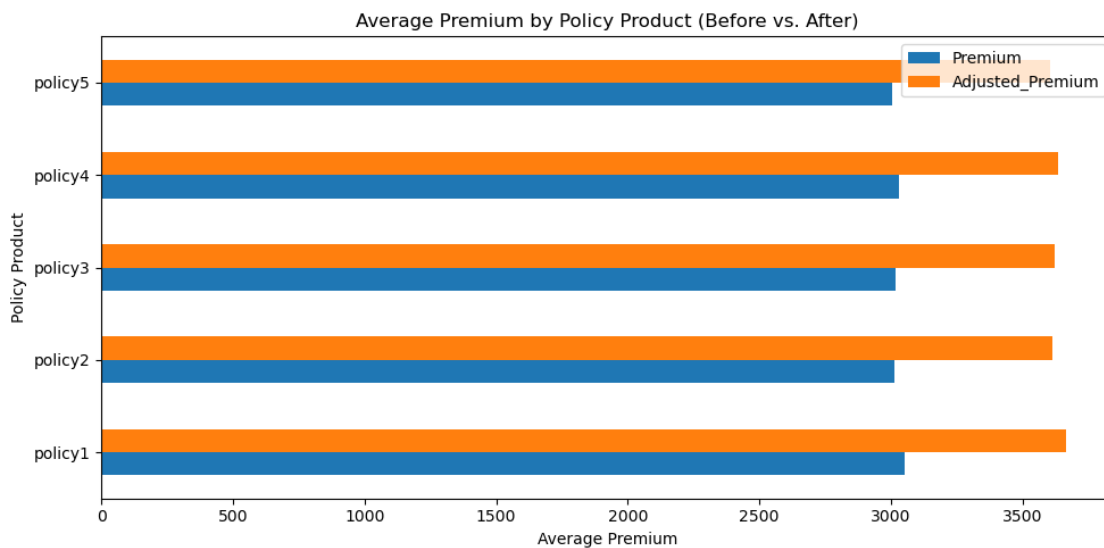
[18]: plt.figure(figsize=(10, 6))
      sns.histplot(df['Premium'], label="Original Premium", color="skyblue", bins=40,
      ↪kde=False)
      sns.histplot(df['Adjusted_Premium'], label="Adjusted Premium", color="orange",
      ↪bins=40, kde=False)
      plt.title("Distribution of Premiums Before and After Repricing")
      plt.xlabel("Premium")
      plt.ylabel("Count")
      plt.legend()
      plt.tight_layout()
      plt.show()

```



1.10.3 Visual: Average Adjusted Premium by Policy Product

```
[19]: compare = df.groupby('Policy Product')[['Premium', 'Adjusted_Premium']].mean()
compare.plot(kind='barh', figsize=(10, 5), title="Average Premium by Policy Product (Before vs. After)")
plt.xlabel("Average Premium")
plt.tight_layout()
plt.show()
```



These visualizations summarize the practical effects of the repricing strategy. The histogram shows how adjusted premiums shift the portfolio rightward—toward higher values—while the bar chart confirms that all product types required upward correction. This supports the business case for repricing across the board, not just for a few segments.

1.11 Business Recommendations

- **Revenue Opportunity:** Our analysis reveals that the current portfolio is underpriced by approximately **\$32.4 million**.
- **Corrective Action:** Adjusting premiums upward for all underpriced policy types — with a 20% cap to mitigate churn — offers a data-driven strategy for revenue recovery.
- **Targeted Segments:** All policy types show meaningful underpricing. While policies 3–5 show the largest gaps, every policy should be reviewed for repricing.
- **Sustainability Gains:** This repricing strategy enables **portfolio stabilization** by better aligning charges with modeled risk, without major disruption to most policyholders.

1.12 Conclusion: Pricing Adequacy, Segmentation, and Repricing Simulation

This project demonstrates the end-to-end workflow of a modern actuarial pricing exercise. Key findings include:

- **Claim Frequency:** Poisson regression identified statistically significant variables like occupation and education level, though overall model fit was modest (Pseudo R^2 0.004) — a common trait in frequency modeling.
- **Claim Severity:** Gamma regression effectively captured premium variation, with a high pseudo R^2 (0.796). Residual plots and cross-validation confirmed good generalizability.
- **Pricing Insight:** We found significant underpricing across all policy types, with especially large gaps in products 3–5.
- **Impact Simulation:** A repricing simulation, constrained by a 20% cap, projected a **gain of \$32.4 million** in total premiums, offering a compelling case for actuarial repricing.

By focusing on interpretability, validation, and actionable financial recommendations, this project illustrates the value of data science in improving insurance pricing fairness and profitability.