**bip.**

HERE TO DARE

# Exemption Code Mapper

# Agenda

1 Introduction

2 Preprocessing

3 Model Deployment

4 User Interface

5 MLOps

6 Discussion

# 1 Introduction

# Introduction & Our goal

The objective is to create a model that predicts the **exemption VAT** code for invoices.

To achieve the objective, we were provided with a dataset concerning invoices and their associated characteristics.

# Null values

**1** Initially, there were numerous columns with over half of their values null.

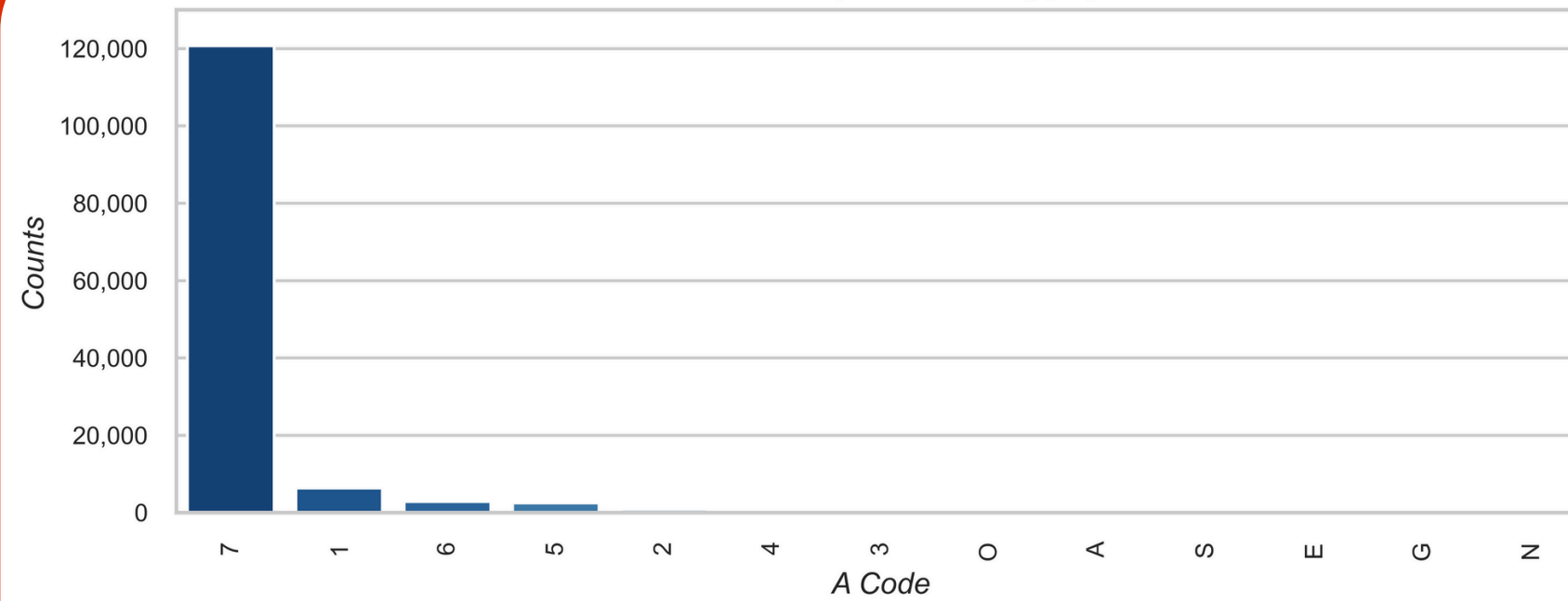**2** We dropped columns with more than 60% null values.

