

## Step 1: Data Import and Initial Review

### 1.1 Import necessary libraries

```
In [1]: import os
import warnings
import logging
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from typing import Tuple, List, Dict, Optional, Callable, Any
warnings.filterwarnings('ignore')

# Configure Logger
LOG_FORMAT = "%(asctime)s - %(levelname)s - %(message)s"
logging.basicConfig(level=logging.INFO, format=LOG_FORMAT)
logger = logging.getLogger(__name__)
```

```
In [2]: # Define file paths
train_file_path = 'dataset/aps_failure_training_set.csv'
test_file_path = 'dataset/aps_failure_test_set.csv'

# Load the datasets
try:
    # Load training and test datasets
    train_data = pd.read_csv(train_file_path, skiprows=20,na_values=["na"]) # Replace 'na' with NaN for consistency
    test_data = pd.read_csv(test_file_path, skiprows=20,na_values=["na"])

    print("Datasets successfully loaded!")
except FileNotFoundError as e:
    print(f"Error: {e}")
    print("Please ensure the file path is correct.")
```

Datasets successfully loaded!

```
In [3]: # Display basic information about the training dataset
print("\n--- Training Dataset Info ---")
print(train_data.info()) # Summary of columns, data types, and non-null counts

# Display basic information about the test dataset
print("\n--- Test Dataset Info ---")
print(test_data.info()) # Summary of columns, data types, and non-null counts

# Display a few rows from the training dataset to understand its structure
print("\n--- First Five Rows of Training Dataset ---")
```

```
print(train_data.head())

# Check for missing values in both datasets
print("\n--- Missing Values in Training Dataset ---")
print(train_data.isnull().sum())

print("\n--- Missing Values in Test Dataset ---")
print(test_data.isnull().sum())
```

```
--- Training Dataset Info ---  
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 60000 entries, 0 to 59999  
Columns: 171 entries, class to eg_000  
dtypes: float64(169), int64(1), object(1)  
memory usage: 78.3+ MB  
None
```

```
--- Test Dataset Info ---  
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 16000 entries, 0 to 15999  
Columns: 171 entries, class to eg_000  
dtypes: float64(169), int64(1), object(1)  
memory usage: 20.9+ MB  
None
```

```
--- First Five Rows of Training Dataset ---
```

	class	aa_000	ab_000	ac_000	ad_000	ae_000	af_000	ag_000	ag_001	\
0	neg	76698	Nan	2.130706e+09	280.0	0.0	0.0	0.0	0.0	
1	neg	33058	Nan	0.000000e+00	NaN	0.0	0.0	0.0	0.0	
2	neg	41040	Nan	2.280000e+02	100.0	0.0	0.0	0.0	0.0	
3	neg	12	0.0	7.000000e+01	66.0	0.0	10.0	0.0	0.0	
4	neg	60874	Nan	1.368000e+03	458.0	0.0	0.0	0.0	0.0	

	ag_002	...	ee_002	ee_003	ee_004	ee_005	ee_006	ee_007	\
0	0.0	...	1240520.0	493384.0	721044.0	469792.0	339156.0	157956.0	
1	0.0	...	421400.0	178064.0	293306.0	245416.0	133654.0	81140.0	
2	0.0	...	277378.0	159812.0	423992.0	409564.0	320746.0	158022.0	
3	0.0	...	240.0	46.0	58.0	44.0	10.0	0.0	
4	0.0	...	622012.0	229790.0	405298.0	347188.0	286954.0	311560.0	

	ee_008	ee_009	ef_000	eg_000
0	73224.0	0.0	0.0	0.0
1	97576.0	1500.0	0.0	0.0
2	95128.0	514.0	0.0	0.0
3	0.0	0.0	4.0	32.0
4	433954.0	1218.0	0.0	0.0

```
[5 rows x 171 columns]
```

```
--- Missing Values in Training Dataset ---
```

	class	0
aa_000	0	
ab_000	46329	
ac_000	3335	
ad_000	14861	
...		
ee_007	671	
ee_008	671	
ee_009	671	
ef_000	2724	

```
eg_000      2723  
Length: 171, dtype: int64
```

```
--- Missing Values in Test Dataset ---
```

```
class          0  
aa_000        0  
ab_000    12363  
ac_000        926  
ad_000    3981  
...  
ee_007      192  
ee_008      192  
ee_009      192  
ef_000      762  
eg_000      762
```

```
Length: 171, dtype: int64
```

```
In [4]: # Section 2: Initial Exploration
```

```
# Display dataset structure  
print("\n--- Dataset Structure (Training Data) ---")  
print(f"Number of rows: {train_data.shape[0]}, Number of columns: {train_data.shape[1]}")  
print(f"Names of Columns/Features:\n {list(train_data.columns)}")  
  
# Generate summary statistics for numerical features  
print("\n--- Summary Statistics (Numerical Features) ---")  
print(train_data.describe()) # Includes mean, median, std, min, max  
  
# Identify categorical features  
categorical_features = train_data.select_dtypes(include=['object']).columns  
print(f"\n--- Categorical Features ---")  
if len(categorical_features) > 0:  
    print(categorical_features)  
else:  
    print("No categorical features found.")
```

--- Dataset Structure (Training Data) ---  
Number of rows: 60000, Number of columns: 171  
Names of Columns/Features:

```
['class', 'aa_000', 'ab_000', 'ac_000', 'ad_000', 'ae_000', 'af_000', 'ag_000', 'ag_001', 'ag_002', 'ag_003', 'ag_004', 'ag_005', 'ag_006', 'ag_007', 'ag_008', 'ag_009', 'ah_000', 'ai_000', 'aj_000', 'ak_000', 'al_000', 'am_0', 'an_000', 'ao_000', 'ap_000', 'aq_000', 'ar_000', 'as_000', 'at_000', 'au_000', 'av_000', 'ax_000', 'ay_000', 'ay_001', 'ay_002', 'ay_003', 'ay_004', 'ay_005', 'ay_006', 'ay_007', 'ay_008', 'ay_009', 'az_000', 'az_001', 'az_002', 'az_003', 'az_004', 'az_005', 'az_006', 'az_007', 'az_008', 'az_009', 'ba_000', 'ba_001', 'ba_002', 'ba_003', 'ba_004', 'ba_005', 'ba_006', 'ba_007', 'ba_008', 'ba_009', 'bb_000', 'bc_000', 'bd_000', 'be_000', 'bf_000', 'bg_000', 'bh_000', 'bi_000', 'bj_000', 'bk_000', 'bl_000', 'bm_000', 'bn_000', 'bo_000', 'bp_000', 'bq_000', 'br_000', 'bs_000', 'bt_000', 'bu_000', 'bv_000', 'bx_000', 'by_000', 'bz_000', 'ca_000', 'cb_000', 'cc_000', 'cd_000', 'ce_000', 'cf_000', 'cg_000', 'ch_000', 'ci_000', 'cj_000', 'ck_000', 'cl_000', 'cm_000', 'cn_000', 'cn_001', 'cn_002', 'cn_003', 'cn_004', 'cn_005', 'cn_006', 'cn_007', 'cn_008', 'cn_009', 'co_000', 'cp_000', 'cq_000', 'cr_000', 'cs_000', 'cs_001', 'cs_002', 'cs_003', 'cs_004', 'cs_005', 'cs_006', 'cs_007', 'cs_008', 'cs_009', 'ct_000', 'cu_000', 'cv_000', 'cx_000', 'cy_000', 'cz_000', 'da_000', 'db_000', 'dc_000', 'dd_000', 'de_000', 'df_000', 'dg_000', 'dh_000', 'di_000', 'dj_000', 'dk_000', 'dl_000', 'dm_000', 'dn_000', 'do_000', 'dp_000', 'dq_000', 'dr_000', 'ds_000', 'dt_000', 'du_000', 'dv_000', 'dx_000', 'dy_000', 'dz_000', 'ea_000', 'eb_000', 'ec_000', 'ed_000', 'ee_000', 'ee_001', 'ee_002', 'ee_003', 'ee_004', 'ee_005', 'ee_006', 'ee_007', 'ee_008', 'ee_009', 'ef_000', 'eg_000']
```

--- Summary Statistics (Numerical Features) ---

	aa_000	ab_000	ac_000	ad_000	ae_000	\
count	6.000000e+04	13671.000000	5.666500e+04	4.513900e+04	57500.000000	
mean	5.933650e+04	0.713189	3.560143e+08	1.906206e+05	6.819130	
std	1.454301e+05	3.478962	7.948749e+08	4.040441e+07	161.543373	
min	0.000000e+00	0.000000	0.000000e+00	0.000000e+00	0.000000	
25%	8.340000e+02	0.000000	1.600000e+01	2.400000e+01	0.000000	
50%	3.077600e+04	0.000000	1.520000e+02	1.260000e+02	0.000000	
75%	4.866800e+04	0.000000	9.640000e+02	4.300000e+02	0.000000	
max	2.746564e+06	204.000000	2.130707e+09	8.584298e+09	21050.000000	

	af_000	ag_000	ag_001	ag_002	ag_003	\
count	57500.000000	5.932900e+04	5.932900e+04	5.932900e+04	5.932900e+04	
mean	11.006817	2.216364e+02	9.757223e+02	8.606015e+03	8.859128e+04	
std	209.792592	2.047846e+04	3.420053e+04	1.503220e+05	7.617312e+05	
min	0.000000	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	
25%	0.000000	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	
50%	0.000000	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	
75%	0.000000	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	
max	20070.000000	3.376892e+06	4.109372e+06	1.055286e+07	6.340207e+07	

	ee_002	ee_003	ee_004	ee_005	\
count	... 5.932900e+04	5.932900e+04	5.932900e+04	5.932900e+04	
mean	... 4.454897e+05	2.111264e+05	4.457343e+05	3.939462e+05	
std	... 1.155540e+06	5.433188e+05	1.168314e+06	1.121044e+06	
min	... 0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	
25%	... 2.936000e+03	1.166000e+03	2.700000e+03	3.584000e+03	
50%	... 2.337960e+05	1.120860e+05	2.215180e+05	1.899880e+05	
75%	... 4.383960e+05	2.182320e+05	4.666140e+05	4.032220e+05	
max	... 7.793393e+07	3.775839e+07	9.715238e+07	5.743524e+07	

	ee_006	ee_007	ee_008	ee_009	ef_000	\
count	5.932900e+04	5.932900e+04	5.932900e+04	5.932900e+04	57276.000000	
mean	3.330582e+05	3.462714e+05	1.387300e+05	8.388915e+03	0.090579	
std	1.069160e+06	1.728056e+06	4.495100e+05	4.747043e+04	4.368855	
min	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000	

```
25% 5.120000e+02 1.100000e+02 0.000000e+00 0.000000e+00 0.000000
50% 9.243200e+04 4.109800e+04 3.812000e+03 0.000000e+00 0.000000
75% 2.750940e+05 1.678140e+05 1.397240e+05 2.028000e+03 0.000000
max 3.160781e+07 1.195801e+08 1.926740e+07 3.810078e+06 482.000000
```

```
eg_000
count 57277.000000
mean 0.212756
std 8.830641
min 0.000000
25% 0.000000
50% 0.000000
75% 0.000000
max 1146.000000
```

[8 rows x 170 columns]

```
--- Categorical Features ---
Index(['class'], dtype='object')
```

In [5]: # Section 3: Check for Missing Data

```
# Calculate percentage of missing values for each feature
missing_percentage = train_data.isnull().mean() * 100
print("\n--- Percentage of Missing Values (Training Data) ---")
print(missing_percentage.sort_values(ascending=False))

print(type(missing_percentage))
print(missing_percentage.index)
missing_percentages = missing_percentage.sort_values(ascending=False)
print(missing_percentages)
```

```

--- Percentage of Missing Values (Training Data) ---
br_000    82.106667
bq_000    81.203333
bp_000    79.566667
bo_000    77.221667
ab_000    77.215000
...
cj_000    0.563333
ci_000    0.563333
bt_000    0.278333
aa_000    0.000000
class     0.000000
Length: 171, dtype: float64
<class 'pandas.core.series.Series'>
Index(['class', 'aa_000', 'ab_000', 'ac_000', 'ad_000', 'ae_000', 'af_000',
       'ag_000', 'ag_001', 'ag_002',
       ...
       'ee_002', 'ee_003', 'ee_004', 'ee_005', 'ee_006', 'ee_007', 'ee_008',
       'ee_009', 'ef_000', 'eg_000'],
      dtype='object', length=171)
br_000    82.106667
bq_000    81.203333
bp_000    79.566667
bo_000    77.221667
ab_000    77.215000
...
cj_000    0.563333
ci_000    0.563333
bt_000    0.278333
aa_000    0.000000
class     0.000000
Length: 171, dtype: float64

```

```

In [6]: # Plotting a graph showing the top 20 features having highest percentage of missing values
sns.set_style(style="whitegrid")
plt.figure(figsize=(20,5))

# Plot top 20 missing values
missing_percentage = missing_percentage[missing_percentage > 0] # Only show features with missing values
missing_percentage = missing_percentage.sort_values(ascending=False)
data_to_plot = missing_percentage.head(20)

print(data_to_plot)

plot = sns.barplot(x=data_to_plot.index, y=data_to_plot.values)

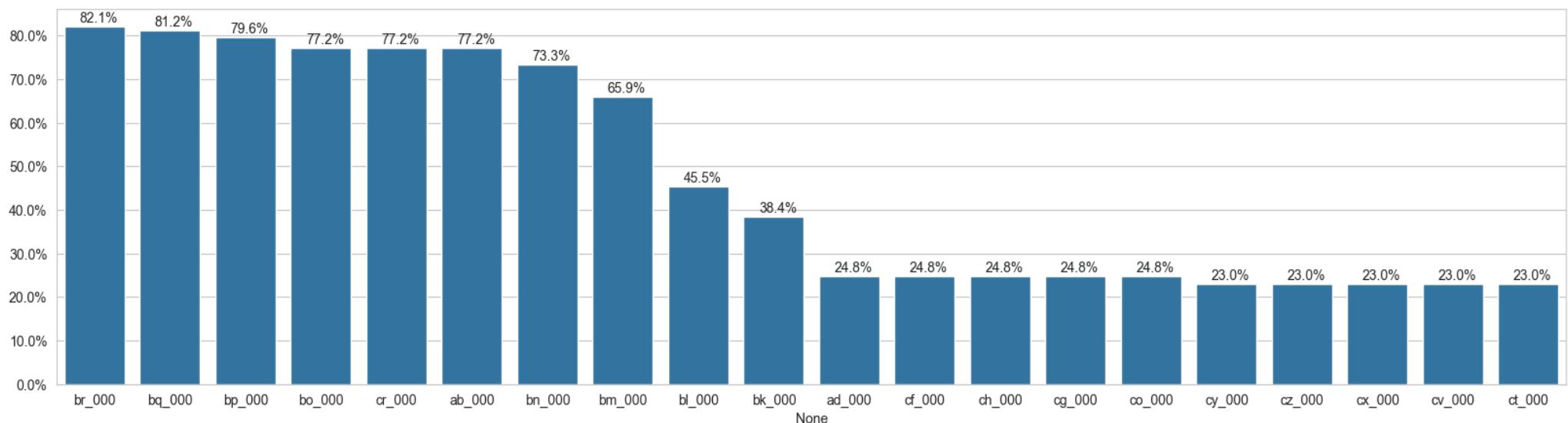
# Add annotations above each bar signifying their value (Shetty, 2021)
for p in plot.patches:
    plot.annotate('{:.1f}%'.format(p.get_height()), (p.get_x()+0.2, p.get_height()+1))

# Make y-axis more interpretable (Shetty, 2021)

```

```
plot.set_yticklabels(map('{:.1f}%'.format, plot.yaxis.get_majorticklocs()))
plt.show()
```

```
br_000    82.106667
bq_000    81.203333
bp_000    79.566667
bo_000    77.221667
cr_000    77.215000
ab_000    77.215000
bn_000    73.348333
bm_000    65.915000
bl_000    45.461667
bk_000    38.390000
ad_000    24.768333
cf_000    24.768333
ch_000    24.768333
cg_000    24.768333
co_000    24.768333
cy_000    23.013333
cz_000    23.013333
cx_000    23.013333
cv_000    23.013333
ct_000    23.013333
dtype: float64
```



In [7]: # Section 4: Class Imbalance Analysis

```
# Calculate class proportions
class_counts = train_data['class'].value_counts()
class_proportions = class_counts / len(train_data) * 100

print("\n--- Class Distribution (Training Data) ---")
print(class_counts)
```

```
print("\n--- Class Proportions (Training Data) ---")
print(class_proportions)
```

```
--- Class Distribution (Training Data) ---
class
```

```
neg    59000
```

```
pos     1000
```

```
Name: count, dtype: int64
```

```
--- Class Proportions (Training Data) ---
class
```

