

# Insurance Claim Status Prediction using ML Classification Models

Machine Learning 1, DS&BA, UW

---

# Agenda

**1** Project Objective

**2** Exploratory Data  
Analysis (EDA)

**3** Feature Selection

**4** Feature Engineering

**5** Models Consideration

**6** Prediction

---

---

# Project Objective

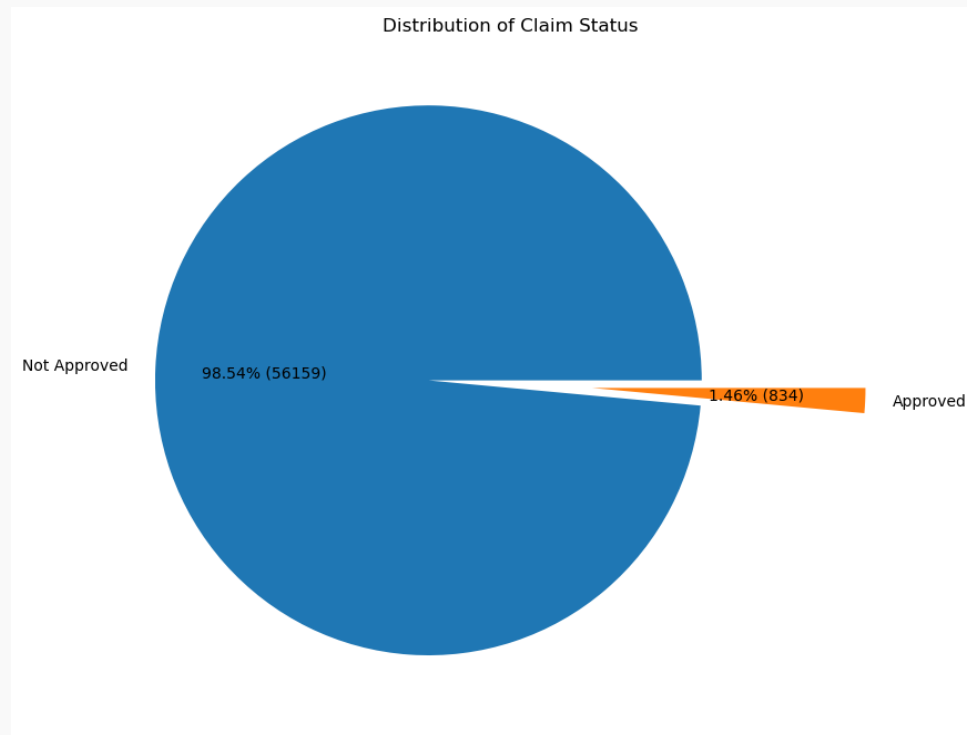
By training on the given sample, best algorithm should be found for predicting the travel insurance claim outcome (target variable `claim_status`: 1 – Approved / 0 – Not approved) on unseen data.

- Training Sample: 56993 observations including target variable
- Test Sample: 6333 observations without target variable
- 14 features in total including target variable

# Exploratory Data Analysis

- No missing values and duplicates
- 3 numeric variables (revenue, reward, customer\_score)
- 3 discrete variables (trip\_length, person\_age, support\_interactions)
- 7 nominal variables (person\_gender, entity\_type, channel, agent\_id, entity\_a, location, product\_id)
- 1 binary variable (claim\_status)

The data is heavily imbalanced.

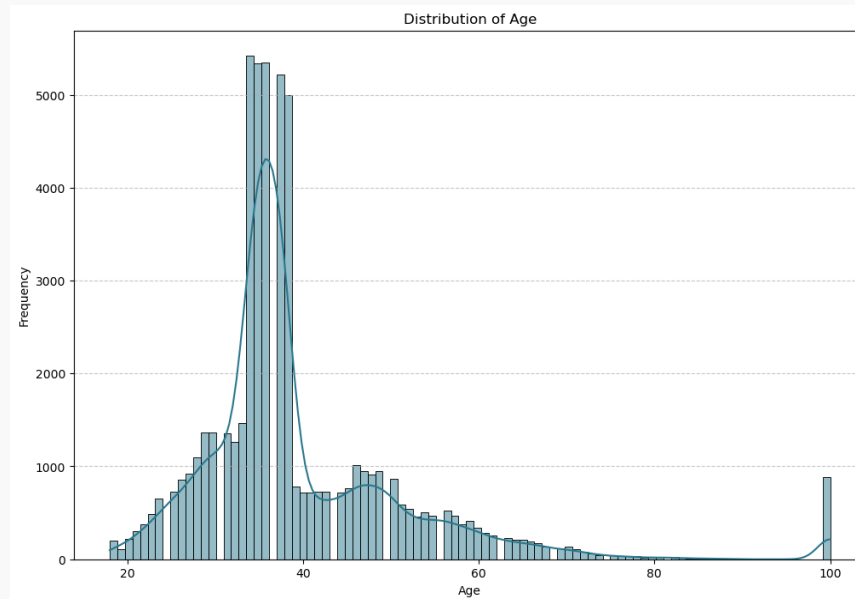
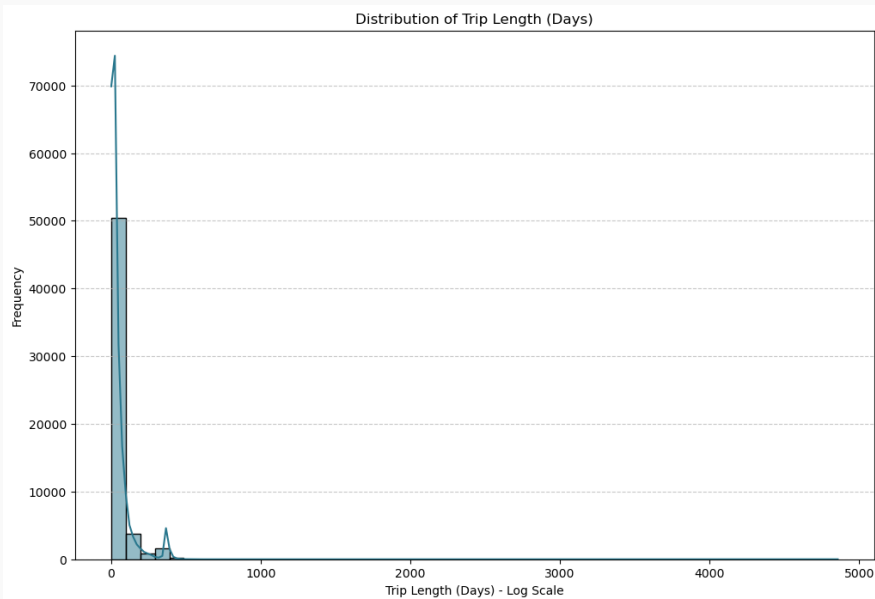


# Exploratory Data Analysis

For trip length and person's age:

- Strongly right-skewed, especially for trip length
- Extreme outliers (e.g. max 4856 and min 1 for trip length).