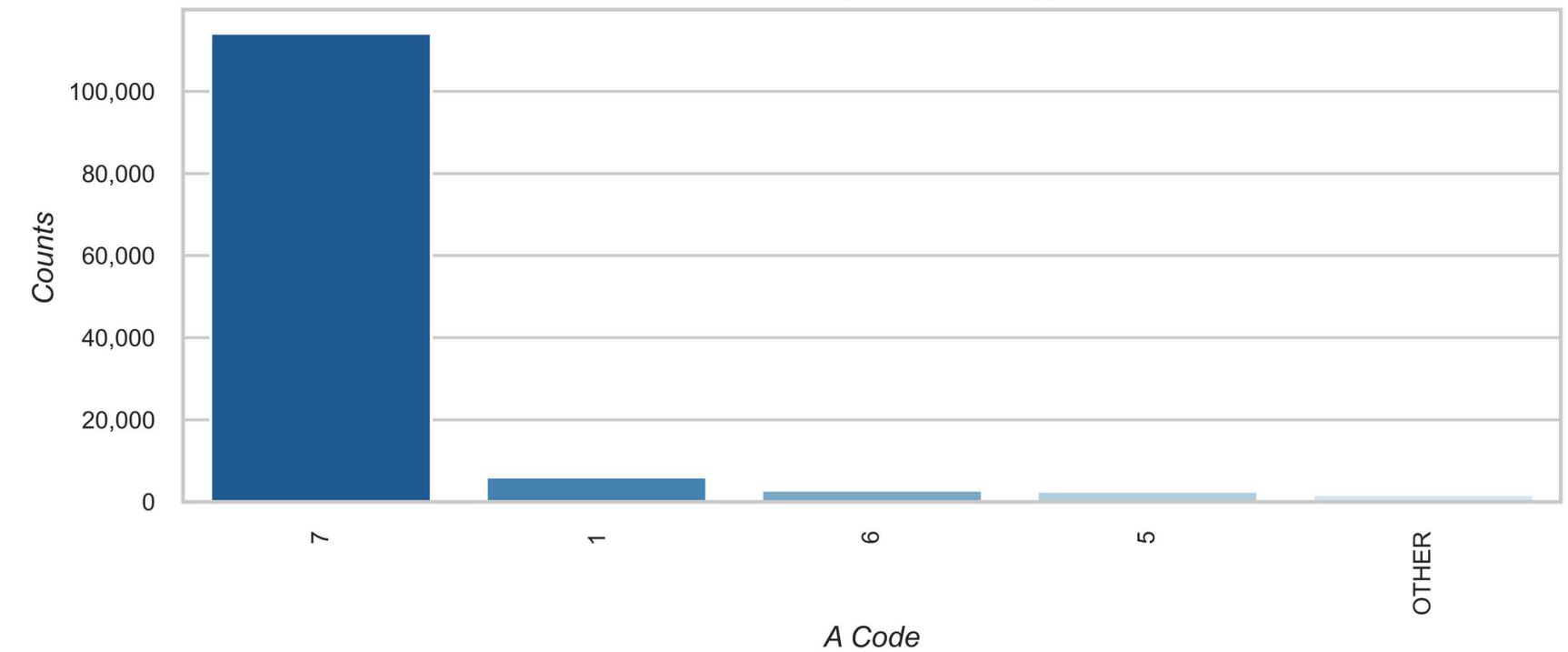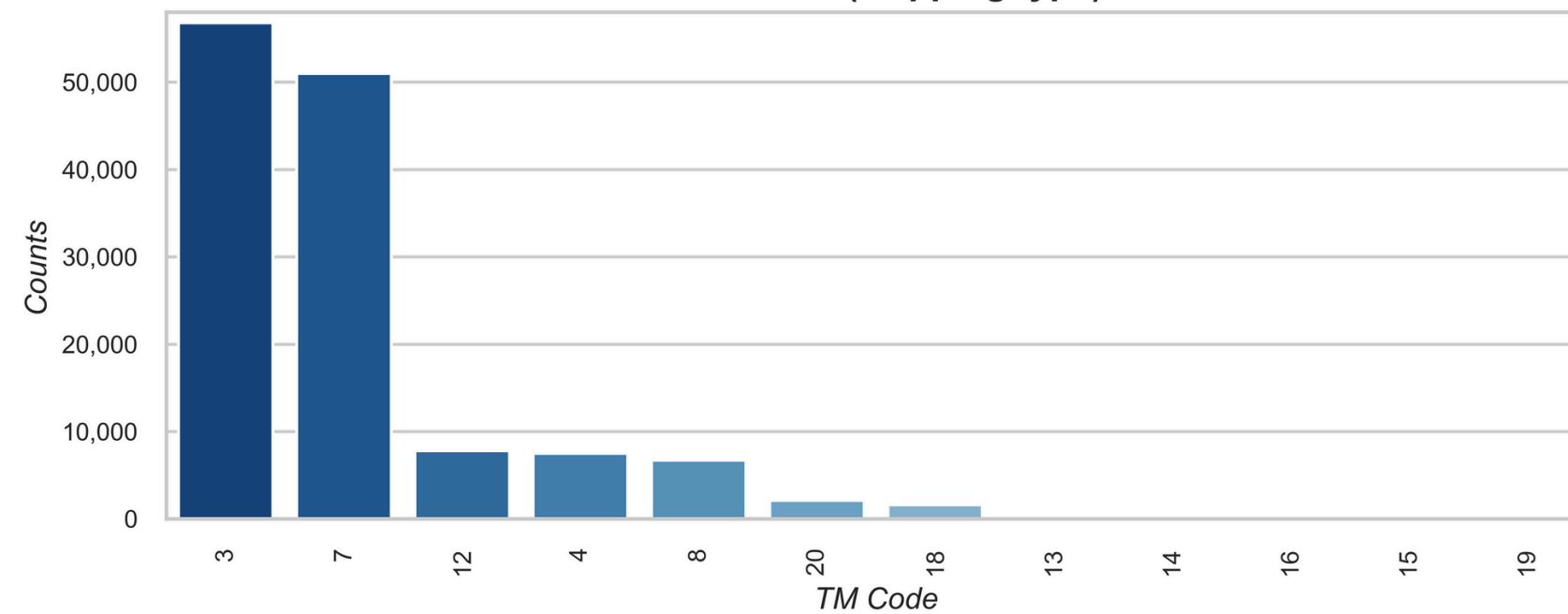**3** The remaining NaN values were filled with the most frequent class within the variable.

Value Counts in A (Business type) - Before
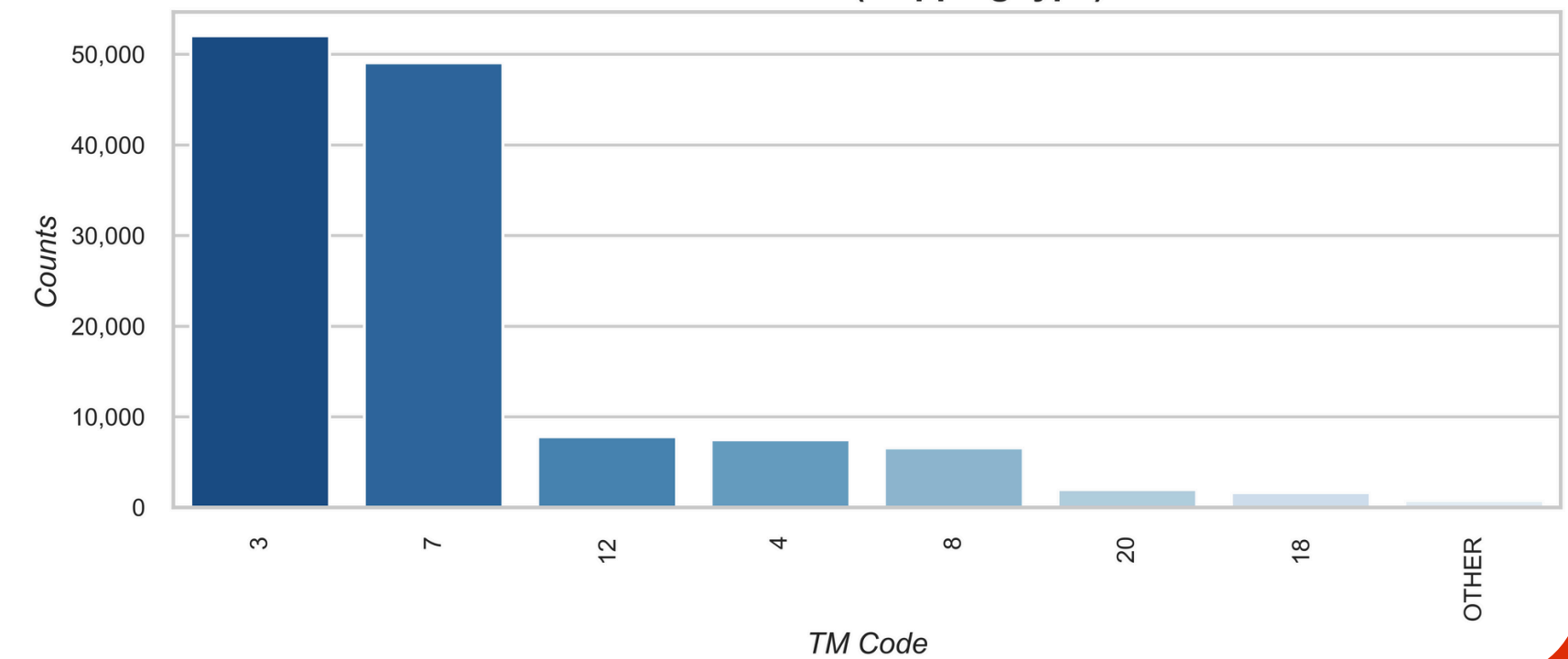
Value Counts in A (Business type) - After