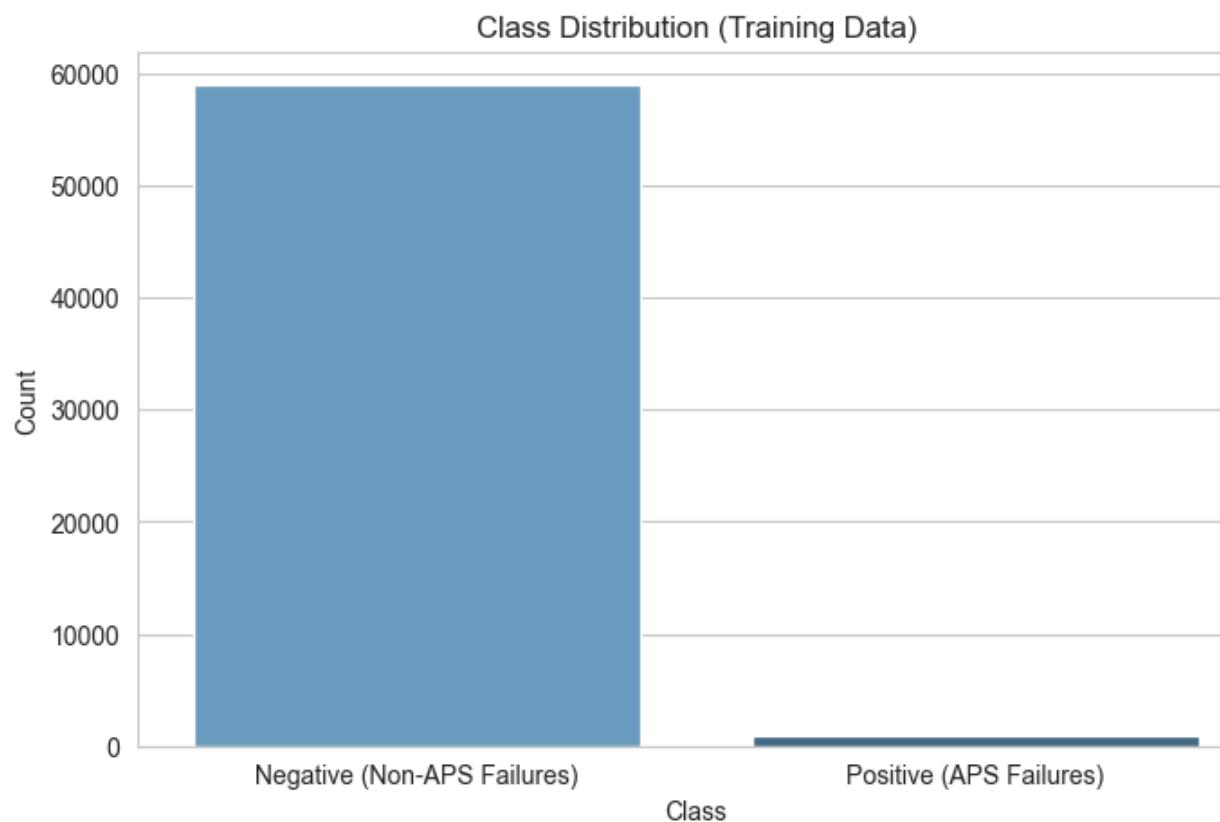
```
neg    98.333333
```

```
pos     1.666667
```

```
Name: count, dtype: float64
```

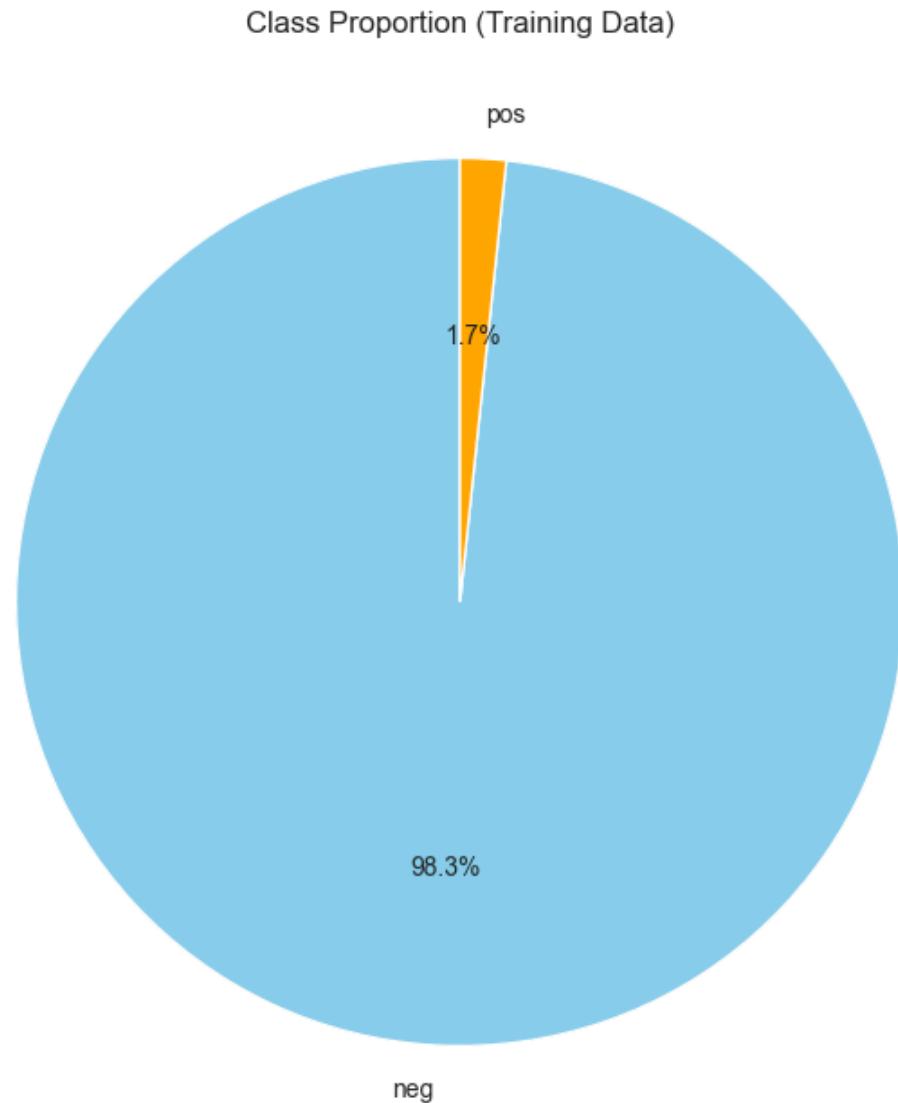
```
In [8]: # Visualize class distribution using a bar plot
```

```
plt.figure(figsize=(8, 5))
sns.barplot(x=class_counts.index, y=class_counts.values, palette="Blues_d")
plt.title("Class Distribution (Training Data)")
plt.ylabel("Count")
plt.xlabel("Class")
plt.xticks(ticks=[0, 1], labels=["Negative (Non-APS Failures)", "Positive (APS Failures)"])
plt.show()
```



In [9]:

```
# Visualize class distribution using a pie chart
plt.figure(figsize=(8, 8))
class_counts.plot.pie(autopct='%1.1f%%', startangle=90, colors=['skyblue', 'orange'])
plt.title("Class Proportion (Training Data)")
plt.ylabel("")
plt.show()
```



## Step 2: Handle Missing Values

In [10]:

```
import json

def get_missing_values(
    data: pd.DataFrame,
    label_columns: List[str],
    low_missing_threshold: int = 5,
    moderate_missing_threshold: int = 15,
    high_missing_threshold: int = 70,
    is_training: bool = False,
    verbose: int = 0,
    save_path: str = "missing_values.json",
) -> Tuple[pd.Series, Dict[str, Tuple[int, List[str]]]]:
    ms_values_dict = {}
    if not is_training:
        logger.warning(
            "Training mode turned off. Loading saved missing values dictionary if available."
        )
        if os.path.exists(save_path):
            with open(save_path, "r") as file:
                json_dict = json.load(file)
            ms_values_dict = {
                k: (threshold, features)
                for k, (threshold, features) in json_dict.items()
            }
        else:
            logger.warning(f"No missing values file found at {save_path}")
            return pd.Series(), {}
    else:
        # Calculate the percentage of missing values, excluding the class column(s)
        missing_percentage = data.drop(columns=label_columns).isnull().mean() * 100

        # Categorize features by the percentage of missing values
        ms_values_dict["low"] = (
            low_missing_threshold,
            list(missing_percentage[missing_percentage < low_missing_threshold].index),
        )
        ms_values_dict["moderate"] = (
            moderate_missing_threshold,
            list(
                missing_percentage[
                    (missing_percentage >= low_missing_threshold)
                    & (missing_percentage < moderate_missing_threshold)
                ].index
            ),
        )
        ms_values_dict["high"] = (
            high_missing_threshold,
            list(
                missing_percentage[
                    (missing_percentage >= moderate_missing_threshold)
                    & (missing_percentage <= high_missing_threshold)
                ].index
            ),
        )
    return pd.Series(), ms_values_dict
```

```

        ].index
    ),
)
ms_values_dict["very_high"] = (
    high_missing_threshold,
    list(missing_percentage[missing_percentage > high_missing_threshold].index)
    + ["cd_000"],
)
# Save dictionary to JSON file
with open(save_path, "w") as file:
    json.dump(ms_values_dict, file)
logger.info(f"Missing values dictionary saved to {save_path}")

# Verbose Logging for insights
if verbose > 0:
    logger.info("\n--- Features Categorized by Missingness ---")
    logger.info(
        f"Low Missing (<{low_missing_threshold}%): "
        f"{len(ms_values_dict['low'][1])} features"
    )
    logger.info(
        f"Moderate Missing ({low_missing_threshold}-{moderate_missing_threshold}%): "
        f"{len(ms_values_dict['moderate'][1])} features"
    )
    logger.info(
        f"High Missing ({moderate_missing_threshold}-{high_missing_threshold}%): "
        f"{len(ms_values_dict['high'][1])} features"
    )
    logger.info(
        f"Very High Missing (>{high_missing_threshold}%): "
        f"{len(ms_values_dict['very_high'][1])} features\n\n"
    )
return ms_values_dict

```

In [11]:

```

import joblib
from sklearn.experimental import enable_iterative_imputer # noqa: F401
from sklearn.impute import SimpleImputer, IterativeImputer

```

```

def remove_missing_values(
    data: pd.DataFrame,
    label_columns: List[str],
    low_missing_threshold: int = 5,
    moderate_missing_threshold: int = 15,
    high_missing_threshold: int = 70,
    is_training: bool = False,
    verbose: int = 0,
    save_path: str = "missing_values.json",
) -> pd.DataFrame:
    """
    """

```

Remove missing values from a DataFrame and apply imputation where necessary.

Args:

```
    data (pd.DataFrame): The input DataFrame.  
    label_columns (List[str]): List of label columns to exclude from missing value operations.  
    low_missing_threshold (int): Threshold for low missing values.  
    moderate_missing_threshold (int): Threshold for moderate missing values.  
    high_missing_threshold (int): Threshold for high missing values.  
    is_training (bool): Whether the function is in training mode.  
    verbose (int): Verbosity level for logging.  
    save_path (str): Path to save or load missing value metadata.
```

Returns:

```
    pd.DataFrame: Cleaned DataFrame with missing values handled.  
    """
```

```
# Assume get_missing_values is defined elsewhere and provides the necessary structure
```

```
missing_values = get_missing_values(  
    data, label_columns, low_missing_threshold, moderate_missing_threshold,  
    high_missing_threshold, is_training, verbose, save_path  
)
```

```
logger.info("Missing values")
```

```
logger.info(f"\n--- Removing Missing Values ---\nInitial Data Shape: {data.shape}")
```

```
# Drop rows with less than 5% missing values
```

```
if missing_values["low"][1]:  
    data_cleaned = data.dropna(subset=missing_values["low"][1], axis=0)  
else:  
    raise ValueError(  
        f"Error dropping row with less than 5% missing values: "  
        f"{missing_values['low'][0]}")  
)
```

```
# Impute for 5%-15% missing values using median/mode imputation
```

```
if is_training:  
    simple_imputer = SimpleImputer(strategy="median")  
    if missing_values["moderate"][1]:  
        data_cleaned[missing_values["moderate"][1]] = simple_imputer.fit_transform(  
            data_cleaned[missing_values["moderate"][1]])  
)
```

```
# Impute for 15%-70% missing values using MICE
```

```
mice_imputer = IterativeImputer(max_iter=10, random_state=42)  
if missing_values["high"][1]:  
    data_cleaned[missing_values["high"][1]] = mice_imputer.fit_transform(  
        data_cleaned[missing_values["high"][1]])  
)
```

```
joblib.dump(simple_imputer, "median_imputer.joblib")
```

```
joblib.dump(mice_imputer, "mice_imputer.joblib")
```

```
else:
```

```

simple_imputer = joblib.load("median_imputer.joblib")
mice_imputer = joblib.load("mice_imputer.joblib")

if missing_values["moderate"][1]:
    data_cleaned[missing_values["moderate"][1]] = simple_imputer.transform(
        data_cleaned[missing_values["moderate"][1]])
)

if missing_values["high"][1]:
    data_cleaned[missing_values["high"][1]] = mice_imputer.transform(
        data_cleaned[missing_values["high"][1]])
)

# Drop features with >70% missing values
if missing_values["very_high"][1]:
    data_cleaned.drop(columns=missing_values["very_high"][1], inplace=True)

if verbose > 0:
    logger.info(
        f"Dropped rows with <{low_missing_threshold}% missing values: "
        f"{missing_values['low'][1]}")
)
logger.info(f"Features with median imputation: {missing_values['moderate'][1]}")
logger.info(f"Features imputed using MICE: {missing_values['high'][1]}")
logger.info(
    f"Dropped features due to high missing values: "
    f"{missing_values['very_high'][1]}")
)
logger.info(f"Final Data Shape: {data_cleaned.shape}")

logger.info(
    f"\n--- Data cleaning complete. Final shape: {data_cleaned.shape} ---\n\n")
)
return data_cleaned

```

In [12]:

```

def remove_outliers(data: pd.DataFrame, label_columns: List[str]) -> (pd.DataFrame, List[str]):
    """
    Removes outliers from a DataFrame by capping values based on z-score thresholds.

    Args:
        data (pd.DataFrame): The input DataFrame containing features and labels.
        label_columns (List[str]): List of columns to exclude from outlier detection.

    Returns:
        pd.DataFrame: DataFrame with outliers capped.
        List[str]: List of numerical feature columns considered for outlier detection.
    """
    # Identify numerical features
    numerical_features = data.drop(columns=label_columns, axis=1).select_dtypes(
        include=["float64", "int64"])
    .columns

```

```

# Use z-scores to identify outliers
z_scores = np.abs(
    (data[numerical_features] - data[numerical_features].mean()) /
    data[numerical_features].std()
)

# Define a threshold for identifying outliers (e.g., z > 3)
outlier_threshold = 3
outliers = (z_scores > outlier_threshold).sum()

logger.info("\n--- Outlier Detection ---")
logger.info(f"Number of outliers per feature (z > {outlier_threshold}): \n{outliers}")

# Decide how to handle outliers (capping in this example)
for feature in numerical_features:
    upper_limit = data[feature].mean() + outlier_threshold * data[feature].std()
    lower_limit = data[feature].mean() - outlier_threshold * data[feature].std()
    data[feature] = np.clip(data[feature], lower_limit, upper_limit)

logger.info(
    f"\n--- Outliers have been capped based on z-score thresholds. "
    f"Final Data Shape: {data.shape} ---\n"
)

```

```
return data, list(numerical_features)
```

```
In [13]: def wrangle_dataset(
    data: pd.DataFrame,
    label_columns: List[str],
    low_missing_threshold: int = 5,
    moderate_missing_threshold: int = 15,
    high_missing_threshold: int = 70,
    is_training: bool = False,
    verbose: int = 0,
    save_path: str = "missing_values.json"
):
    cleaned_data = remove_missing_values(data, label_columns, low_missing_threshold,
                                          moderate_missing_threshold, high_missing_threshold,
                                          is_training, verbose, save_path)

    cleaned_data, numerical_features = remove_outliers(data=cleaned_data, label_columns=label_columns)
    return cleaned_data, numerical_features
```

```
# Ensure train_data is a valid DataFrame before running this
train_data_cleaned, numerical_features = wrangle_dataset(
    data=train_data,
    label_columns=['class'],
    is_training=True,
    verbose=1
```