Further binning and transformation have to be considered for modeling.

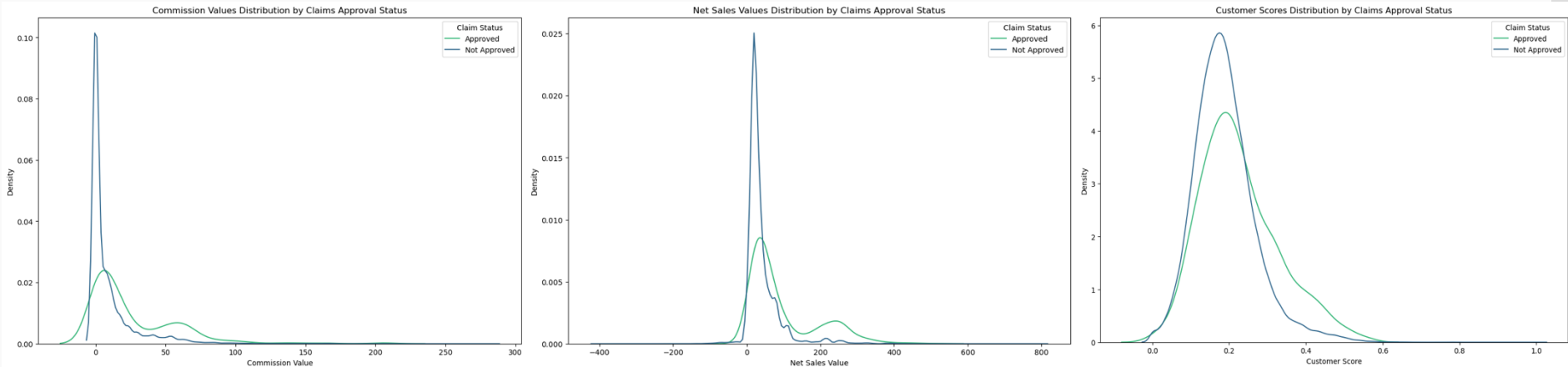


# Exploratory Data Analysis

For commission values, net sales values and customer scores:

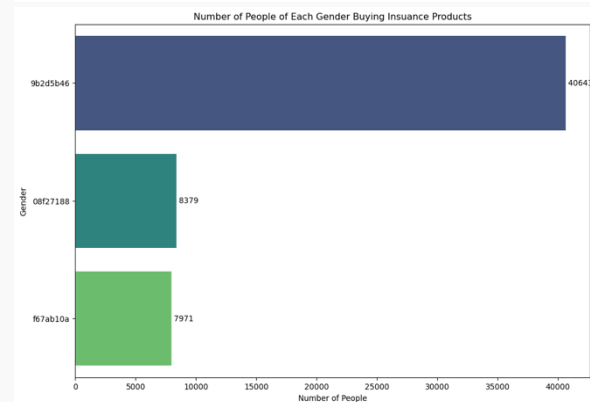
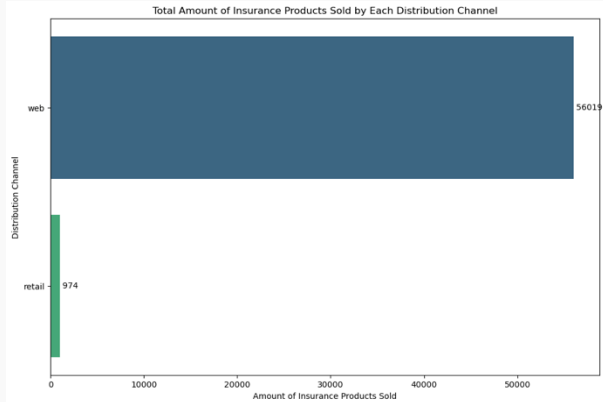
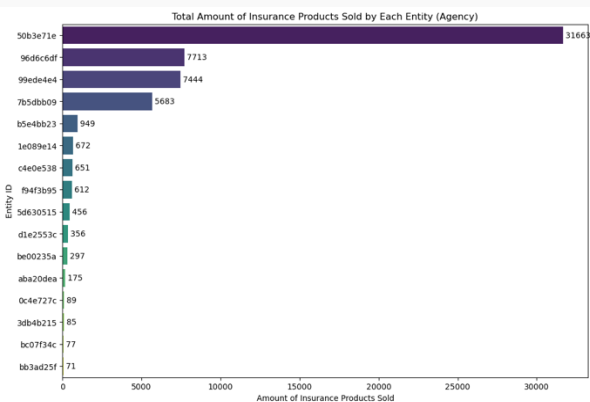
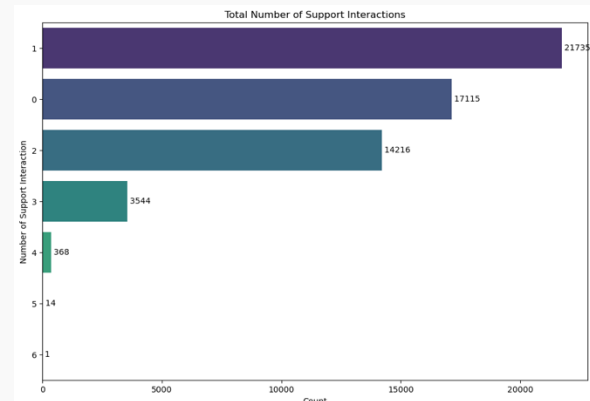
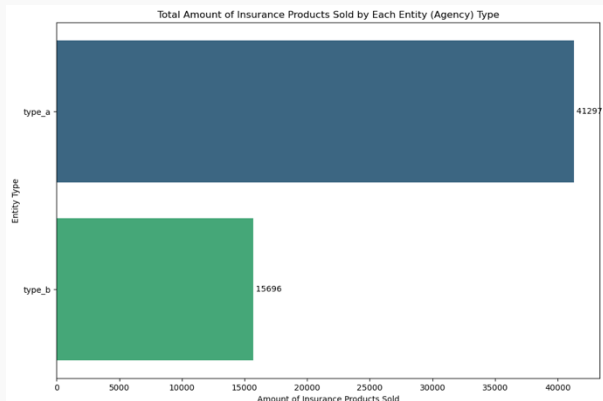
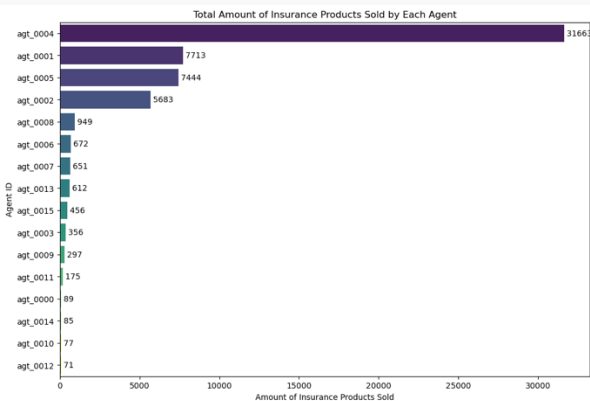
- Unevenly distributed referring to claim status