Value Counts in TM (Mapping type) - Before

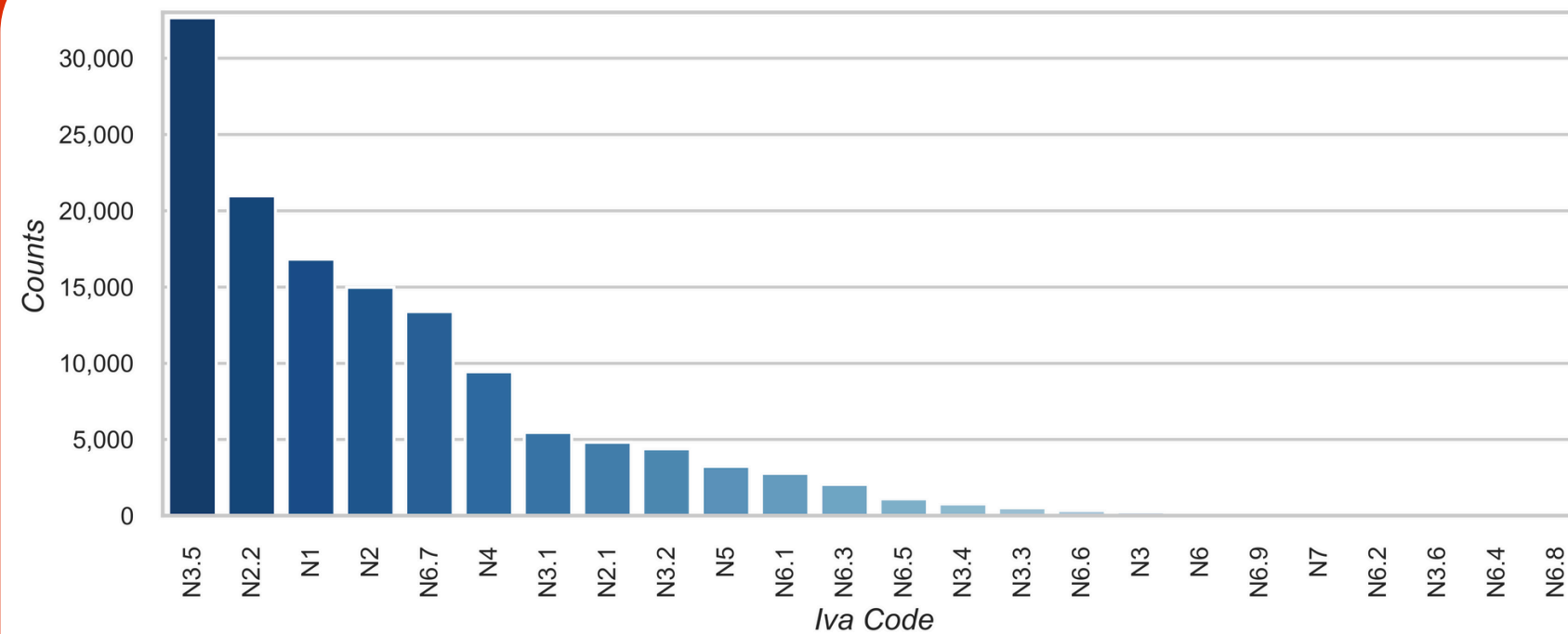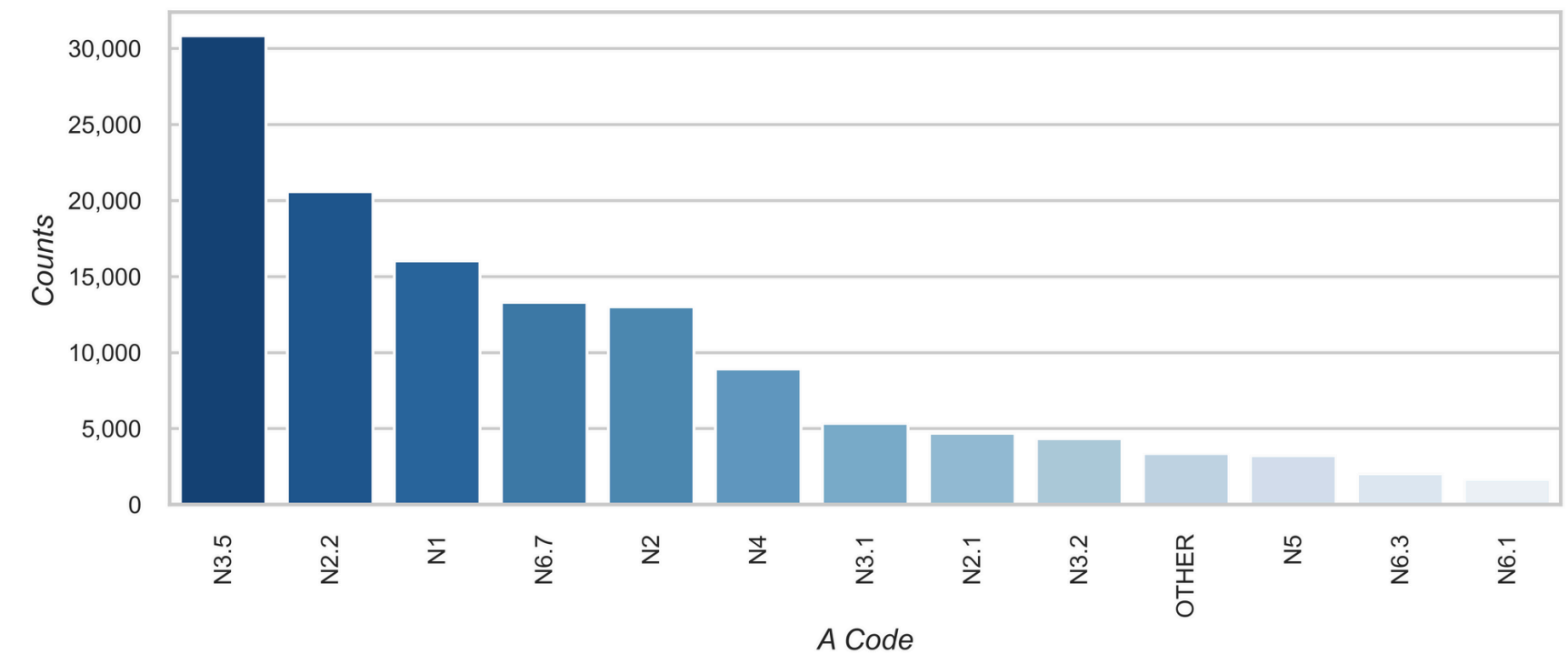