)

```
train_data_cleaned.info()
```

```
2025-02-09 07:23:35,577 - INFO - Missing values dictionary saved to missing_values.json
2025-02-09 07:23:35,579 - INFO -
--- Features Categorized by Missingness ---
2025-02-09 07:23:35,581 - INFO - Low Missing (<5%): 128 features
2025-02-09 07:23:35,583 - INFO - Moderate Missing (5-15%): 14 features
2025-02-09 07:23:35,585 - INFO - High Missing (15-70%): 21 features
2025-02-09 07:23:35,587 - INFO - Very High Missing (>70%): 8 features

2025-02-09 07:23:35,591 - INFO - Missing_values
2025-02-09 07:23:35,594 - INFO -
--- Removing Missing Values ---
Initial Data Shape: (60000, 171)
2025-02-09 07:23:55,701 - INFO - Dropped rows with <5% missing values: ['aa_000', 'ae_000', 'af_000', 'ag_000', 'ag_001', 'ag_002', 'ag_003', 'ag_004', 'ag_005', 'ag_006', 'ag_007', 'ag_008', 'ag_009', 'ah_000', 'ai_000', 'aj_000', 'al_000', 'am_0', 'an_000', 'ao_000', 'ap_000', 'aq_000', 'ar_000', 'as_000', 'at_000', 'au_000', 'av_000', 'ax_000', 'ay_000', 'ay_001', 'ay_002', 'ay_003', 'ay_004', 'ay_005', 'ay_006', 'ay_007', 'ay_008', 'ay_009', 'az_000', 'az_001', 'az_002', 'az_003', 'az_004', 'az_005', 'az_006', 'az_007', 'az_008', 'az_009', 'ba_000', 'ba_001', 'ba_002', 'ba_003', 'ba_004', 'ba_005', 'ba_006', 'ba_007', 'ba_008', 'ba_009', 'bb_000', 'bc_000', 'bd_000', 'be_000', 'bf_000', 'bg_000', 'bh_000', 'bi_000', 'bj_000', 'bs_000', 'bt_000', 'bu_000', 'bv_000', 'by_000', 'bz_000', 'cb_000', 'cd_000', 'ce_000', 'ci_000', 'cj_000', 'ck_000', 'cn_000', 'cn_001', 'cn_002', 'cn_003', 'cn_004', 'cn_005', 'cn_006', 'cn_007', 'cn_008', 'cn_009', 'cp_000', 'cq_000', 'cs_000', 'cs_001', 'cs_002', 'cs_003', 'cs_004', 'cs_005', 'cs_006', 'cs_007', 'cs_008', 'cs_009', 'dd_000', 'de_000', 'dn_000', 'do_000', 'dp_000', 'dq_000', 'dr_000', 'ds_000', 'dt_000', 'du_000', 'dv_000', 'dx_000', 'dy_000', 'dz_000', 'ea_000', 'ee_000', 'ee_001', 'ee_002', 'ee_003', 'ee_004', 'ee_005', 'ee_006', 'ee_007', 'ee_008', 'ef_000', 'eg_000']
2025-02-09 07:23:55,702 - INFO - Features with median imputation: ['ac_000', 'ak_000', 'bx_000', 'ca_000', 'cc_000', 'df_000', 'dg_000', 'dh_000', 'di_000', 'dj_000', 'dk_000', 'd1_000', 'dm_000', 'eb_000']
2025-02-09 07:23:55,703 - INFO - Features imputed using MICE: ['ad_000', 'bk_000', 'bl_000', 'bm_000', 'cf_000', 'cg_000', 'ch_000', 'cl_000', 'cm_000', 'co_000', 'ct_000', 'cu_000', 'cv_000', 'cx_000', 'cy_000', 'cz_000', 'da_000', 'db_000', 'dc_000', 'ec_00', 'ed_000']
2025-02-09 07:23:55,708 - INFO - Dropped features due to high missing values: ['ab_000', 'bn_000', 'bo_000', 'bp_000', 'bq_000', 'br_000', 'cr_000', 'cd_000']
2025-02-09 07:23:55,709 - INFO - Final Data Shape: (55936, 163)
2025-02-09 07:23:55,710 - INFO -
--- Data cleaning complete. Final shape: (55936, 163) ---
```

```
2025-02-09 07:23:56,399 - INFO -
--- Outlier Detection ---
2025-02-09 07:23:56,401 - INFO - Number of outliers per feature (z > 3):
aa_000      1072
ac_000       0
ad_000       1
ae_000      166
af_000      224
...
ee_007      576
ee_008      654
ee_009      591
ef_000       33
eg_000       77
Length: 162, dtype: int64
2025-02-09 07:23:57,073 - INFO -
--- Outliers have been capped based on z-score thresholds. Final Data Shape: (55936, 163) ---
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 55936 entries, 0 to 59999
Columns: 163 entries, class to eg_000
dtypes: float64(162), object(1)
memory usage: 70.0+ MB
```

```
In [14]: # Example Usage
# Ensure train_data is a valid DataFrame before running this
test_data_cleaned, numerical_features = wrangle_dataset(
    data=test_data,
    label_columns=['class'],
    is_training=False,
    verbose=1
)

test_data_cleaned.info()
```

```
2025-02-09 07:23:57,124 - WARNING - Training mode turned off. Loading saved missing values dictionary if available.
2025-02-09 07:23:57,150 - INFO -
--- Features Categorized by Missingness ---
2025-02-09 07:23:57,152 - INFO - Low Missing (<5%): 128 features
2025-02-09 07:23:57,156 - INFO - Moderate Missing (5-15%): 14 features
2025-02-09 07:23:57,158 - INFO - High Missing (15-70%): 21 features
2025-02-09 07:23:57,160 - INFO - Very High Missing (>70%): 8 features

2025-02-09 07:23:57,161 - INFO - Missing_values
2025-02-09 07:23:57,162 - INFO -
--- Removing Missing Values ---
Initial Data Shape: (16000, 171)
2025-02-09 07:23:57,853 - INFO - Dropped rows with <5% missing values: ['aa_000', 'ae_000', 'af_000', 'ag_000', 'ag_001', 'ag_002', 'ag_003', 'ag_004', 'ag_005', 'ag_006', 'ag_007', 'ag_008', 'ag_009', 'ah_000', 'ai_000', 'aj_000', 'al_000', 'am_0', 'an_000', 'ao_000', 'ap_000', 'aq_000', 'ar_000', 'as_000', 'at_000', 'au_000', 'av_000', 'ax_000', 'ay_000', 'ay_001', 'ay_002', 'ay_003', 'ay_004', 'ay_005', 'ay_006', 'ay_007', 'ay_008', 'ay_009', 'az_000', 'az_001', 'az_002', 'az_003', 'az_004', 'az_005', 'az_006', 'az_007', 'az_008', 'az_009', 'ba_000', 'ba_001', 'ba_002', 'ba_003', 'ba_004', 'ba_005', 'ba_006', 'ba_007', 'ba_008', 'ba_009', 'bb_000', 'bc_000', 'bd_000', 'be_000', 'bf_000', 'bg_000', 'bh_000', 'bi_000', 'bj_000', 'bs_000', 'bt_000', 'bu_000', 'bv_000', 'by_000', 'bz_000', 'cb_000', 'cd_000', 'ce_000', 'ci_000', 'cj_000', 'ck_000', 'cn_000', 'cn_001', 'cn_002', 'cn_003', 'cn_004', 'cn_005', 'cn_006', 'cn_007', 'cn_008', 'cn_009', 'cp_000', 'cq_000', 'cs_000', 'cs_001', 'cs_002', 'cs_003', 'cs_004', 'cs_005', 'cs_006', 'cs_007', 'cs_008', 'cs_009', 'dd_000', 'de_000', 'dn_000', 'do_000', 'dp_000', 'dq_000', 'dr_000', 'ds_000', 'dt_000', 'du_000', 'dv_000', 'dx_000', 'dy_000', 'dz_000', 'ea_000', 'ee_000', 'ee_001', 'ee_002', 'ee_003', 'ee_004', 'ee_005', 'ee_006', 'ee_007', 'ee_008', 'ef_000', 'eg_000']
2025-02-09 07:23:57,856 - INFO - Features with median imputation: ['ac_000', 'ak_000', 'bx_000', 'ca_000', 'cc_000', 'df_000', 'dg_000', 'dh_000', 'di_000', 'dj_000', 'dk_000', 'dl_000', 'dm_000', 'eb_000']
2025-02-09 07:23:57,857 - INFO - Features imputed using MICE: ['ad_000', 'bk_000', 'bl_000', 'bm_000', 'cf_000', 'cg_000', 'ch_000', 'cl_000', 'cm_000', 'co_000', 'ct_000', 'cu_000', 'cv_000', 'cx_000', 'cy_000', 'cz_000', 'da_000', 'db_000', 'dc_000', 'ec_00', 'ed_000']
2025-02-09 07:23:57,859 - INFO - Dropped features due to high missing values: ['ab_000', 'bn_000', 'bo_000', 'bp_000', 'bq_000', 'br_000', 'cr_000', 'cd_000']
2025-02-09 07:23:57,861 - INFO - Final Data Shape: (14891, 163)
2025-02-09 07:23:57,862 - INFO -
--- Data cleaning complete. Final shape: (14891, 163) ---
```

```
2025-02-09 07:23:58,033 - INFO -
--- Outlier Detection ---
2025-02-09 07:23:58,036 - INFO - Number of outliers per feature (z > 3):
aa_000      290
ac_000       0
ad_000       0
ae_000      68
af_000      76
...
ee_007     193
ee_008     115
ee_009     115
ef_000       8
eg_000      11
Length: 162, dtype: int64
2025-02-09 07:23:58,393 - INFO -
--- Outliers have been capped based on z-score thresholds. Final Data Shape: (14891, 163) ---
```

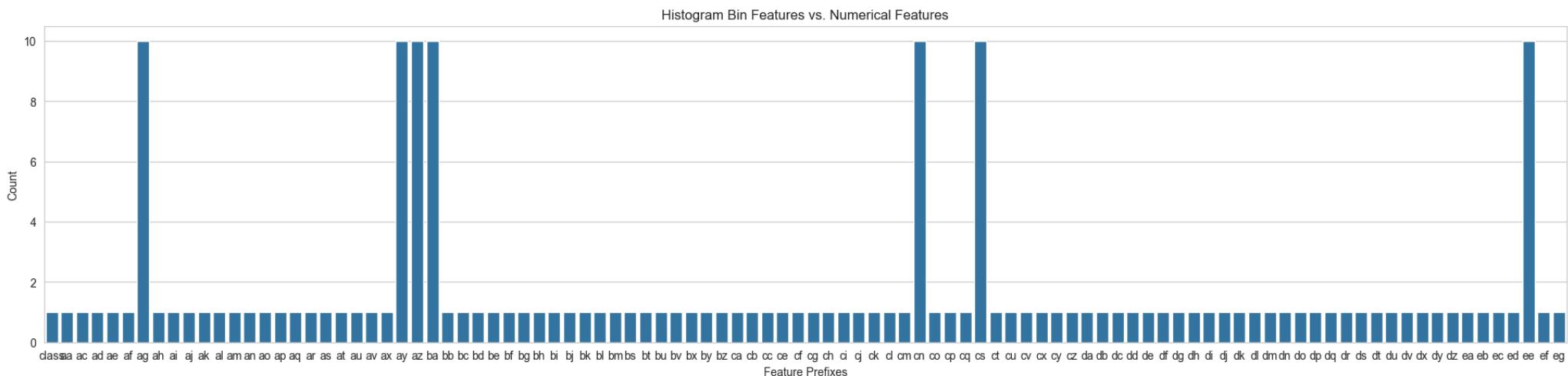
```
<class 'pandas.core.frame.DataFrame'>
Index: 14891 entries, 0 to 15999
Columns: 163 entries, class to eg_000
dtypes: float64(162), object(1)
memory usage: 18.6+ MB
```

## Step 3: Exploratory Data Analysis (EDA)

```
In [15]: from collections import Counter
# Extract feature prefixes
feature_prefix_counts = Counter([name.split('_')[0] for name in train_data_cleaned.columns])
```

```
In [16]: # Get unique feature prefixes and their counts
feature_prefixes = list(feature_prefix_counts.keys())
bin_counts = list(feature_prefix_counts.values())
```

```
In [17]: # Plot the feature distribution
plt.figure(figsize=(24, 5))
sns.barplot(x=feature_prefixes, y=bin_counts)
plt.xlabel("Feature Prefixes")
plt.ylabel("Count")
plt.title("Histogram Bin Features vs. Numerical Features")
plt.show()
```



```
In [18]: # Identify histogram features based on bin count
hist_identifiers = [prefix for prefix, count in zip(feature_prefixes, bin_counts) if count == 10]
print("The Histogram Identifiers are:", hist_identifiers)

# Extract feature names containing histogram bin information
hist_features = [col for col in train_data_cleaned.columns if col.split('_')[0] in hist_identifiers]
print(f"\nThere are {len(hist_features)} features that contain histogram bin information:\n{hist_features}")
```

```
The Histogram Identifiers are: ['ag', 'ay', 'az', 'ba', 'cn', 'cs', 'ee']
```

There are 70 features that contain histogram bin information:

```
['ag_000', 'ag_001', 'ag_002', 'ag_003', 'ag_004', 'ag_005', 'ag_006', 'ag_007', 'ag_008', 'ag_009', 'ay_000', 'ay_001', 'ay_002', 'ay_003', 'ay_004', 'ay_005', 'ay_006', 'ay_007', 'ay_008', 'ay_009', 'az_000', 'az_001', 'az_002', 'az_003', 'az_004', 'az_005', 'az_006', 'az_007', 'az_008', 'az_009', 'ba_000', 'ba_001', 'ba_002', 'ba_003', 'ba_004', 'ba_005', 'ba_006', 'ba_007', 'ba_008', 'ba_009', 'cn_000', 'cn_001', 'cn_002', 'cn_003', 'cn_004', 'cn_005', 'cn_006', 'cn_007', 'cn_008', 'cn_009', 'cs_000', 'cs_001', 'cs_002', 'cs_003', 'cs_004', 'cs_005', 'cs_006', 'cs_007', 'cs_008', 'cs_009', 'ee_000', 'ee_001', 'ee_002', 'ee_003', 'ee_004', 'ee_005', 'ee_006', 'ee_007', 'ee_008', 'ee_009']
```

```
In [19]: # Convert the target column ('class') to binary (0 for 'neg', 1 for 'pos')
train_data_cleaned['class'] = train_data_cleaned['class'].map({'neg':0, 'pos':1})
train_data_cleaned['class'].value_counts()
```

```
Out[19]: class
0    55340
1     596
Name: count, dtype: int64
```

```
In [20]: # Define feature matrix (X) and target variable (y)
X_train = train_data_cleaned.drop('class', axis=1)
y_train = train_data_cleaned['class']
```

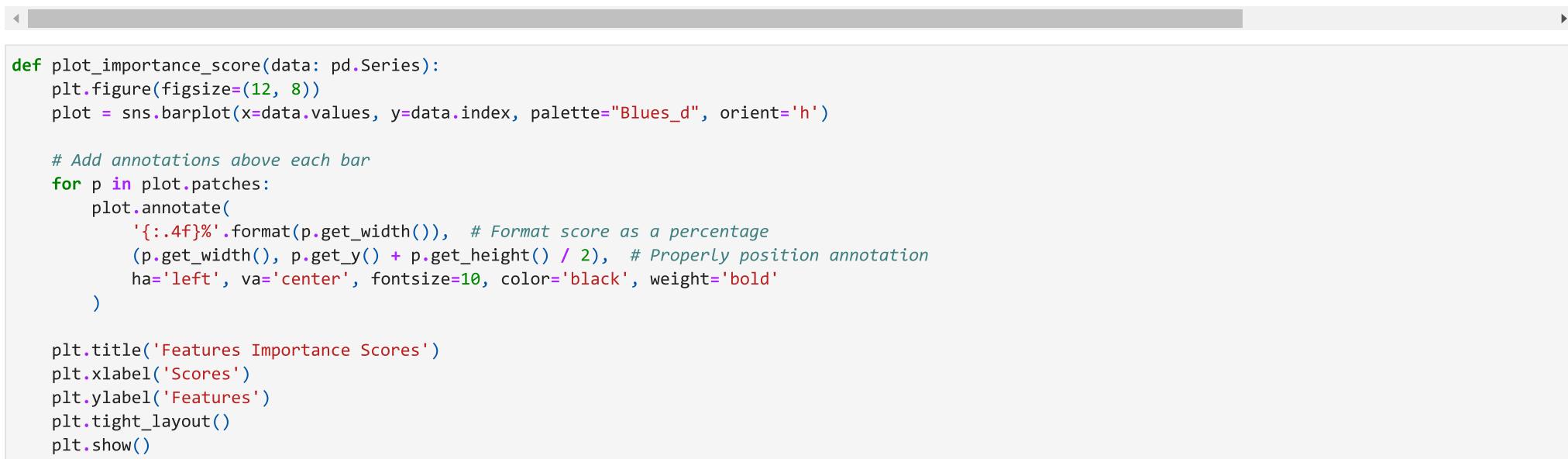
```
In [21]: # Separating the top features from the datasets
histogram_data = X_train[hist_features]
x_without_hist = X_train.drop(hist_features, axis=1)

histogram_data
```

Out[21]:

	ag_000	ag_001	ag_002	ag_003	ag_004	ag_005	ag_006	ag_007	ag_008	ag_009	...	ee_000	ee_001	ee_002	ee_003	ee_004	ee_005
0	0.0	0.0	0.0	0.0	37250.0	1432864.0	3664156.0	1007684.0	25896.0	0.0	...	965866.0	1706908.0	1240520.0	493384.0	721044.0	469792.0
1	0.0	0.0	0.0	0.0	18254.0	653294.0	1720800.0	516724.0	31642.0	0.0	...	664504.0	824154.0	421400.0	178064.0	293306.0	245416.0
2	0.0	0.0	0.0	0.0	1648.0	370592.0	1883374.0	292936.0	12016.0	0.0	...	262032.0	453378.0	277378.0	159812.0	423992.0	409564.0
3	0.0	0.0	0.0	318.0	2212.0	3232.0	1872.0	0.0	0.0	0.0	...	5670.0	1566.0	240.0	46.0	58.0	44.0
4	0.0	0.0	0.0	0.0	43752.0	1966618.0	1800340.0	131646.0	4588.0	0.0	...	404740.0	904230.0	622012.0	229790.0	405298.0	347188.0
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
59995	0.0	0.0	0.0	2564.0	59100.0	1603216.0	6015982.0	1968266.0	164972.0	12560.0	...	1608808.0	1479066.0	998500.0	566884.0	1290398.0	1218244.0
59996	0.0	0.0	0.0	0.0	104.0	99186.0	36564.0	128.0	0.0	0.0	...	13934.0	15024.0	10578.0	6760.0	21126.0	68424.0
59997	0.0	0.0	0.0	0.0	28.0	11592.0	11538.0	0.0	0.0	0.0	...	15876.0	2740.0	792.0	386.0	452.0	144.0
59998	0.0	0.0	0.0	0.0	330.0	202498.0	3251010.0	2061456.0	360436.0	59754.0	...	1180714.0	1709450.0	699352.0	222654.0	347378.0	225724.0
59999	0.0	0.0	0.0	0.0	1226.0	46284.0	1901140.0	855376.0	61744.0	6318.0	...	409798.0	686416.0	440066.0	183200.0	344546.0	254068.0