Further transformation have to be considered for modeling.



# Exploratory Data Analysis

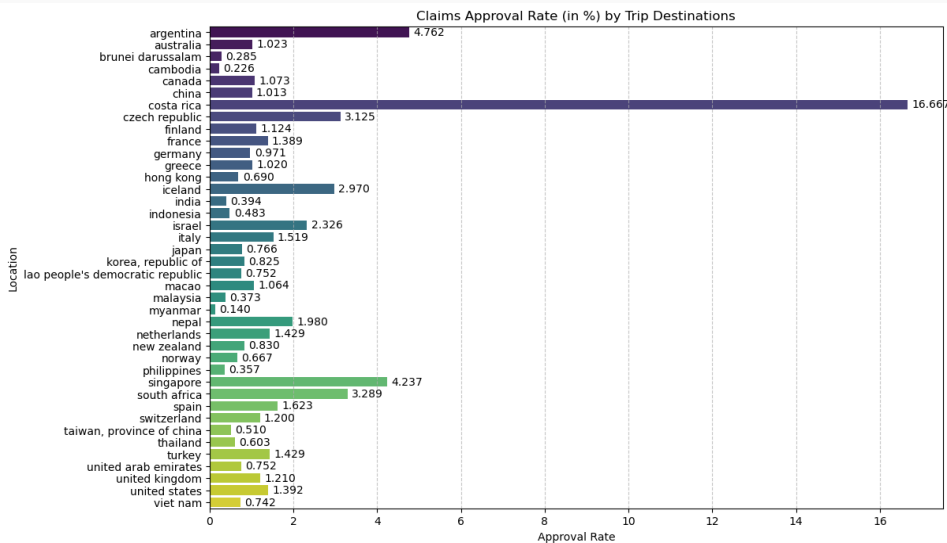
For agent ID, entity type, entity ID, support interactions, gender and channel, it seems there are duplicated, meaningless and unimpactful variables to target.



# Exploratory Data Analysis

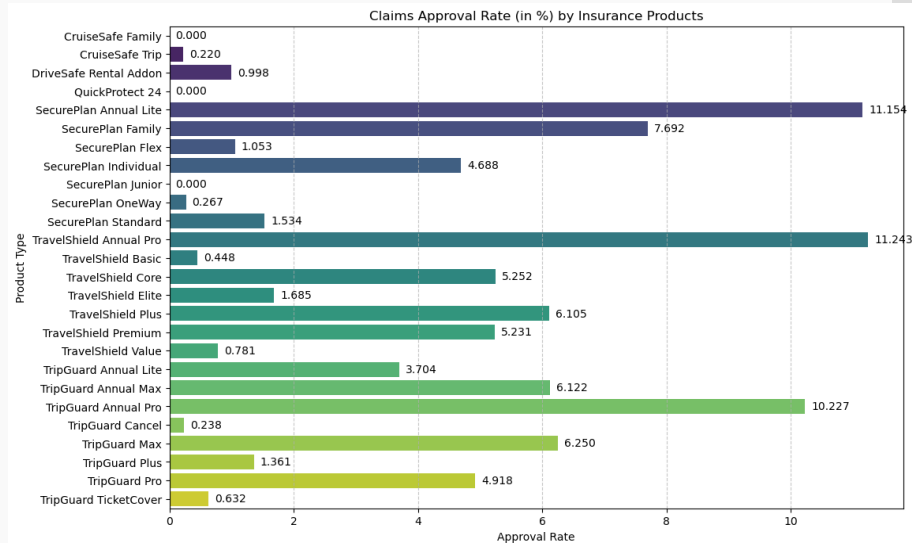
For location:

- Over 140+ countries in observations
- Only 20+ being approved for claims
- Rare distributions (e.g. 11942 observed for Singapore and more than half are below 100 observed)



For insurance product:

- Around 20+ types in observations
- Rare distributions (e.g. 16795 for most popular choices but some plans are only few observed)
- Approval ratio is random