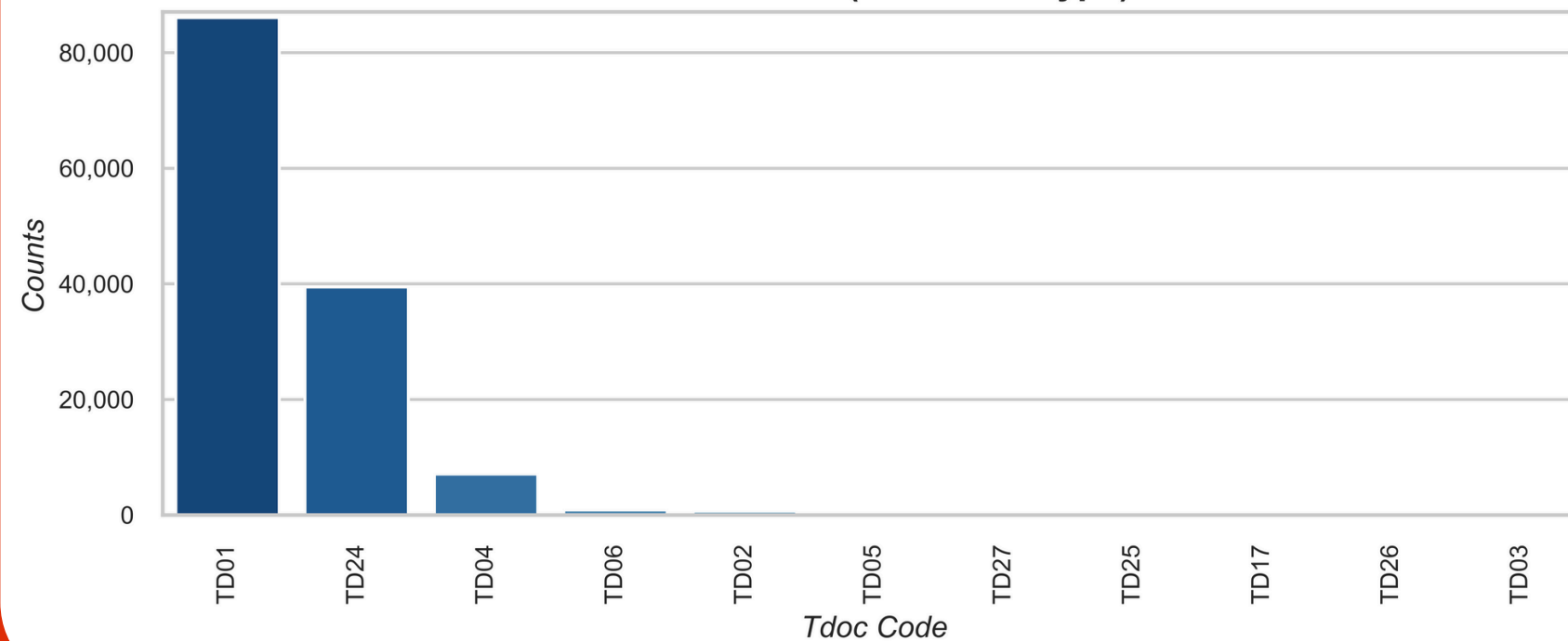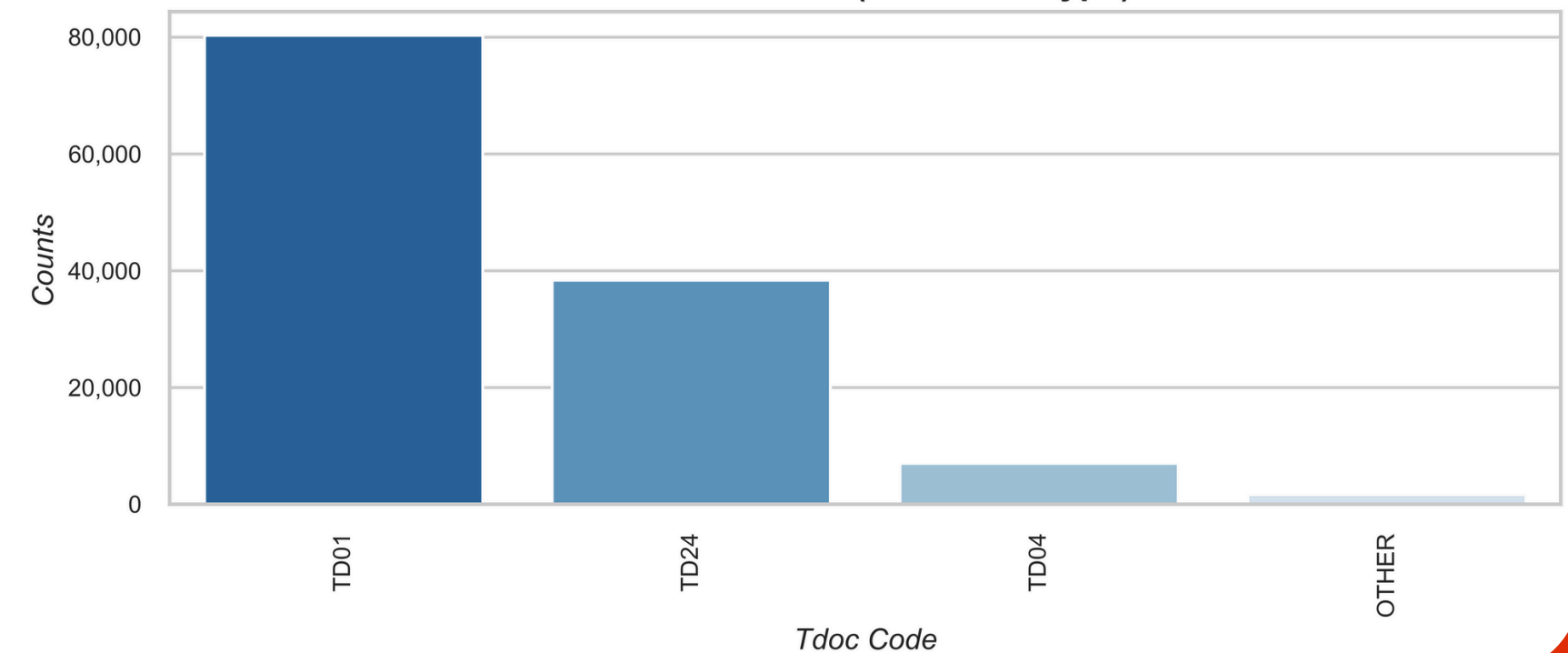Value Counts in TM (Mapping type) - After