55936 rows × 70 columns



In [23]:

```
import pandas as pd
import logging
from sklearn.feature_selection import RFE, mutual_info_classif, SelectKBest
from sklearn.tree import DecisionTreeClassifier

def select_top_k_features(data: pd.DataFrame, y: pd.Series, n_selection: int,
                           plot_scores: bool = True, verbose: int = 3) -> pd.Series:
```

```

logger.info("Starting feature selection process...")

# Mutual Information Scores
mutual_info_scores = mutual_info_classif(data, y)
mutual_info_scores = pd.Series(data=mutual_info_scores, index=data.columns, name='Mutual_Info_Score')
mutual_info_scores.sort_values(ascending=False, inplace=True)
logger.info(f"--- Mutual Information Scores ---\n{mutual_info_scores}")

# Recursive Feature Elimination (RFE) with Decision Tree
estimator = DecisionTreeClassifier(random_state=42)
rfe_selector = RFE(
    estimator=estimator,
    n_features_to_select=n_selection,
    step=1,
    verbose=verbose
)
rfe_selector.fit(data, y)

# Get top selected features from RFE
rfe_selected_features = data.columns[rfe_selector.support_].tolist()
logger.info(f"\n--- Top Features Selected by RFE ---\n{rfe_selected_features}")

# SelectKBest for Mutual Information Scores
kbest_selector = SelectKBest(mutual_info_classif, k=n_selection)
kbest_selector.fit(data, y)
selected_kbest_cols = data.columns[kbest_selector.get_support()]
logger.info(f"\n--- Top Features Selected by SelectKBest ---\n{selected_kbest_cols}")

# Combine selected features from both methods
selected_features = list(set(rfe_selected_features).union(selected_kbest_cols))

logger.info(f"\n--- Combined Selected Features ---\nTotal: {len(selected_features)}\n{selected_features}")

selected_features_df = mutual_info_scores.loc[selected_features]
if plot_scores:
    if "plot_importance_score" in globals():
        plot_importance_score(selected_features_df)
    else:
        logger.warning("plot_importance_score function is not defined!")

return selected_features_df

# Call the function
selected_features_df = select_top_k_features(data=histogram_data, y=y_train, n_selection=15)
selected_features = selected_features_df.index.tolist()

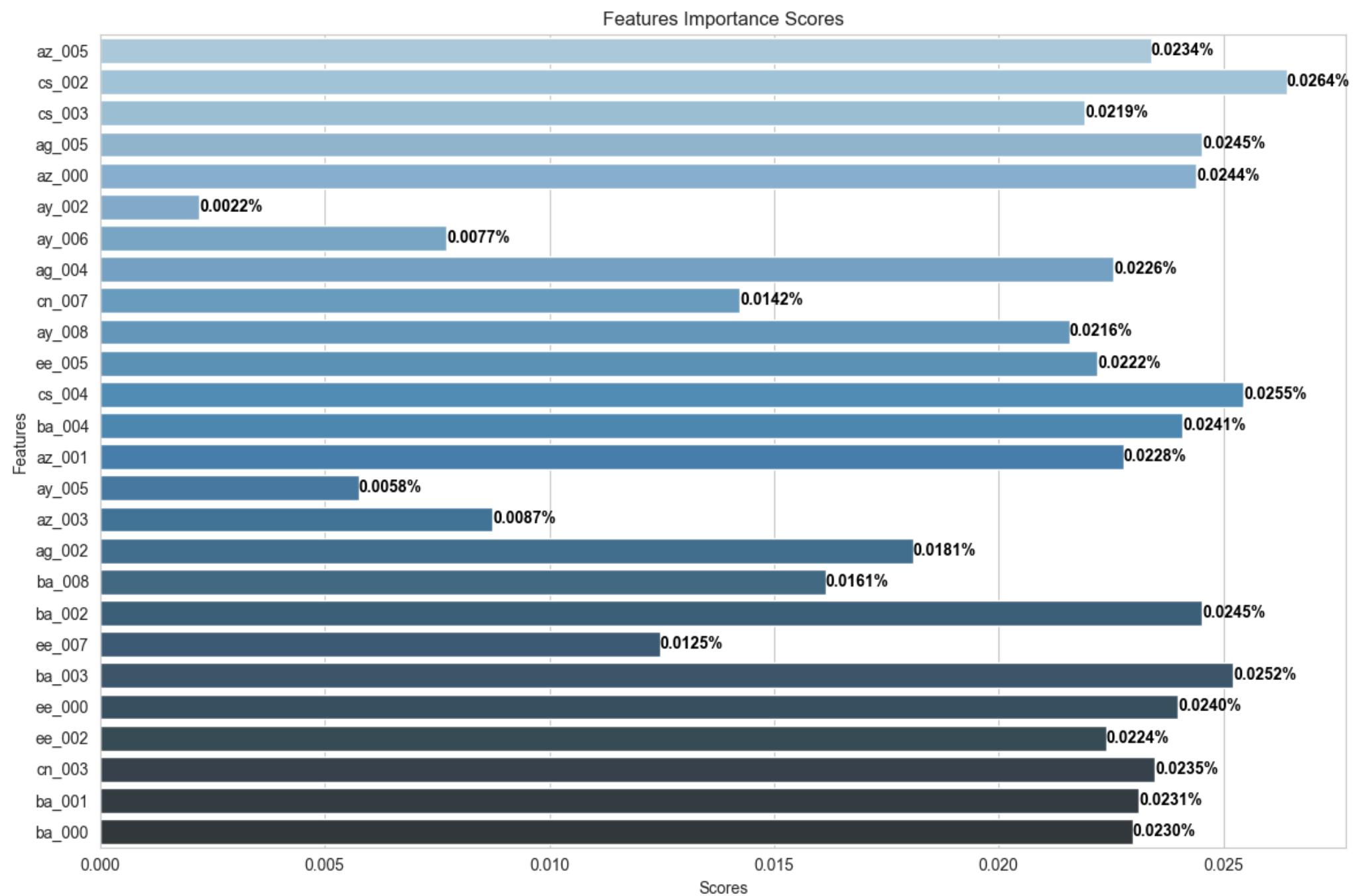
```

```
2025-02-09 07:24:01,109 - INFO - Starting feature selection process...
2025-02-09 07:24:45,263 - INFO - --- Mutual Information Scores ---
cs_002      0.026409
cs_004      0.025452
ba_003      0.025214
ba_002      0.024534
ag_005      0.024531
...
cs_008      0.002008
az_009      0.001608
ay_001      0.001363
ay_000      0.000803
cs_009      0.000000
Name: Mutual_Info_Score, Length: 70, dtype: float64
```

Fitting estimator with 70 features.  
Fitting estimator with 69 features.  
Fitting estimator with 68 features.  
Fitting estimator with 67 features.  
Fitting estimator with 66 features.  
Fitting estimator with 65 features.  
Fitting estimator with 64 features.  
Fitting estimator with 63 features.  
Fitting estimator with 62 features.  
Fitting estimator with 61 features.  
Fitting estimator with 60 features.  
Fitting estimator with 59 features.  
Fitting estimator with 58 features.  
Fitting estimator with 57 features.  
Fitting estimator with 56 features.  
Fitting estimator with 55 features.  
Fitting estimator with 54 features.  
Fitting estimator with 53 features.  
Fitting estimator with 52 features.  
Fitting estimator with 51 features.  
Fitting estimator with 50 features.  
Fitting estimator with 49 features.  
Fitting estimator with 48 features.  
Fitting estimator with 47 features.  
Fitting estimator with 46 features.  
Fitting estimator with 45 features.  
Fitting estimator with 44 features.  
Fitting estimator with 43 features.  
Fitting estimator with 42 features.  
Fitting estimator with 41 features.  
Fitting estimator with 40 features.  
Fitting estimator with 39 features.  
Fitting estimator with 38 features.  
Fitting estimator with 37 features.  
Fitting estimator with 36 features.  
Fitting estimator with 35 features.  
Fitting estimator with 34 features.  
Fitting estimator with 33 features.  
Fitting estimator with 32 features.  
Fitting estimator with 31 features.  
Fitting estimator with 30 features.  
Fitting estimator with 29 features.  
Fitting estimator with 28 features.  
Fitting estimator with 27 features.  
Fitting estimator with 26 features.  
Fitting estimator with 25 features.  
Fitting estimator with 24 features.  
Fitting estimator with 23 features.  
Fitting estimator with 22 features.  
Fitting estimator with 21 features.  
Fitting estimator with 20 features.

```
Fitting estimator with 19 features.  
Fitting estimator with 18 features.  
Fitting estimator with 17 features.  
Fitting estimator with 16 features.
```

```
2025-02-09 07:32:58,964 - INFO -  
--- Top Features Selected by RFE ---  
['ag_002', 'ag_004', 'ay_002', 'ay_005', 'ay_006', 'ay_008', 'az_001', 'az_003', 'ba_002', 'ba_008', 'cn_007', 'cs_002', 'ee_002', 'ee_005', 'ee_007']  
2025-02-09 07:33:40,786 - INFO -  
--- Top Features Selected by SelectKBest ---  
Index(['ag_004', 'ag_005', 'az_000', 'az_001', 'az_005', 'ba_000', 'ba_001',  
       'ba_002', 'ba_003', 'ba_004', 'cn_003', 'cs_002', 'cs_003', 'cs_004',  
       'ee_000'],  
      dtype='object')  
2025-02-09 07:33:40,797 - INFO -  
--- Combined Selected Features ---  
Total: 26  
['az_005', 'cs_002', 'cs_003', 'ag_005', 'az_000', 'ay_002', 'ay_006', 'ag_004', 'cn_007', 'ay_008', 'ee_005', 'cs_004', 'ba_004', 'az_001', 'ay_005', 'az_003', 'ag_002', 'ba_008', 'ba_002', 'ee_007', 'ba_003', 'ee_000', 'ee_002', 'cn_003', 'ba_001', 'ba_000']
```



### Binned Selected Features Analysis

In [24]: `def univariate_analysis(df, target_column):`

`'''`

        Perform univariate analysis on selected numerical features.

```

- Print mean and standard deviation for each class.
- Plot PDF, CDF, and boxplots for each feature.
"""

for feature in df.columns:
    if feature != target_column:
        # Descriptive statistics by class
        describe_0 = df[df[target_column] == 0][feature].describe()
        describe_1 = df[df[target_column] == 1][feature].describe()

        print(f"\nFeature: {feature}")
        print(f"Class 0 - Mean: {round(describe_0['mean'], 2)}, Std Dev: {round(describe_0['std'], 2)}")
        print(f"Class 1 - Mean: {round(describe_1['mean'], 2)}, Std Dev: {round(describe_1['std'], 2)})")

        # Plot PDF, CDF, and boxplot
        fig, ax = plt.subplots(1, 3, figsize=(15, 4))
        sns.kdeplot(df[df[target_column] == 0][feature], ax=ax[0], label="Class 0", shade=True)
        sns.kdeplot(df[df[target_column] == 1][feature], ax=ax[0], label="Class 1", shade=True)
        ax[0].set_title(f"PDF of {feature}")

        sns.ecdfplot(data=df, x=feature, hue=target_column, ax=ax[1])
        ax[1].set_title(f"CDF of {feature}")

        sns.boxplot(x=target_column, y=feature, data=df, ax=ax[2])
        ax[2].set_title(f"Boxplot of {feature}")

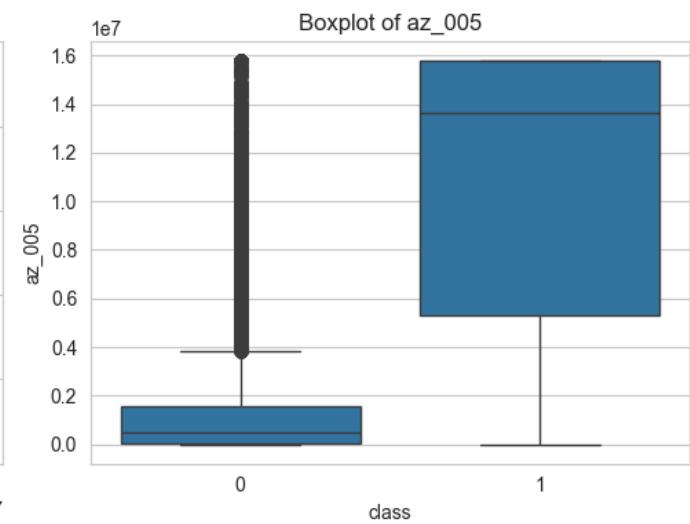
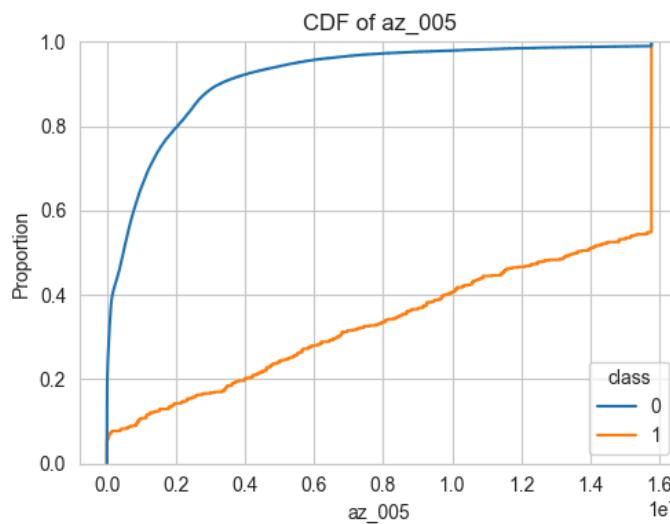
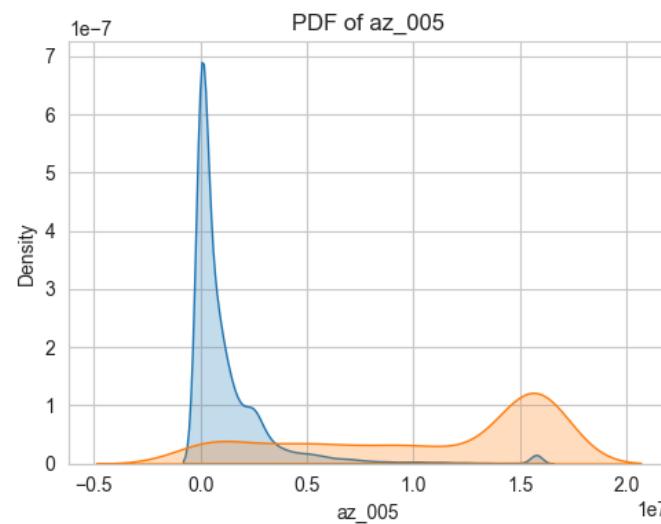
    plt.tight_layout()
    plt.show()

# Perform univariate analysis on selected features
univariate_analysis(train_data_cleaned[selected_features + ['class']], 'class')

```

Feature: az\_005  
 Class 0 - Mean: 1327084.52, Std Dev: 2412810.4  
 Class 1 - Mean: 10594228.32, Std Dev: 5859895.62

2025-02-09 07:33:43,432 - INFO - Using categorical units to plot a list of strings that are all parseable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.  
 2025-02-09 07:33:43,585 - INFO - Using categorical units to plot a list of strings that are all parseable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