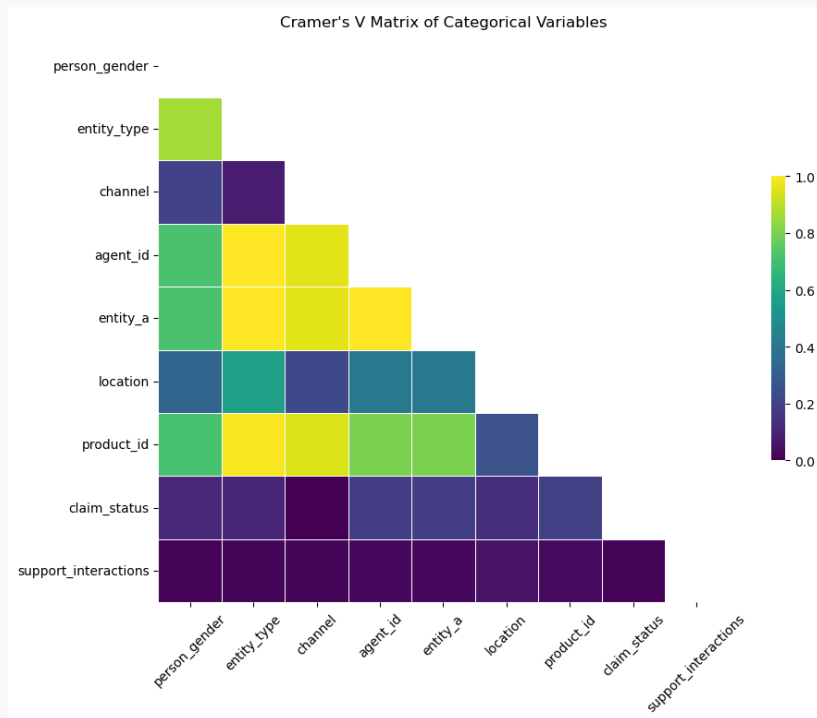


# Train / Validation / Test split on training sample (insurance\_train.csv)

Before feature selection & engineering:

- Train Set (60%), Validation Set (20%), and Internal Test Set (20%) split
- Avoid data leakage while feature selected and engineered based on train split only
- stratify = insurance['claim\_status'] used to ensure all splits maintain the same proportion of 0s and 1s
- Validation and interval testing helping better estimate of generalization and stable / reliable model on minority class

# Feature Selection



Cramer's V for categorical target and categorical inputs:

- support\_interactions -> nearly 0 importance to all
- entity\_a, entity\_type and agent\_id -> identical, perfectly one-to-one and consistent (coefficient = 1)
- channel -> highly associated to product\_id and entity (coefficient > 0.95) and 0 association to target

# Feature Selection

<b><u>Variable</u></b>	<b>p-value</b>
<u>revenue</u>	2.986475e-140
<u>reward</u>	4.529682e-71
<u>trip_length</u>	6.592646e-42
<u>customer_score</u>	4.212666e-23
person_age	3.990284e-04

ANOVA for categorical target  
and continuous inputs:

All features are statistically significant.

In short, following variables will be removed:

entity\_a  
entity\_type  
channel  
support\_interactions

# Feature Engineering

No zero values

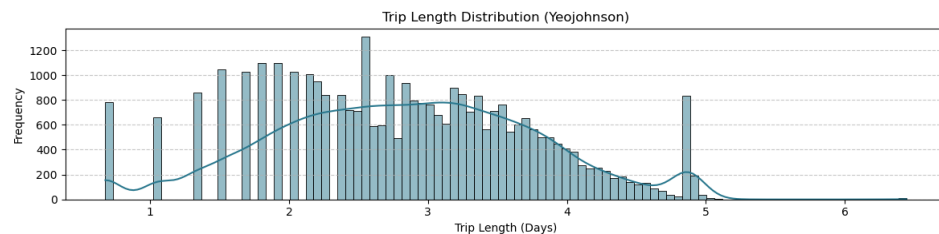
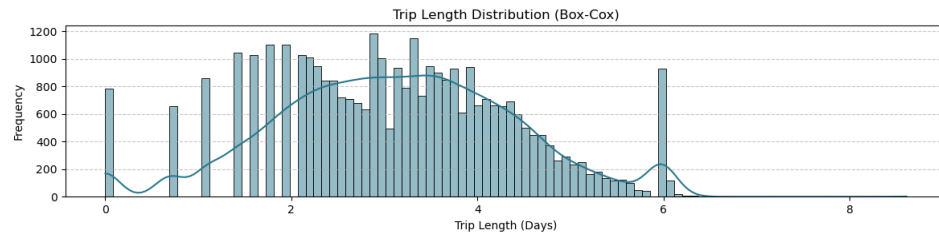
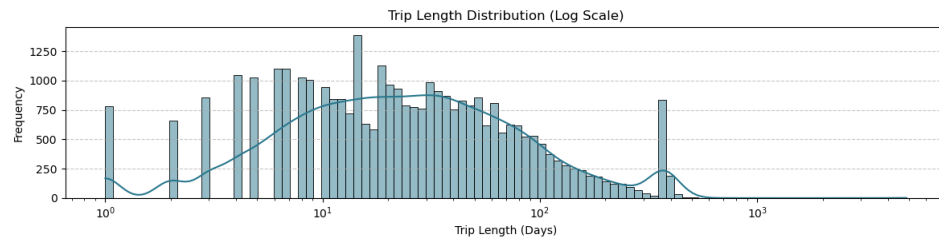
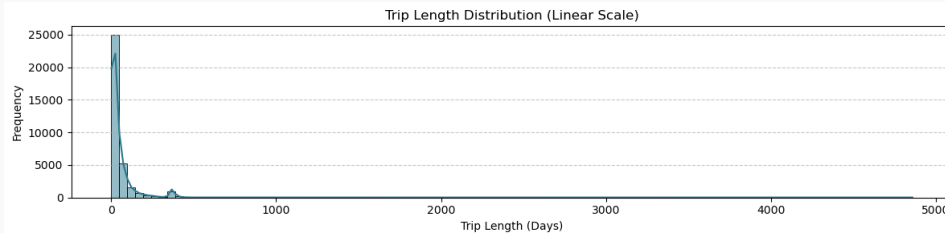
Original skewness of trip\_length: 23.84