Value Counts in Iva Code - Before

Value Counts in Iva Code - After

Value Counts in Tdoc (Document type) - Before

Value Counts in Tdoc (Document type) - After
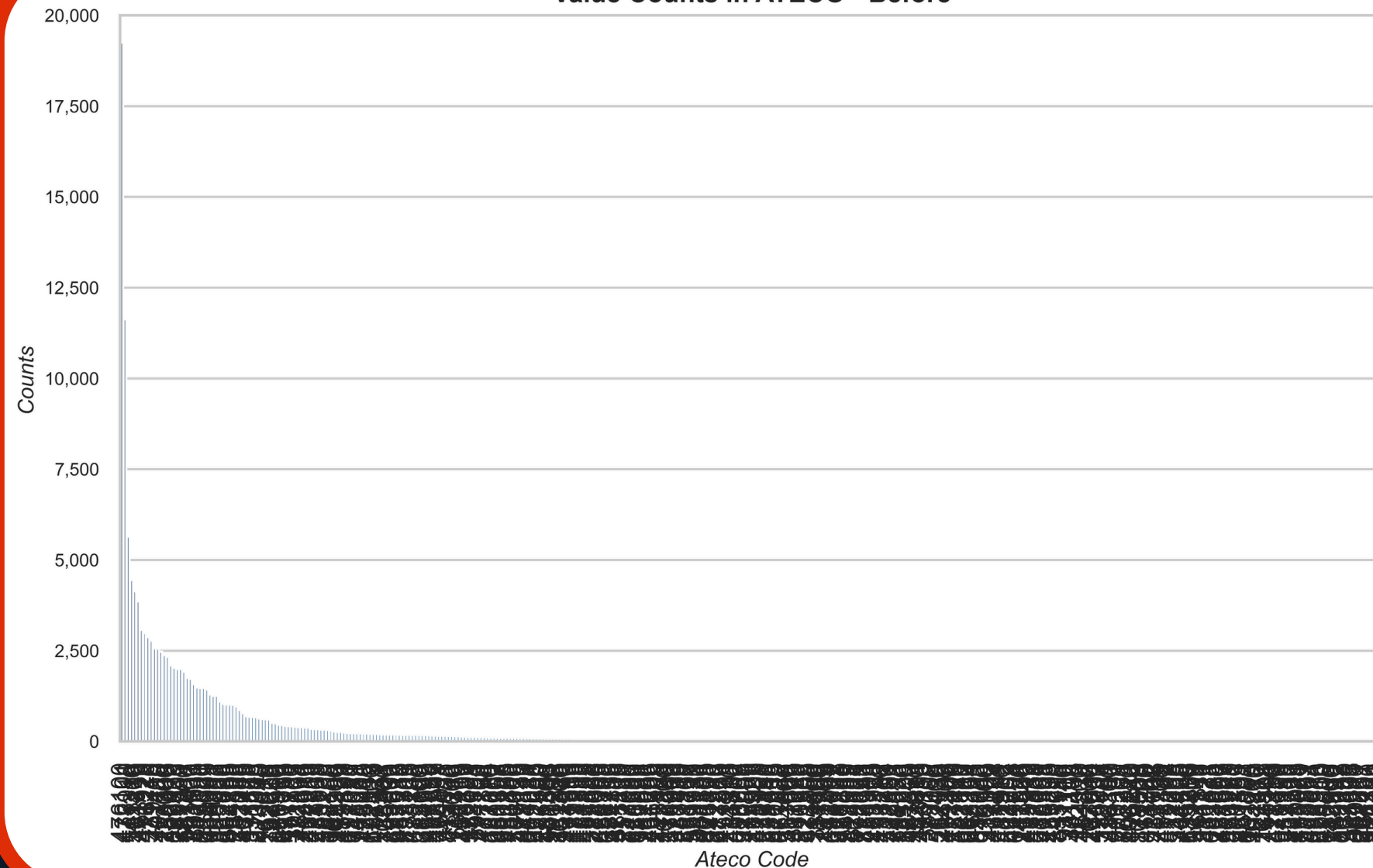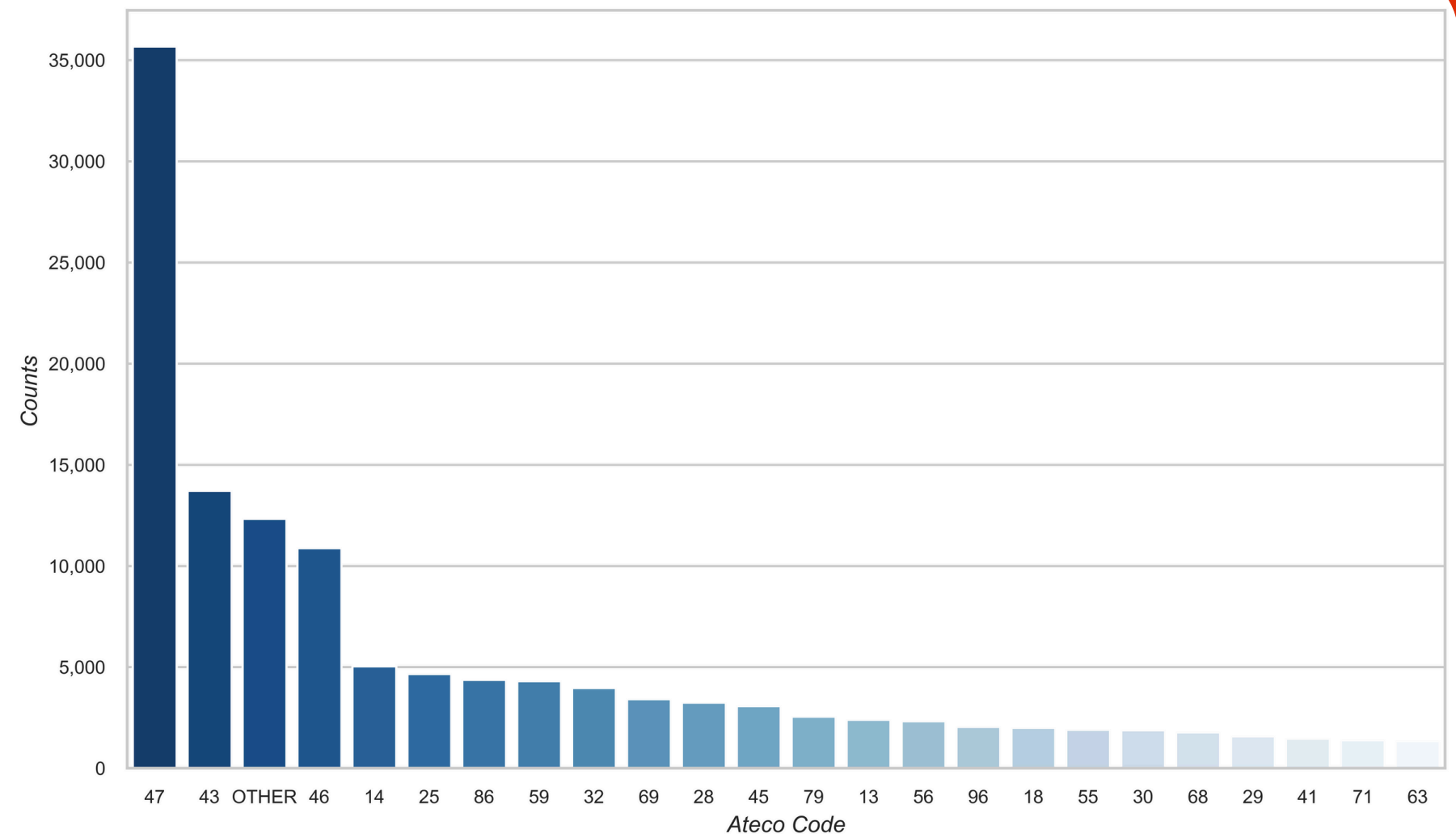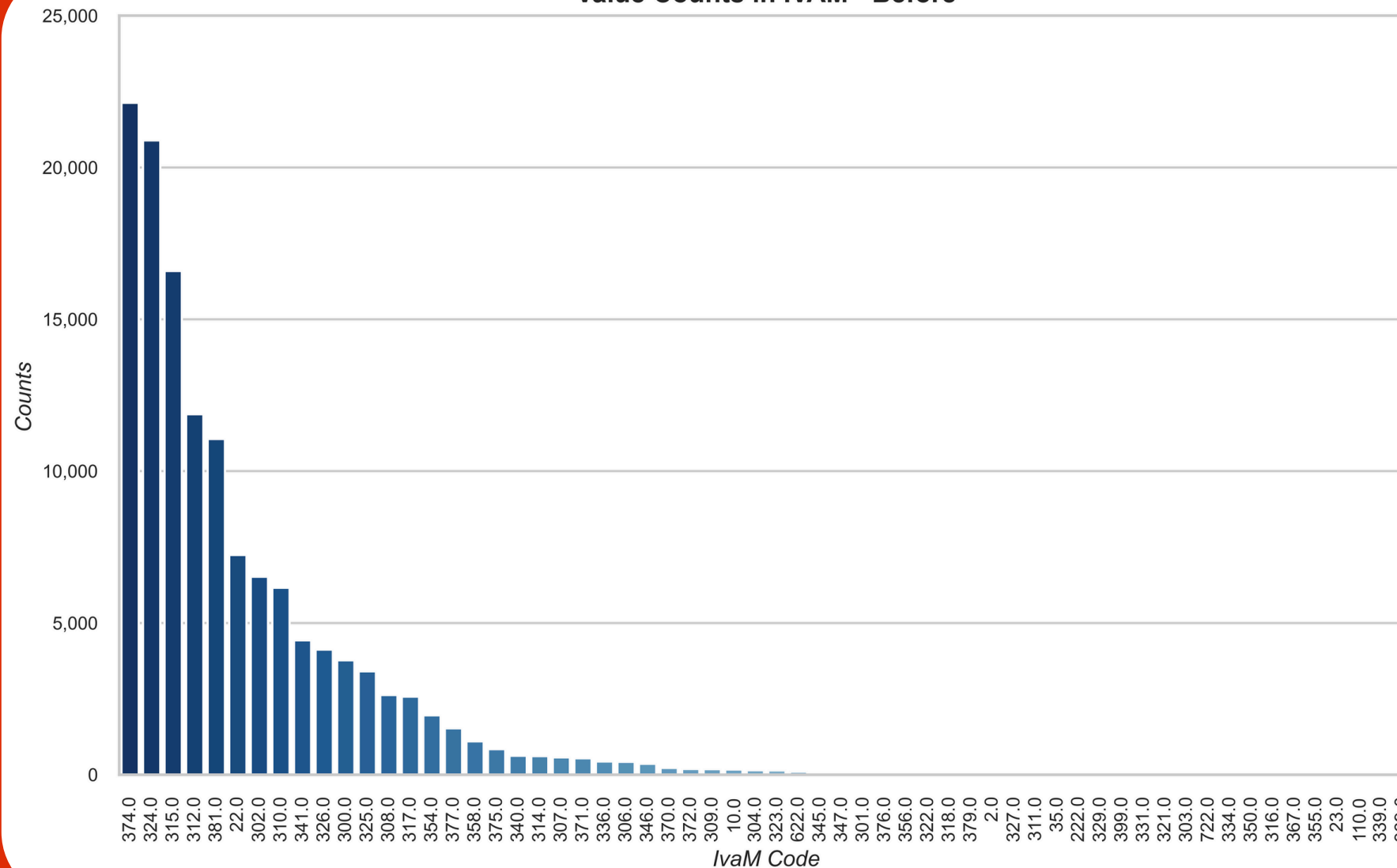
# Focus on Ateco