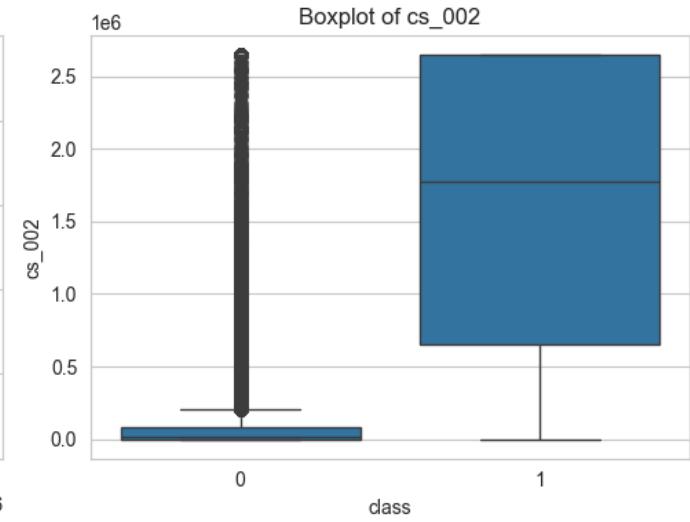
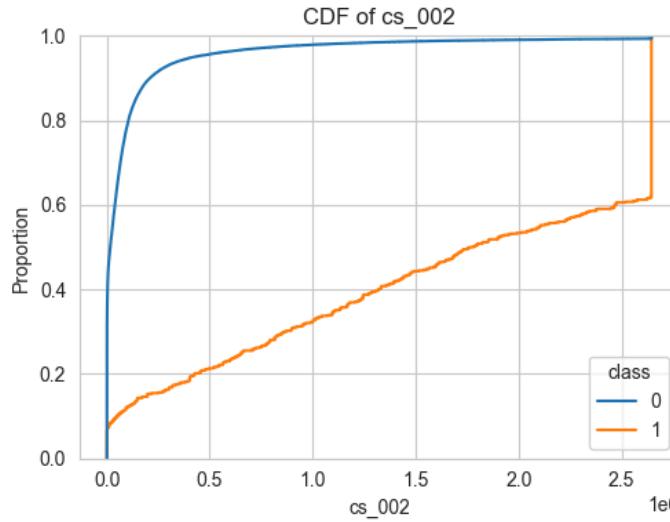
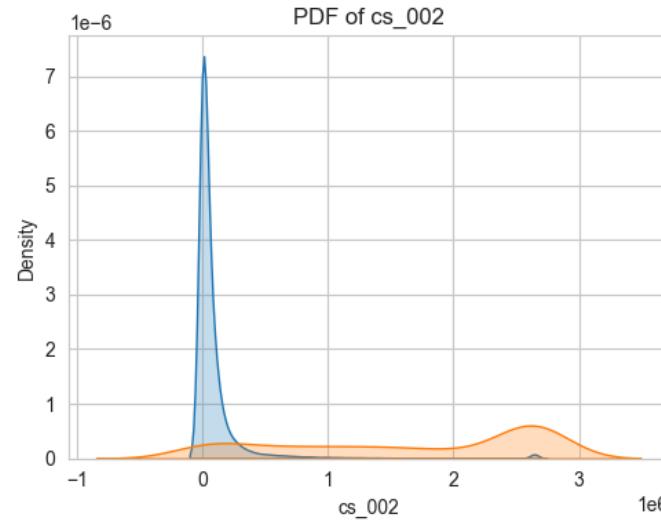
Feature: cs\_002

Class 0 - Mean: 104252.27, Std Dev: 302517.17

Class 1 - Mean: 1624227.23, Std Dev: 1012979.6

2025-02-09 07:33:45,725 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:33:45,778 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



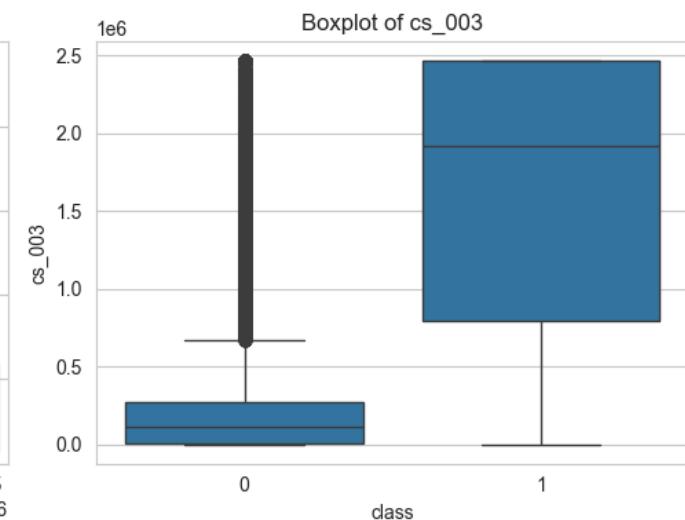
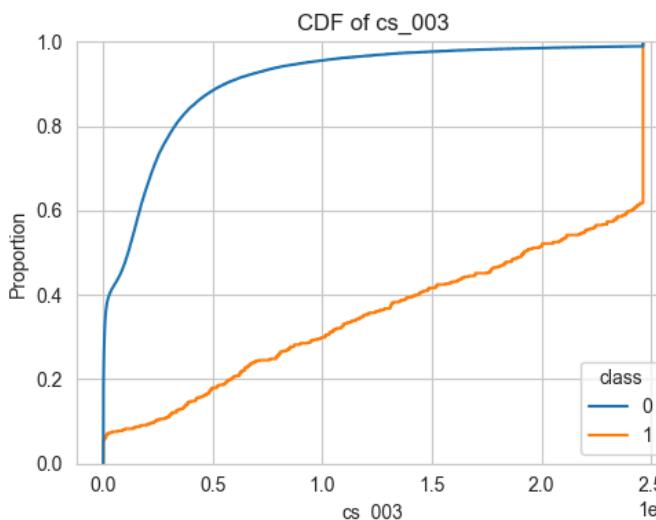
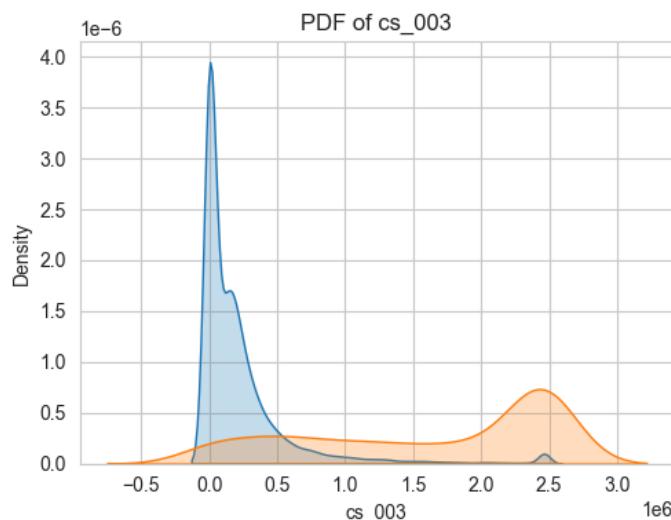
Feature: cs\_003

Class 0 - Mean: 227171.95, Std Dev: 388066.59

Class 1 - Mean: 1615576.08, Std Dev: 900390.88

2025-02-09 07:33:47,594 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:33:47,660 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



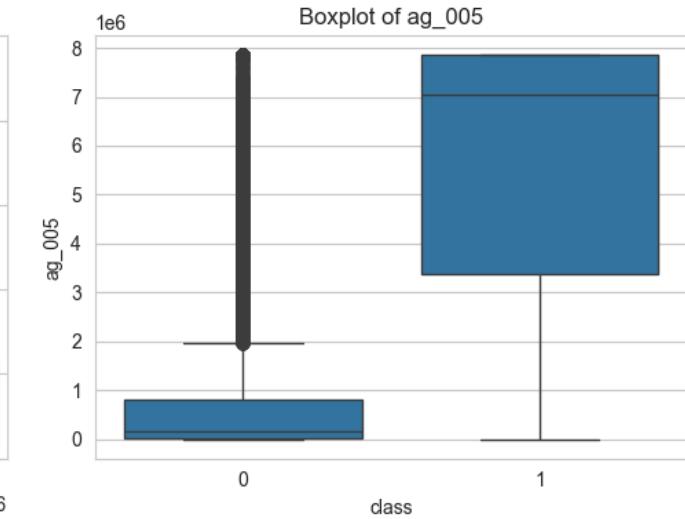
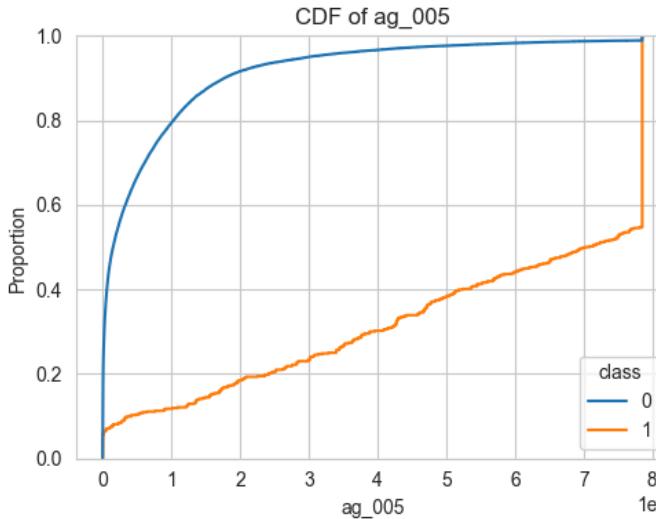
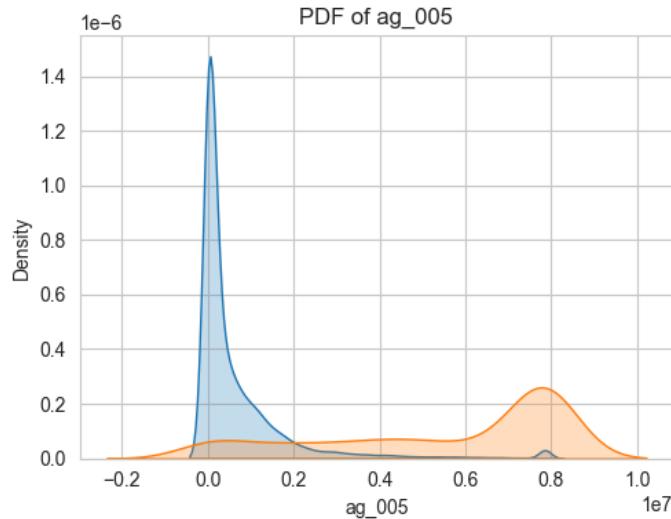
Feature: ag\_005

Class 0 - Mean: 673667.09, Std Dev: 1268774.82

Class 1 - Mean: 5467078.98, Std Dev: 2821820.41

2025-02-09 07:33:49,806 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:33:49,890 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



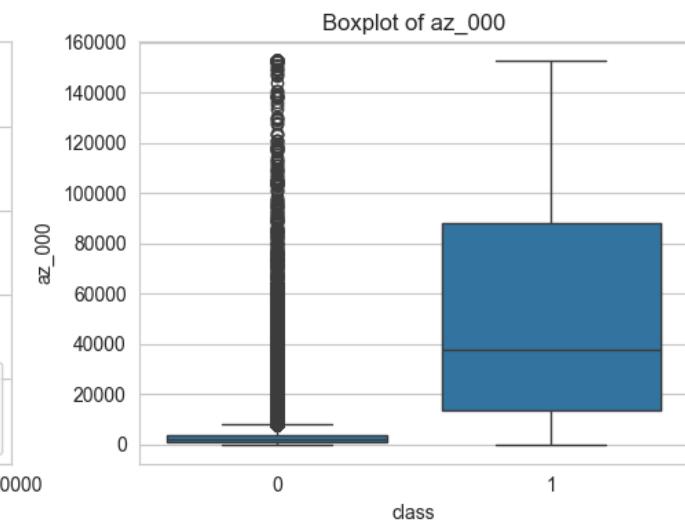
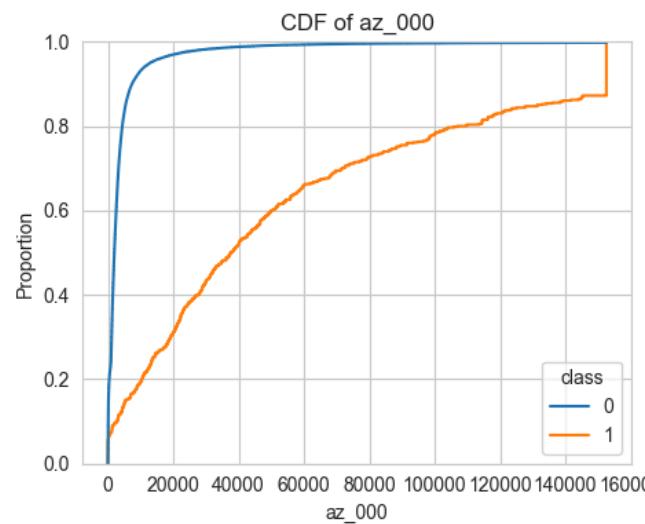
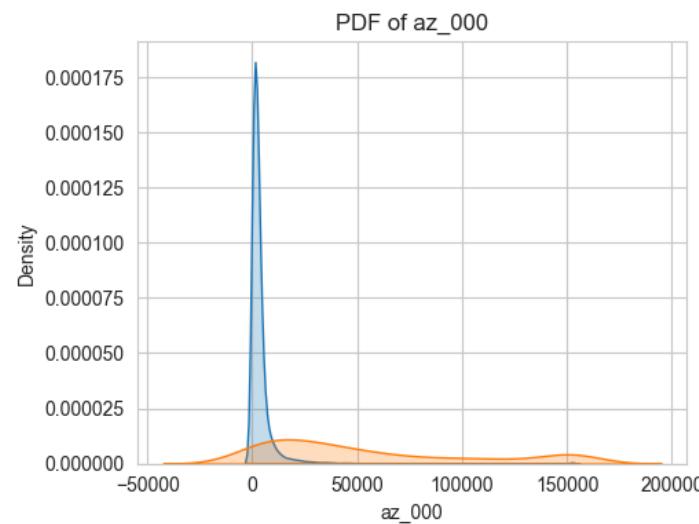
Feature: az\_000

Class 0 - Mean: 4113.24, Std Dev: 10272.3

Class 1 - Mean: 55453.29, Std Dev: 50907.67

2025-02-09 07:33:51,535 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:33:51,590 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



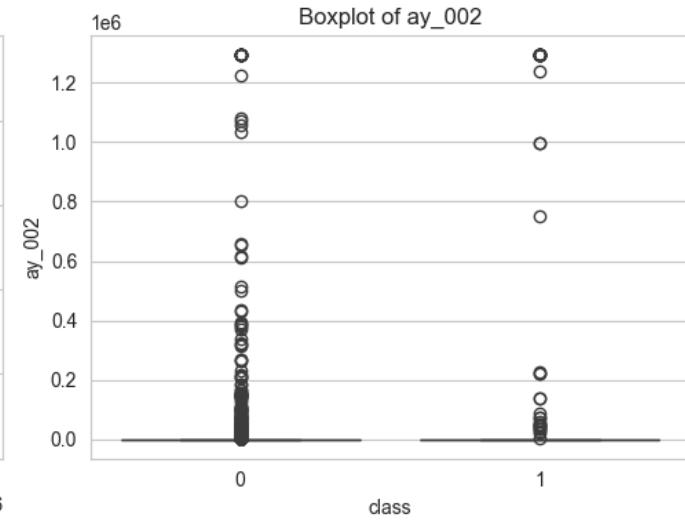
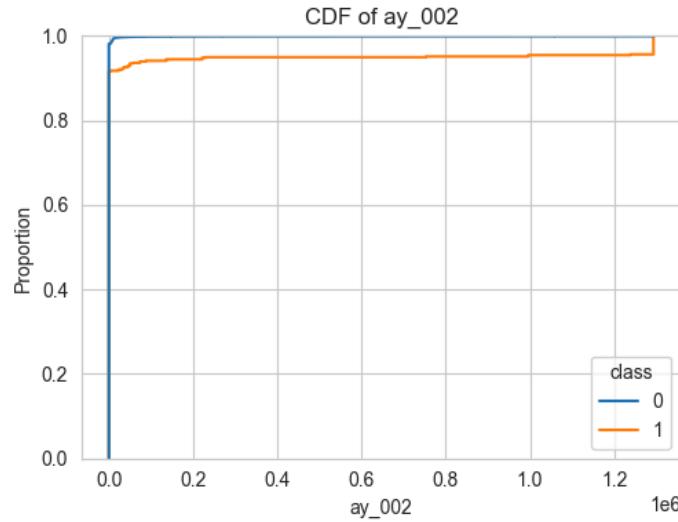
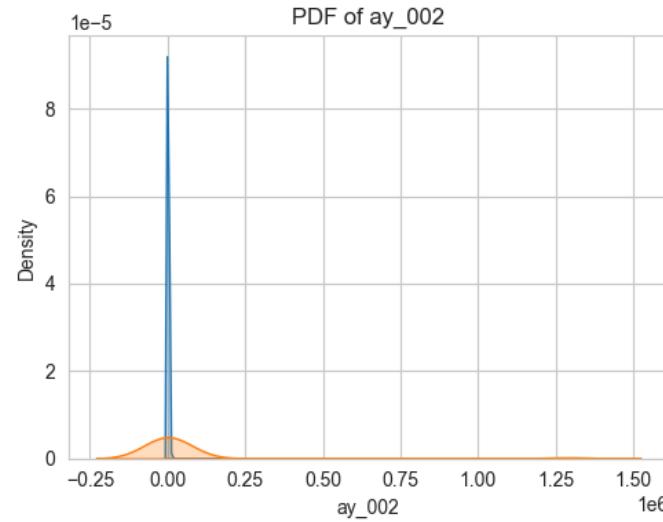
Feature: ay\_002

Class 0 - Mean: 979.71, Std Dev: 26902.12

Class 1 - Mean: 65842.22, Std Dev: 275662.32

2025-02-09 07:33:52,893 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:33:52,945 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