Log Skewness: -0.01

Box-Cox Skewness ( $\lambda = 0.00$ ): -0.00

Yeo-Johnson Skewness ( $\lambda = -0.07$ ): 0.01

**Log transformation** is chosen.



# Feature Engineering

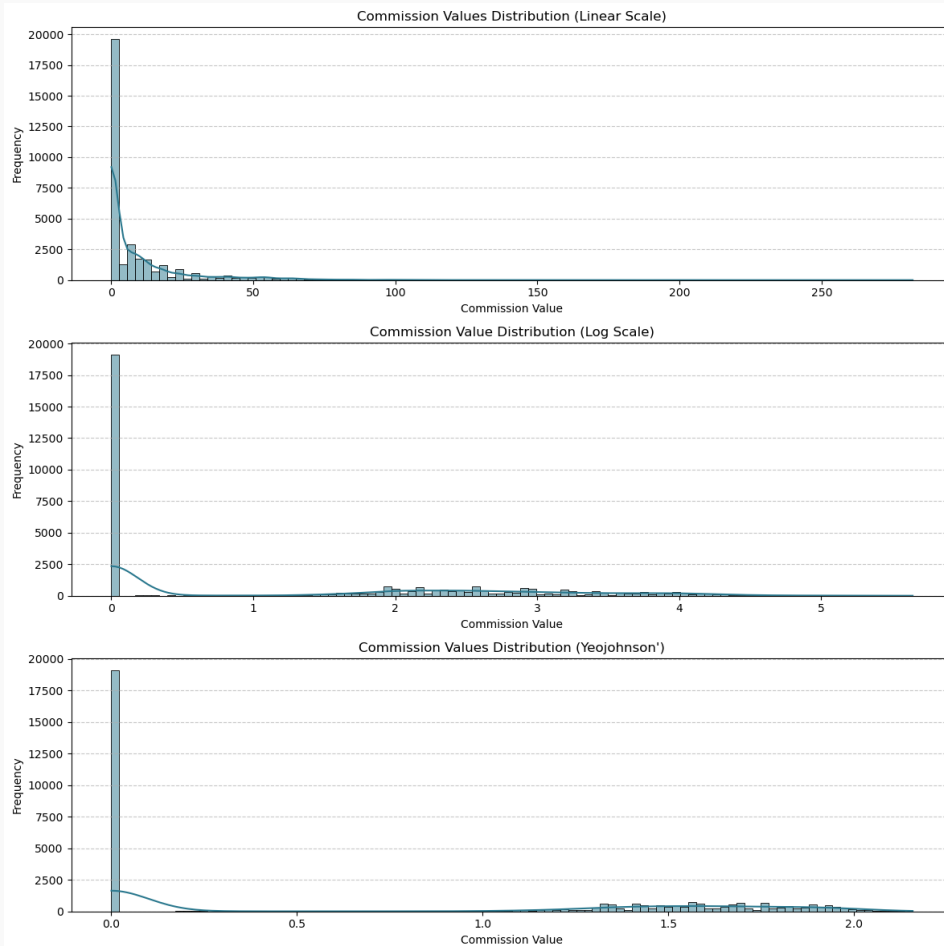
Including zero values

Original skewness of reward: 3.97

Log Skewness: 0.67

Yeo-Johnson Skewness ( $\lambda = -0.42$ ): 0.39

**Yeo-Johnson transformation** is chosen.



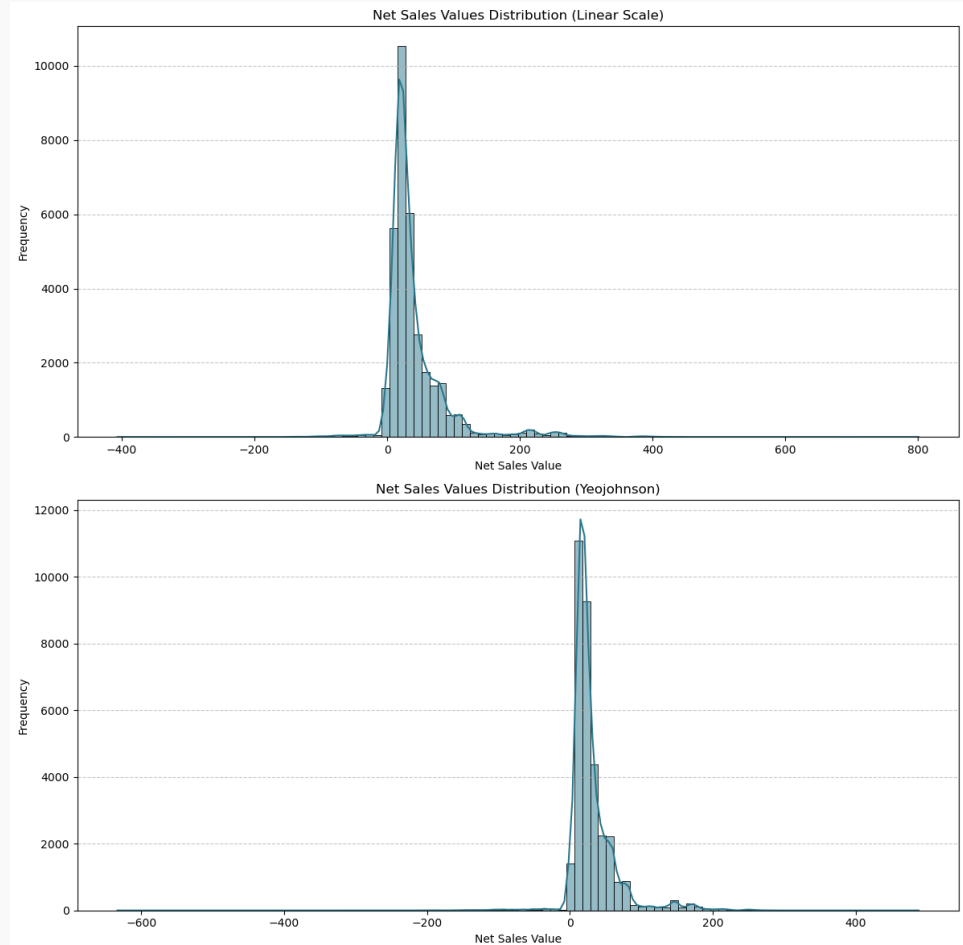
# Feature Engineering

Including zero and negative values

Original skewness of revenue: 3.21

Yeo-Johnson Skewness ( $\lambda = 0.91$ ): 1.33

**Yeo-Johnson transformation** is chosen.



# Feature Engineering

Based on EDA, **Jenk's natural break** for person\_age to 8 classes as ordinal (last group is merged with previous group due to small size)

<b>Age binning</b>	<b>Counts</b>	<b>Encoded as</b>
<u>17.999, 26.0</u>	2301	1
<u>26.0, 32.0</u>	4457	2
<u>32.0, 36.0</u>	10600	3
<u>36.0, 42.0</u>	7845	4
<u>42.0, 51.0</u>	4494	5
<u>51.0, 62.0</u>	2783	6
<u>62.0, 100.0</u>	1715	7