For the ATECO code, our approach involves two steps: first, extracting the initial two digits of each code, and then consolidating classes with limited observations into an 'others' category.
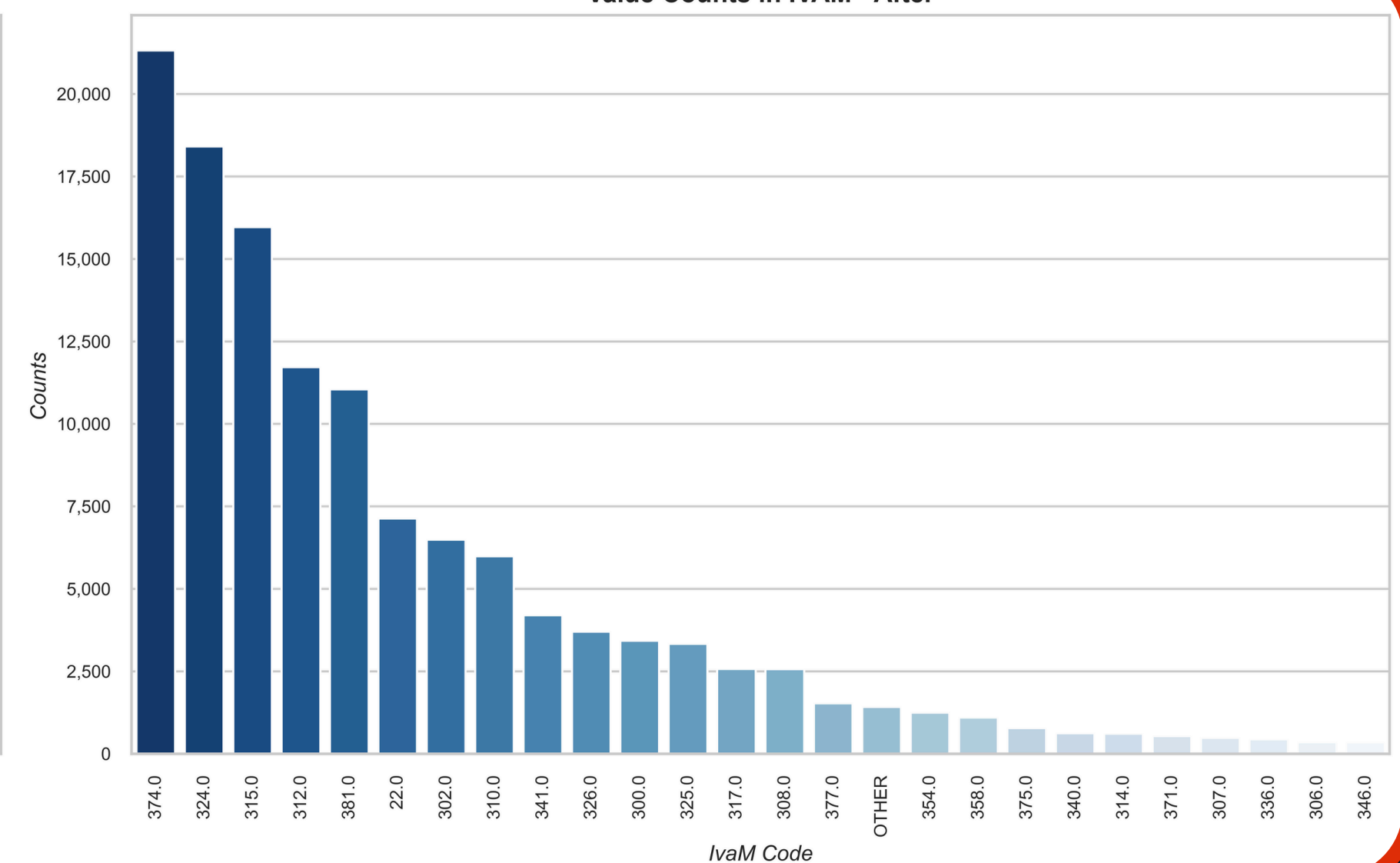
# Focus on IvaM