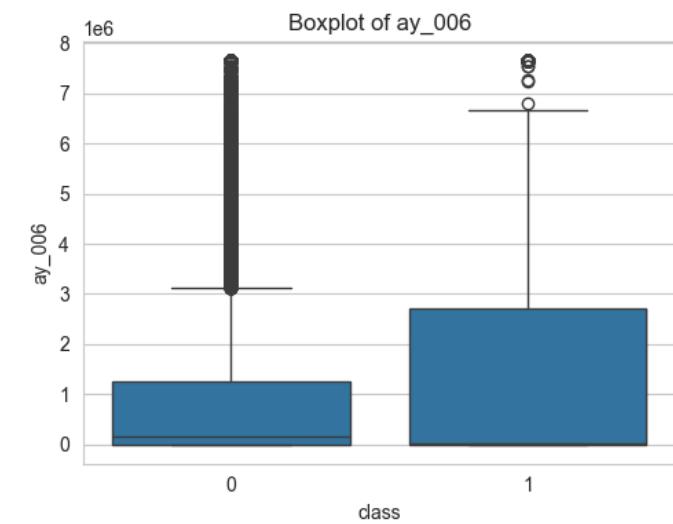
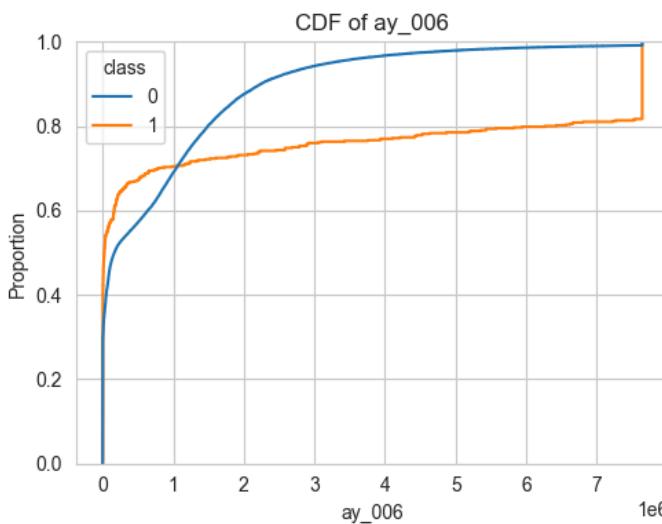
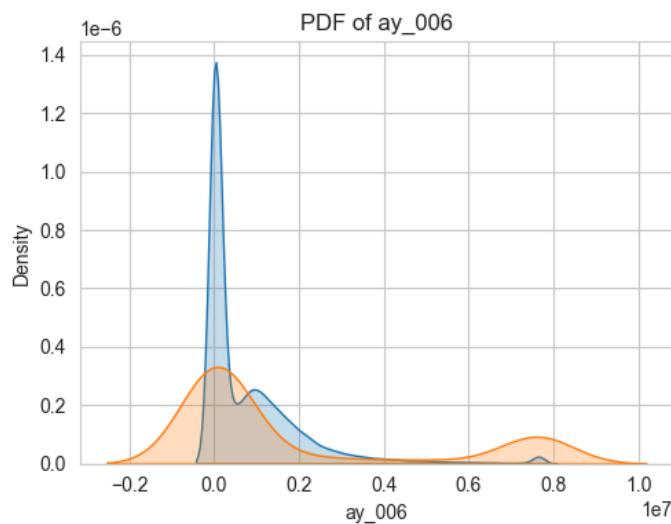
Feature: ay\_006

Class 0 - Mean: 827244.52, Std Dev: 1283793.02

Class 1 - Mean: 1860523.81, Std Dev: 3035459.06

2025-02-09 07:33:54,758 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:33:54,813 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



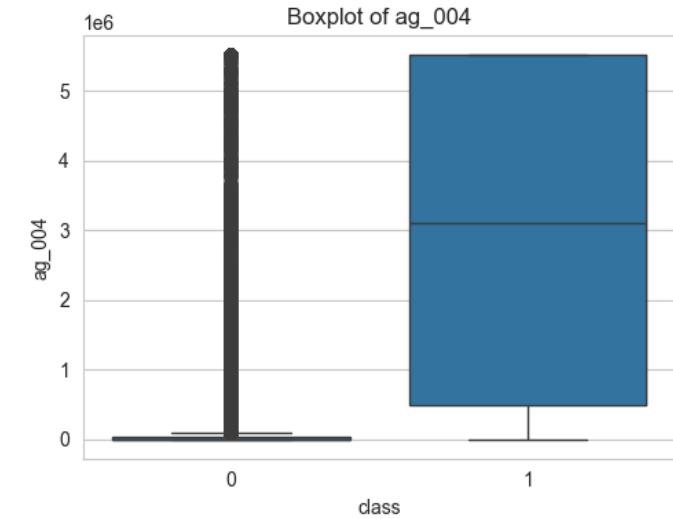
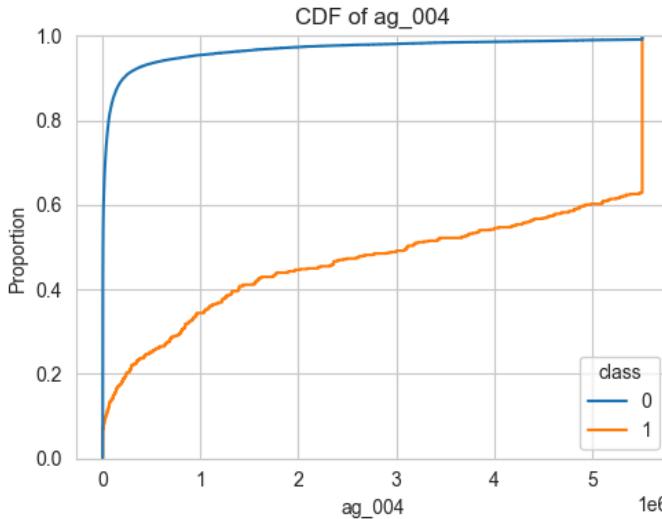
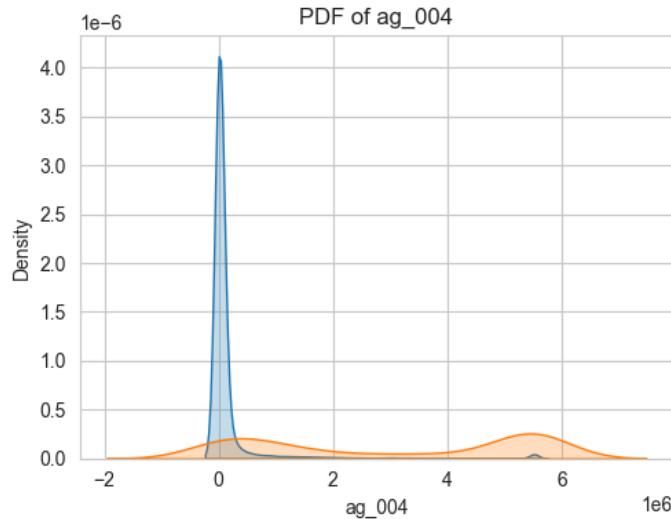
Feature: ag\_004

Class 0 - Mean: 177864.92, Std Dev: 711312.84

Class 1 - Mean: 2990868.85, Std Dev: 2342599.3

2025-02-09 07:33:56,425 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:33:56,481 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



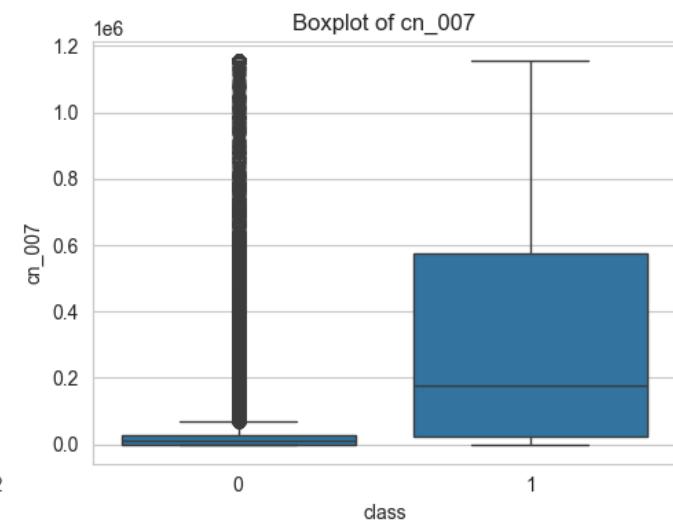
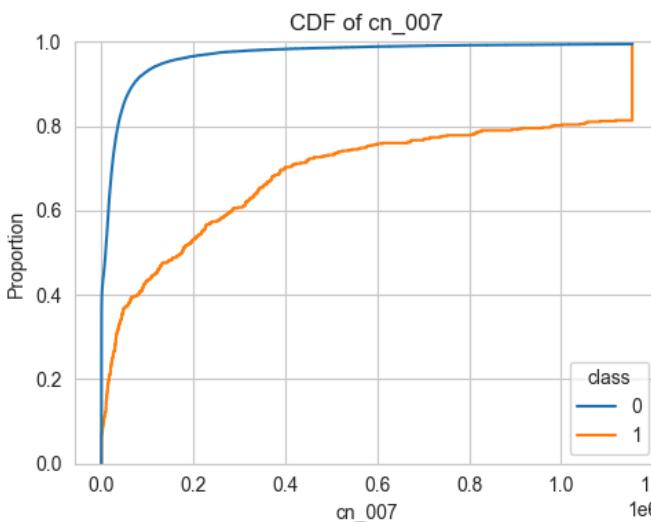
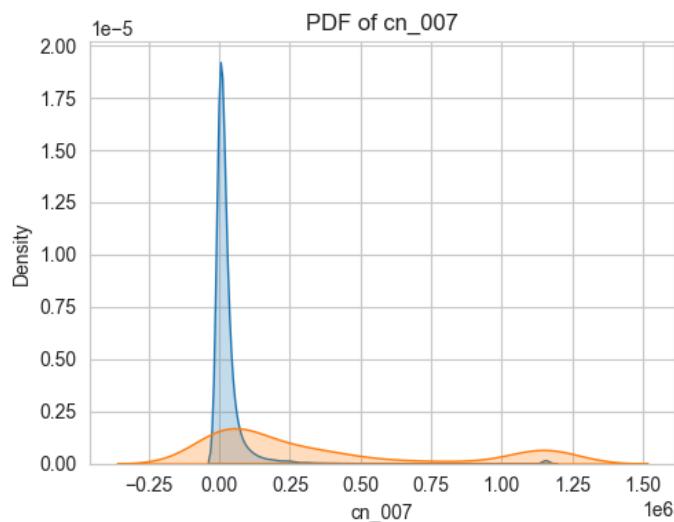
Feature: cn\_007

Class 0 - Mean: 37302.56, Std Dev: 117119.6

Class 1 - Mean: 367807.99, Std Dev: 433203.27

2025-02-09 07:33:58,014 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:33:58,077 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



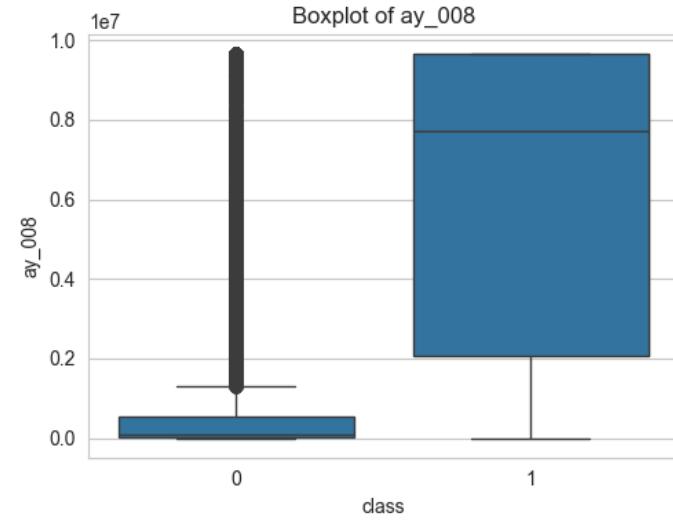
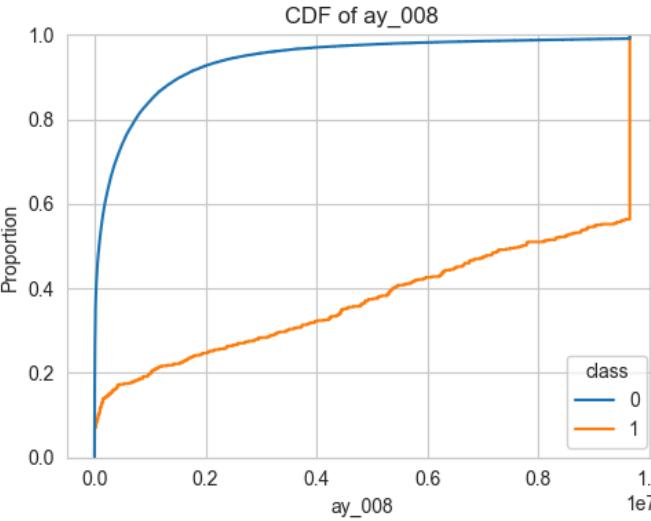
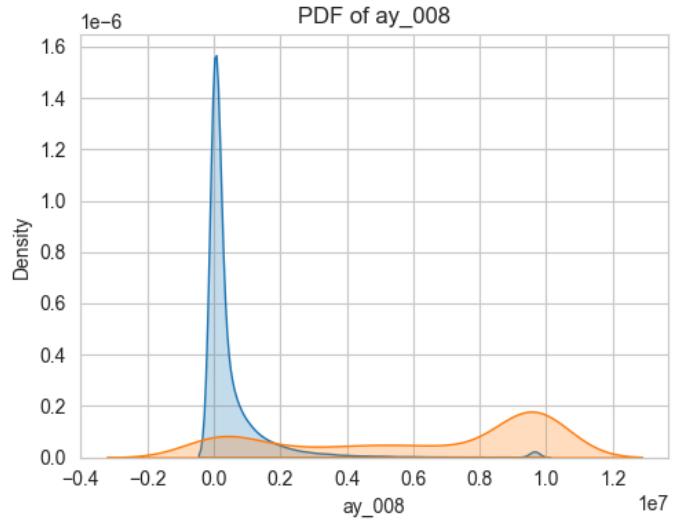
Feature: ay\_008

Class 0 - Mean: 586480.31, Std Dev: 1367972.68

Class 1 - Mean: 6120848.48, Std Dev: 3859848.57

2025-02-09 07:34:00,090 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:00,166 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



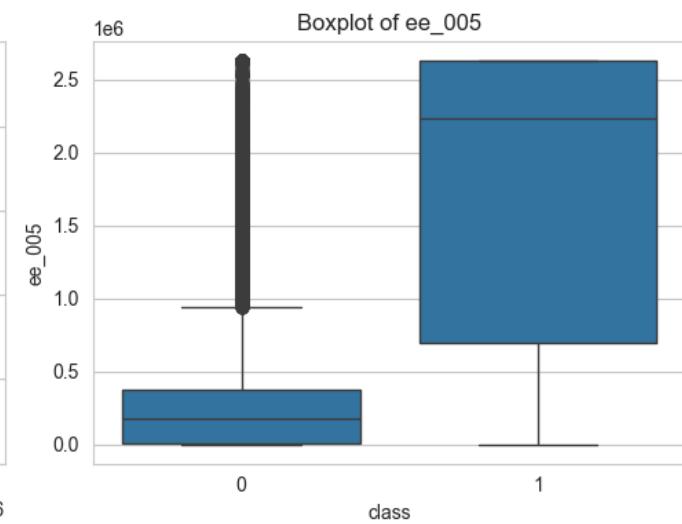
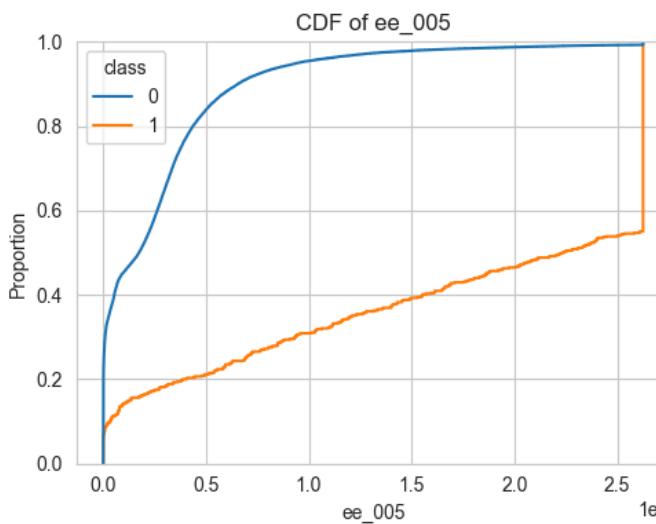
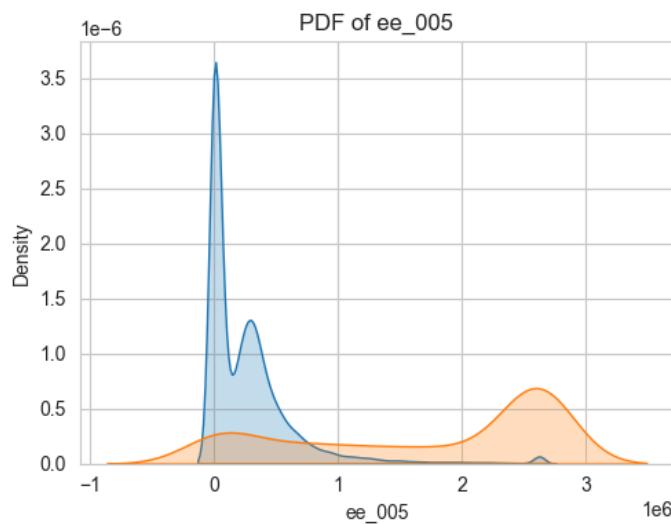
Feature: ee\_005