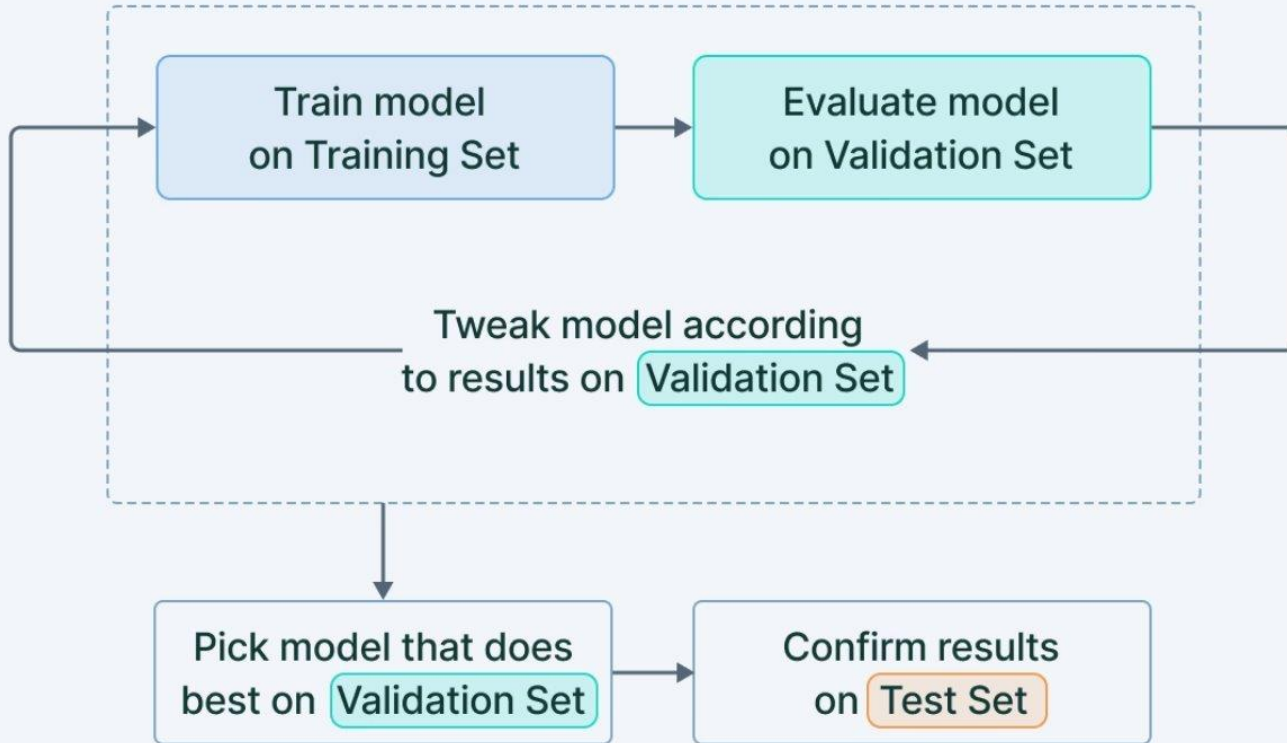
Based on EDA, **frequency binning** as ordinal for location to 4 groups of popularity

<b>Location binning</b>	<b>Counts</b>	<b>Encoded as</b>
<u>Most popular</u>	$\geq 1000$	4
<u>Moderate popular</u>	$\geq 100$	3
<u>Less popular</u>	$\geq 10$	2
<u>Least popular</u>	$< 10$	1

<b>One-hot Encoding as Dummy</b>	person_gender, agent_id, product_id (names with "Family" grouped as "Other" due to small size)
----------------------------------	--

# Models Training





# Model Consideration

4 continuous variable, 2 ordinal variables, and 41 dummy variables

## SVM

- RBF
- Polynomial
- Linear

### StandardScaler()

- perform better if 0 mean and unit variance

## KNN

### MinMaxScaler()

- Deal with the sensitivity of distance in the same range [0, 1]

## Logistic Reg. (Elastic Net)

- L1 Lasso (1)
- L2 Ridge (0)


### StandardScaler()

- help with magnitude

# Model Consideration

To handle imbalance dataset, four sampling techniques are considered and only fit to training split.

Each model will try with each sampling techniques to compare the differences.

SMOTE	SMOTETomek	Over Sampling	Under Sampling
<ul style="list-style-type: none"><li>• Create synthetic samples for class 1</li><li>• Keep class 0</li><li>• Better generalization</li></ul>	<ul style="list-style-type: none"><li>• Same functions as SMOTE</li><li>• Clean overlapping samples from both classes</li><li>• Reduce noise in the decision boundary area</li></ul>	<ul style="list-style-type: none"><li>• Undersample class 0 which will cause severe info. Loss</li><li>• Not good to SVM/KNN if requires sufficient data</li></ul>	<ul style="list-style-type: none"><li>• Risk of overfitting due to simple oversampling existing class 1</li></ul>
 <p>Not preferred but still add for comparison</p>			