Value Counts in IVAM - Before

Value Counts in IVAM - After

For our response variable, we opted to group classes with fewer observations, employing a lower threshold of 250 to maintain high sensitivity in the model.

# ③ Model Deployment

# Model comparison



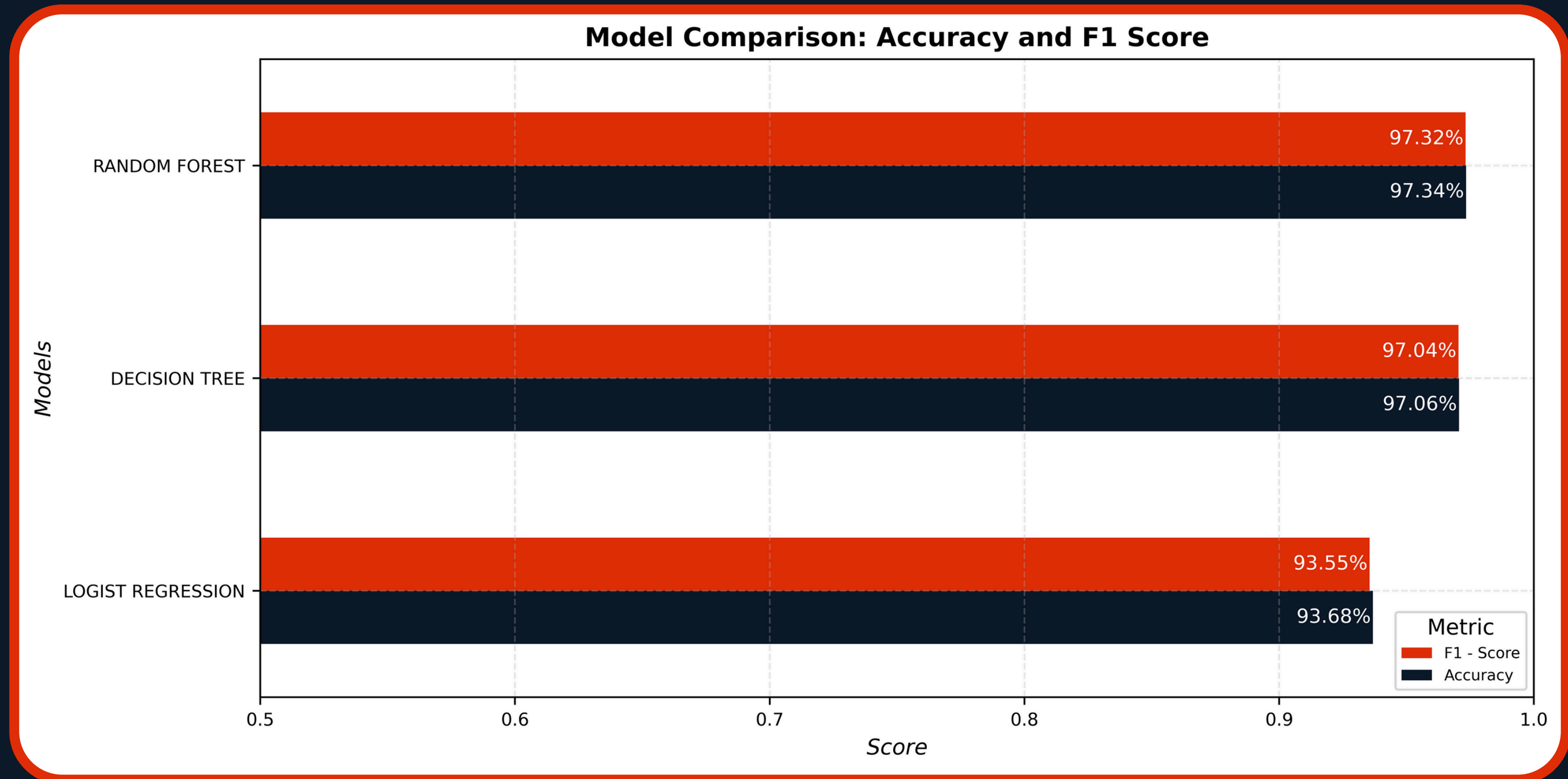**Model Comparison: Accuracy and F1 Score**

RANDOM FOREST
- F1: 97.32%
- Accuracy: 97.34%

DECISION TREE
- F1: 97.04%
- Accuracy: 97.06%

LOGIST REGRESSION
- F1: 93.55%
- Accuracy: 93.68%

Metric
- F1 - Score
- Accuracy

# Final approach - Random Forest



Predictions - Random Forest

4 User Interface

# User interface

Trade-off between precision of the model and usability of the interface

Performed feature selection to to attain the optimal model, leveraging only four features.

http://127.0.0.1:7867

# User interface

5 MLOps

# MLOps - Our Proposal

**Data pipeline**

1. Automation of data collection.

2. Automation of data preprocessing.

3. Implementation of new data to continuously train the model.

# Discussion

Overall we are satisfied with the performance of our model. But we are aware of its inability to predict the exemption code included into the class "Others", however we would be able to fix this by using a more balanced dataset.

# Thank You!

for the Attention

Vincenzo Camerlengo    773731
Daniele De Robertis     787291
Raffaele Torelli        775831

**HERE TO DARE**