Class 0 - Mean: 274215.65, Std Dev: 393568.27

Class 1 - Mean: 1699401.77, Std Dev: 1035584.43

2025-02-09 07:34:02,913 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:02,974 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



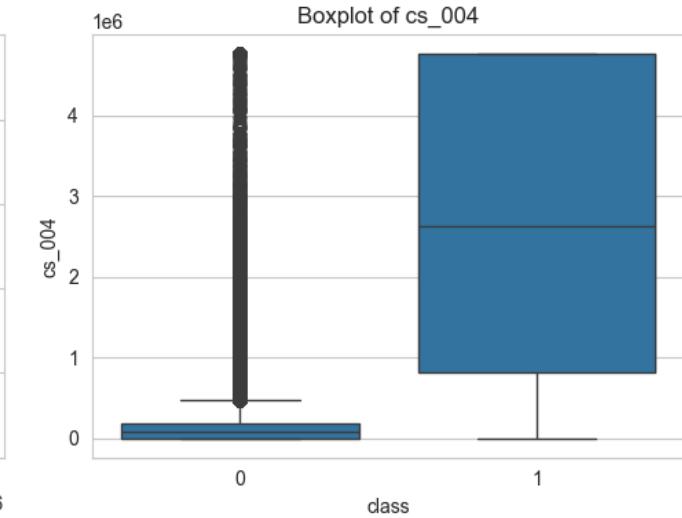
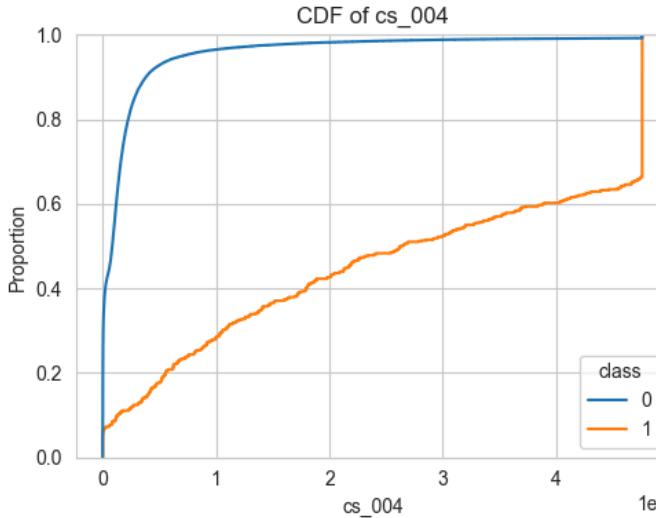
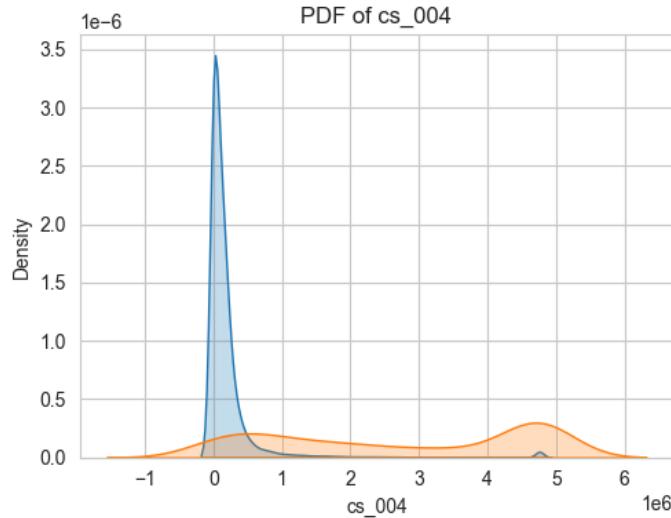
Feature: cs\_004

Class 0 - Mean: 205035.99, Std Dev: 540146.55

Class 1 - Mean: 2687322.51, Std Dev: 1865010.05

2025-02-09 07:34:04,526 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:04,593 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



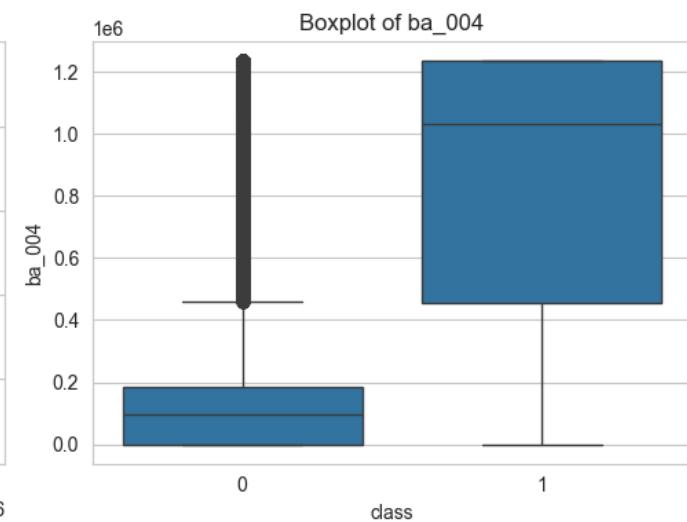
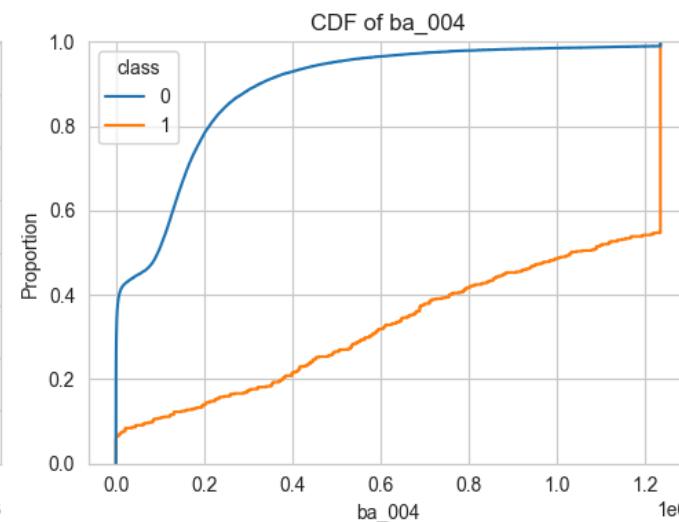
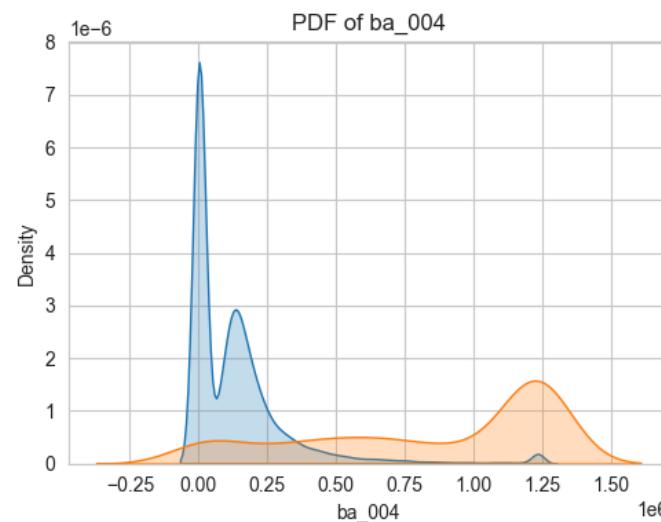
Feature: ba\_004

Class 0 - Mean: 135852.42, Std Dev: 200144.14

Class 1 - Mean: 838808.52, Std Dev: 446858.1

2025-02-09 07:34:06,438 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:06,496 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

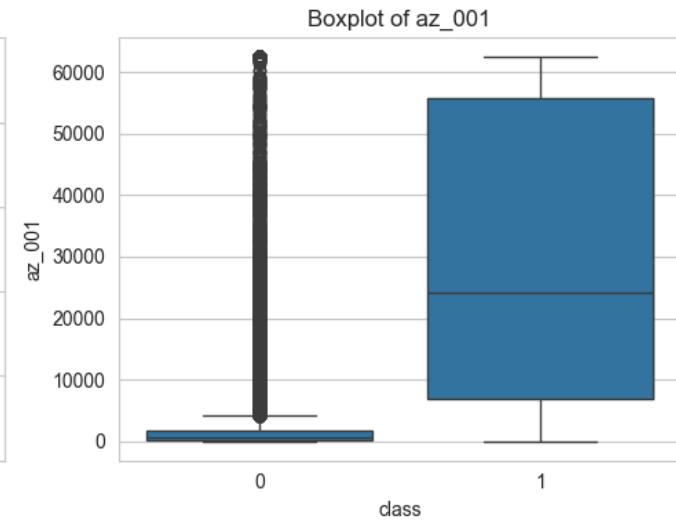
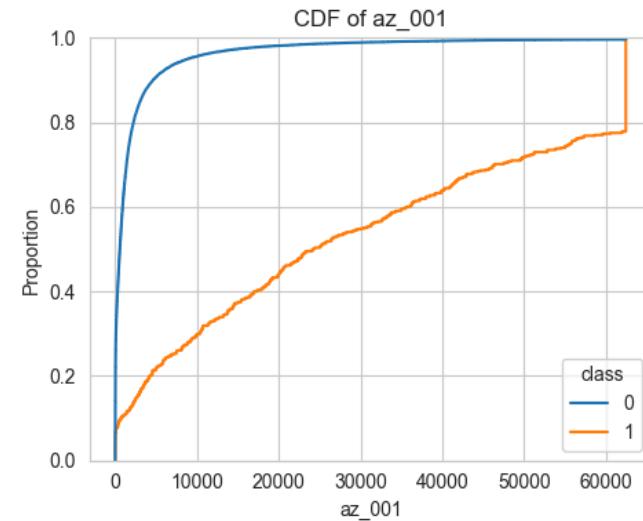
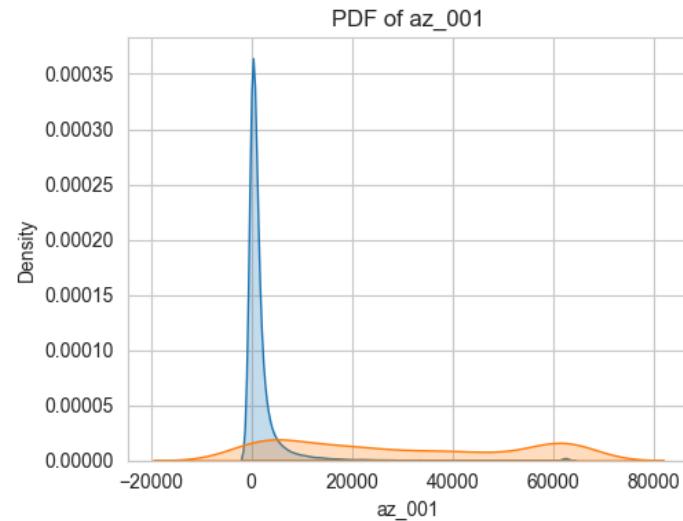


Feature: az\_001

Class 0 - Mean: 2171.77, Std Dev: 5947.25  
Class 1 - Mean: 29546.83, Std Dev: 23398.07

2025-02-09 07:34:08,066 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:08,125 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

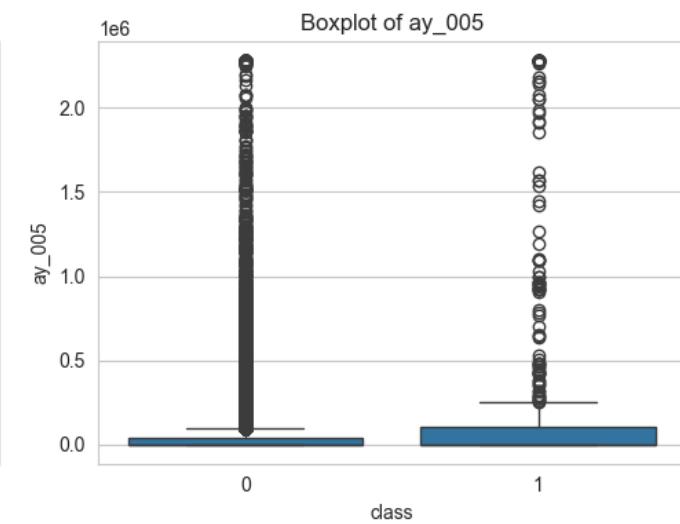
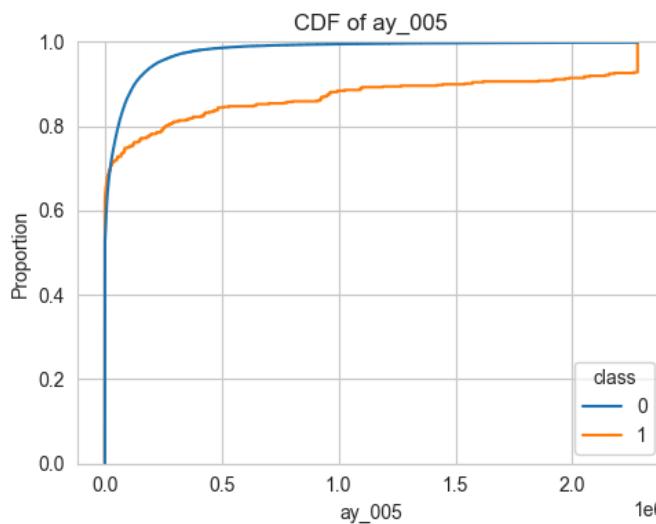
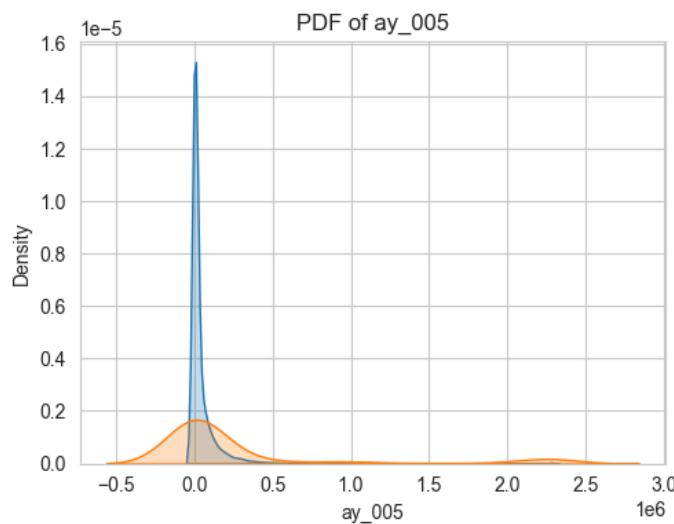


Feature: ay\_005

Class 0 - Mean: 48905.89, Std Dev: 154233.71  
Class 1 - Mean: 301377.55, Std Dev: 673758.72

2025-02-09 07:34:09,637 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:09,695 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



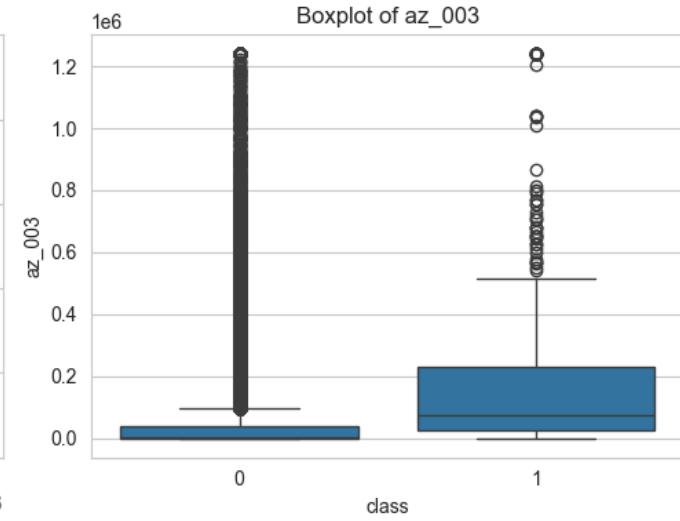
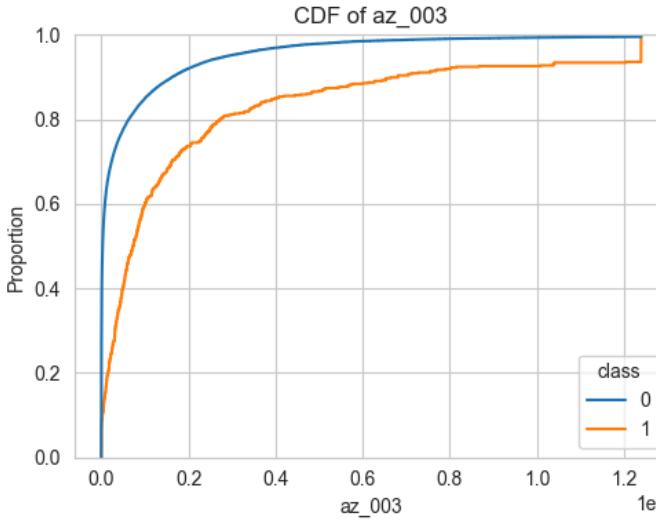
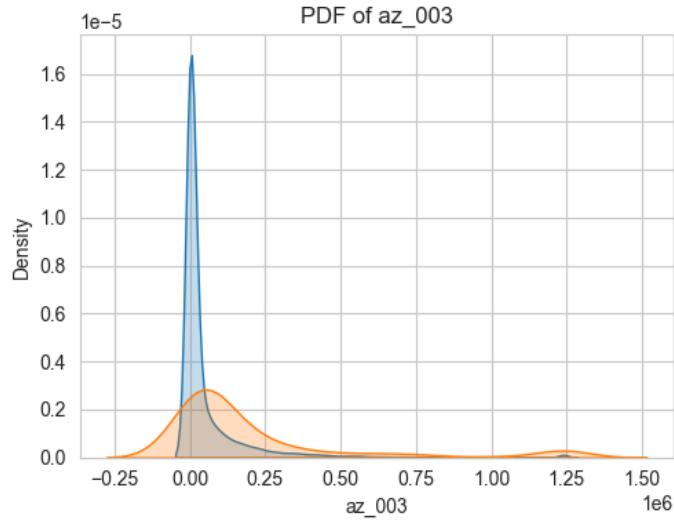
Feature: az\_003