# Model Consideration

1. Always apply the sampling method before scaling

**Sampler**



**Preprocessor**  
**Scaler + Remainder**

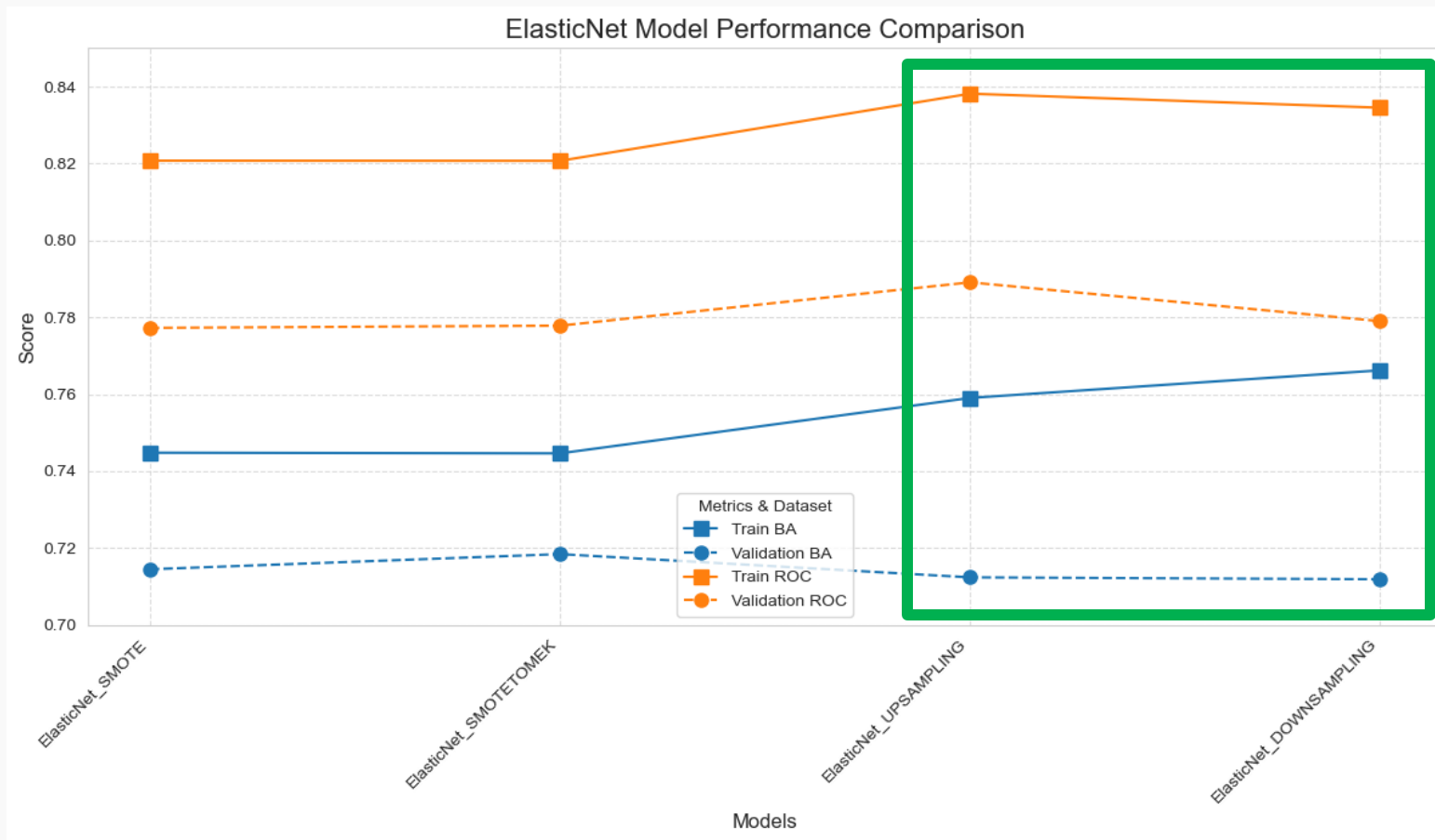
2. Add a preprocessor to ensure only the numeric features being scaled



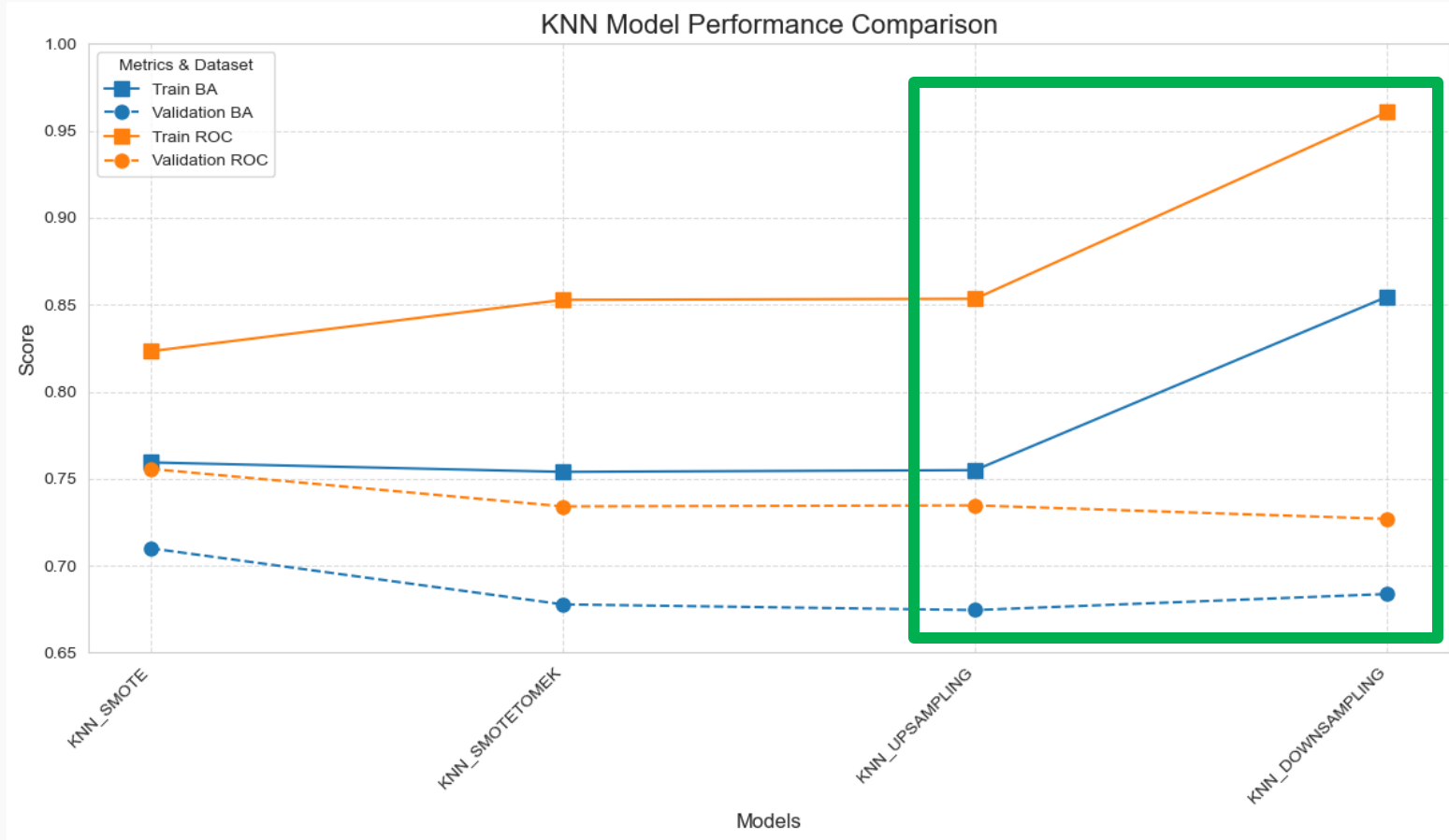
**Model**

3. Fit the models with modified grids
  - RandomizedSearchCV(): reduce computational time
  - StratifiedKFold(n = 5): ensure each fold has the same proportion of observations with target
  - Scoring = balanced accuracy (prioritized)

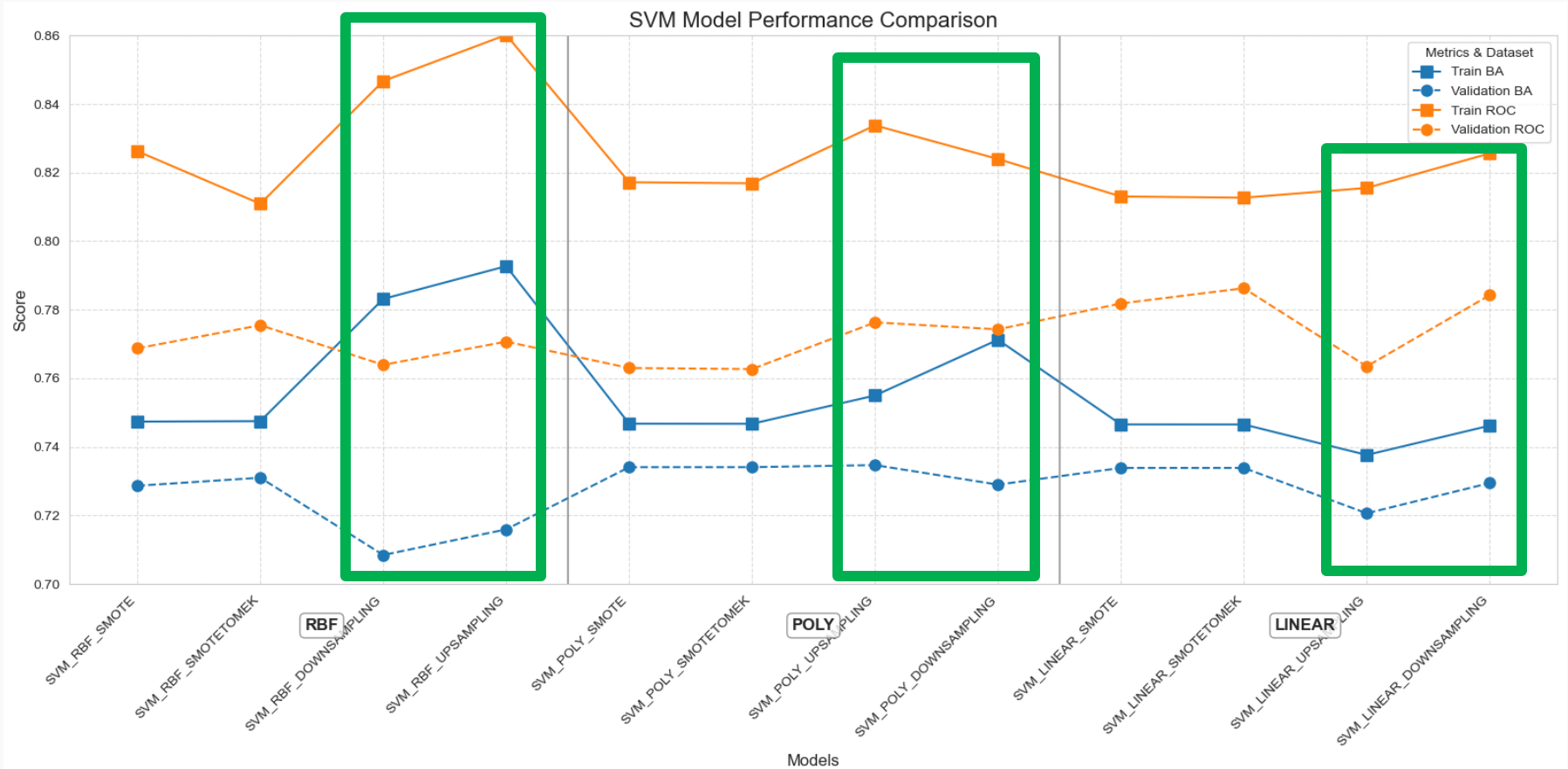
# Model Consideration



# Model Consideration



# Model Consideration



# Model Consideration

*Using SMOTETomek:*

<b>Logistic Regression (Elastic Net)</b>	<b>KNN(<i>SMOTE</i>)</b>	<b>SVM(Poly)</b>	<b>SVM(Linear)</b>	<b>SVM (RBF)</b>
Train: - AUC: 0.82	Train: - AUC: 0.82	Train: - AUC: 0.82	Train: - AUC: 0.81	Train: - AUC: 0.81
Validation: - BA: 0.72 - AUC: 0.78	Validation - BA: 0.71 - AUC: 0.76	Validation: - BA: 0.73 - AUC: 0.76	Validation: - BA: 0.73 - AUC: 0.79	Validation: - BA: 0.73 - AUC: 0.78
<b>Expected Value for each model (mean cross-validated BA scores)</b>				
0.74	0.75	0.75	0.75	0.75

# Model Consideration

By considering:

- Model overall stability
- Computational efficiency
- Scalability (data with heavy dummies and less numeric features)
- Interpretability
- Comparison of other metrics (Recall and F1-score for class 1)
- Risk of overfitting (prone to overfit minority class)

**Final expected value / BA : 0.75**

Selected model: **SVM (Linear)**

On Interval Test split:

- Train BA: 0.75
- Internal Test BA: 0.75
- Train AUC: 0.81
- Internal Test AUC: 0.80

Classification Report based on Test:					
	precision	recall	f1-score	support	
0	0.99	0.84	0.91	11232	
1	0.06	0.65	0.11	167	
accuracy			0.84	11399	
macro avg	0.53	0.75	0.51	11399	
weighted avg	0.98	0.84	0.90	11399	



---

# Prediction

*Based on same feature engineering methods and pipeline:*

<b><u>Predicted Claim Status</u></b>	<b>Counts</b>	<b>Ratio</b>
<u>Class 0</u>	5266	83.15%
<u>Class 1</u>	1067	16.85%

A decorative gray vertical bar is on the left side of the slide. A thin horizontal line extends from the top of this bar across the top of the slide. Another thin horizontal line is positioned below the text.

**THANK YOU**