Class 0 - Mean: 55657.64, Std Dev: 143415.85

Class 1 - Mean: 210253.21, Std Dev: 330013.06

2025-02-09 07:34:11,902 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:11,970 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



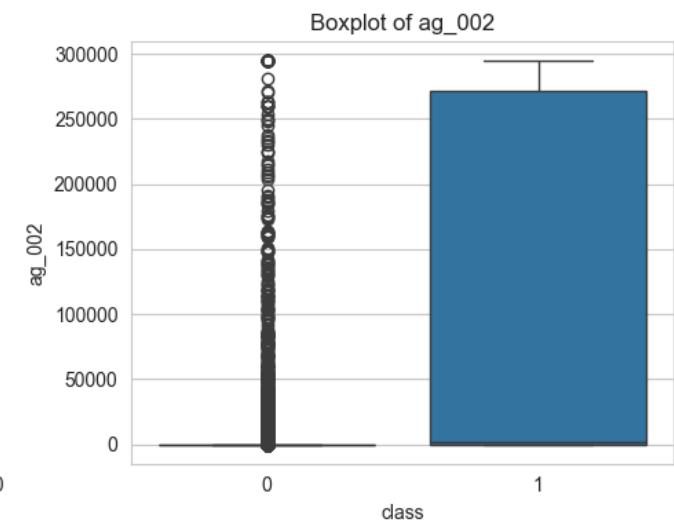
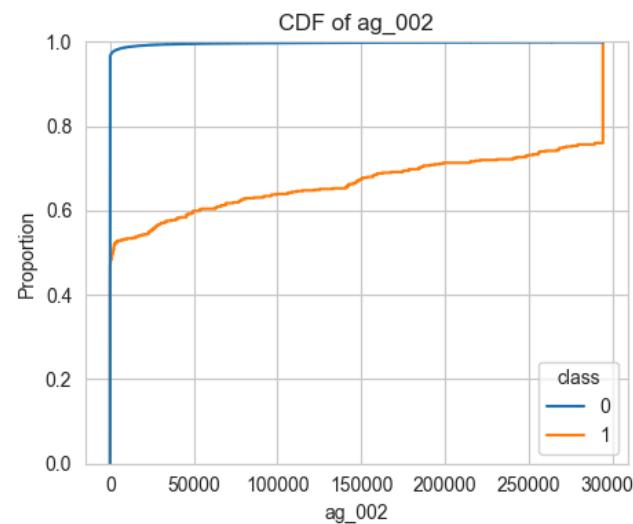
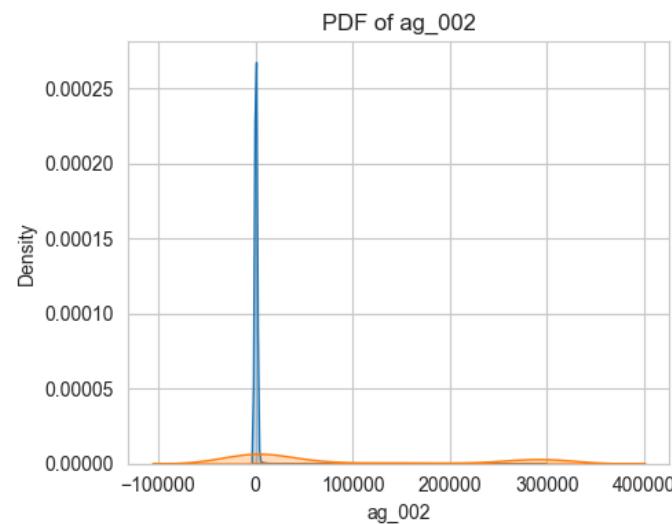
Feature: ag\_002

Class 0 - Mean: 921.23, Std Dev: 11790.91

Class 1 - Mean: 99295.65, Std Dev: 127166.33

2025-02-09 07:34:13,300 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:13,355 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

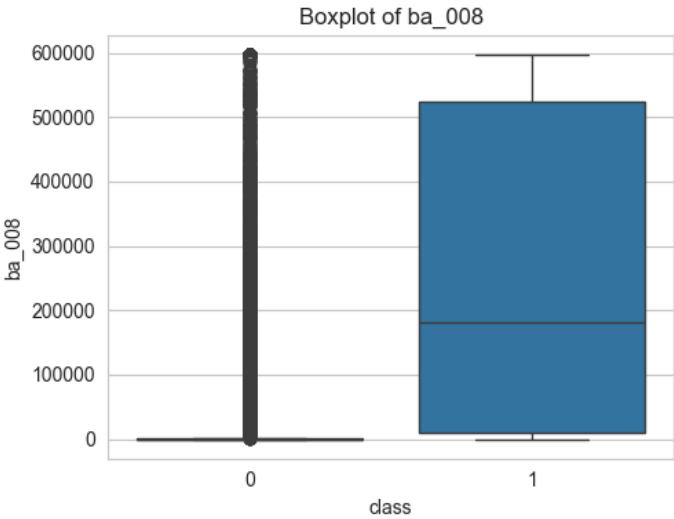
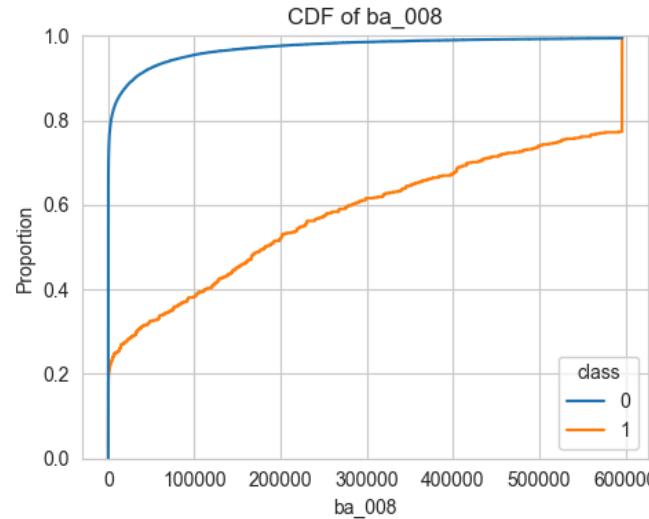
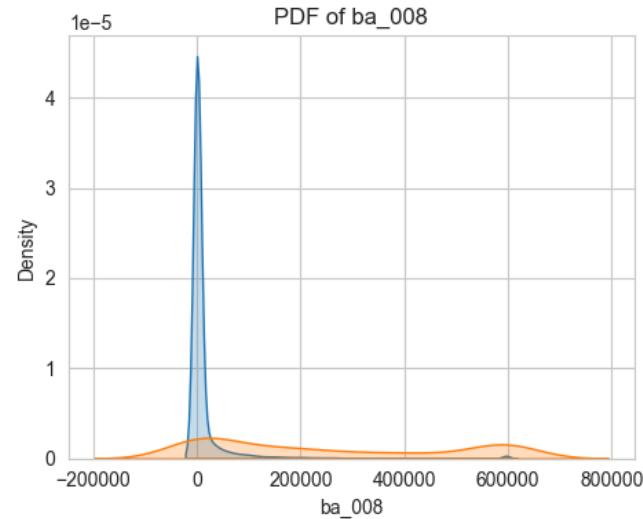


Feature: ba\_008

Class 0 - Mean: 16728.03, Std Dev: 66532.6  
Class 1 - Mean: 253312.74, Std Dev: 236171.96

2025-02-09 07:34:14,947 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:15,008 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

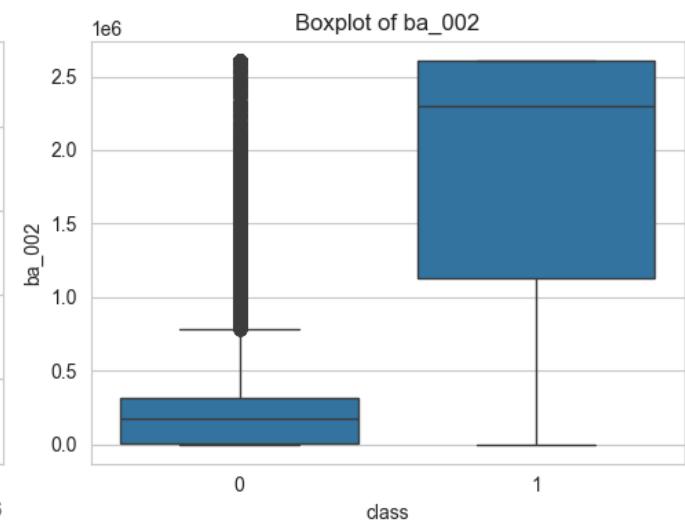
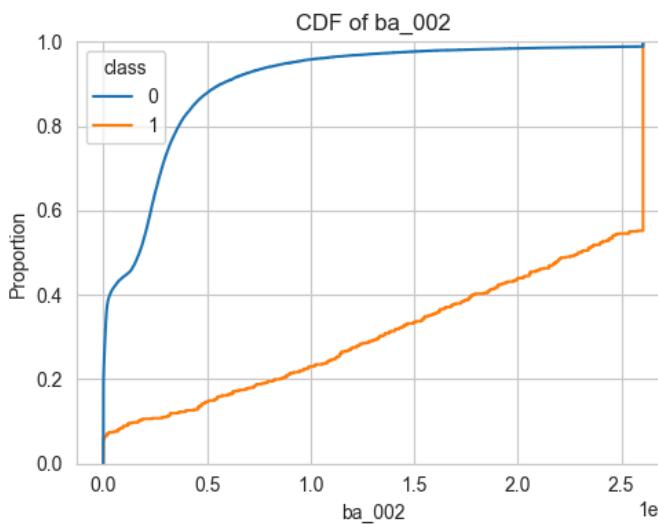
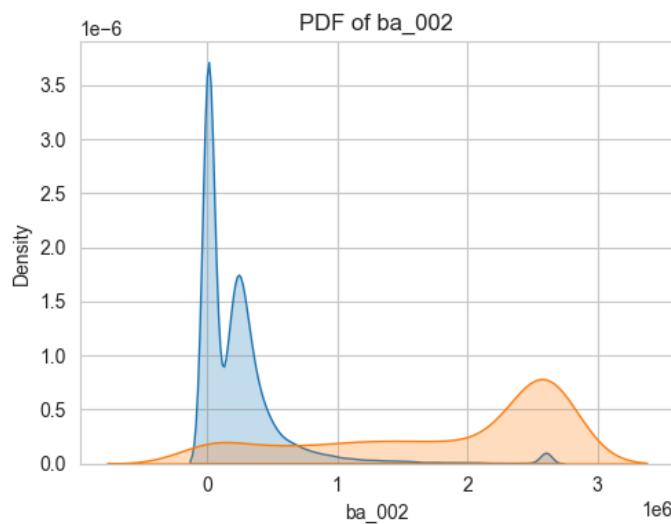


Feature: ba\_002

Class 0 - Mean: 252747.15, Std Dev: 399245.77  
Class 1 - Mean: 1821248.6, Std Dev: 925573.73

2025-02-09 07:34:16,871 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:16,937 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



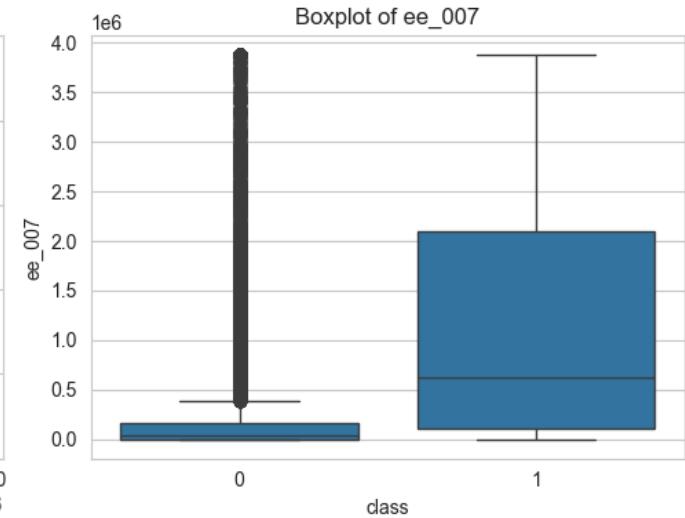
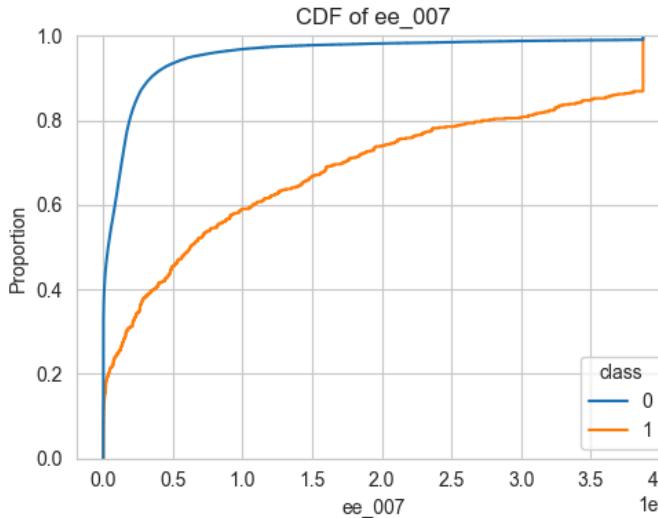
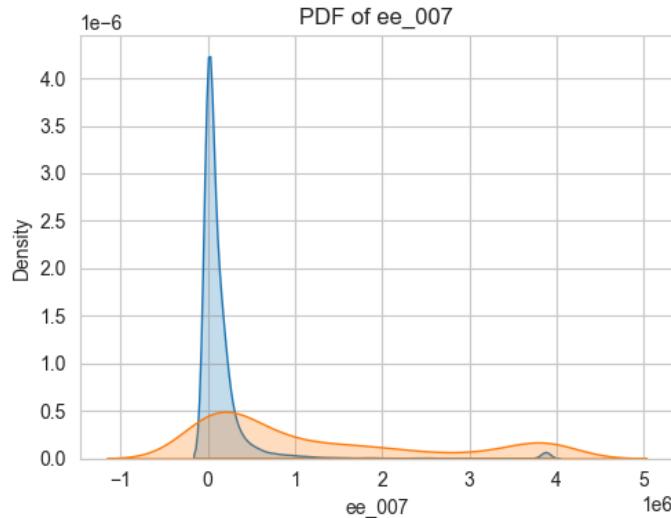
Feature: ee\_007

Class 0 - Mean: 174434.37, Std Dev: 484325.76

Class 1 - Mean: 1260042.51, Std Dev: 1384523.08

2025-02-09 07:34:18,806 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:18,873 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



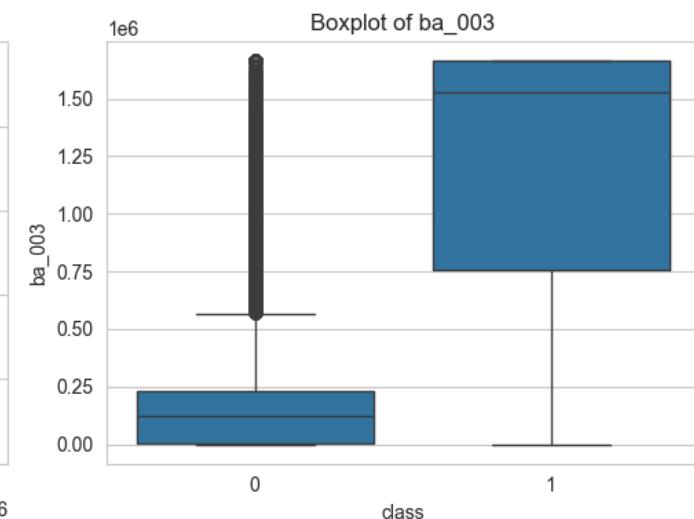
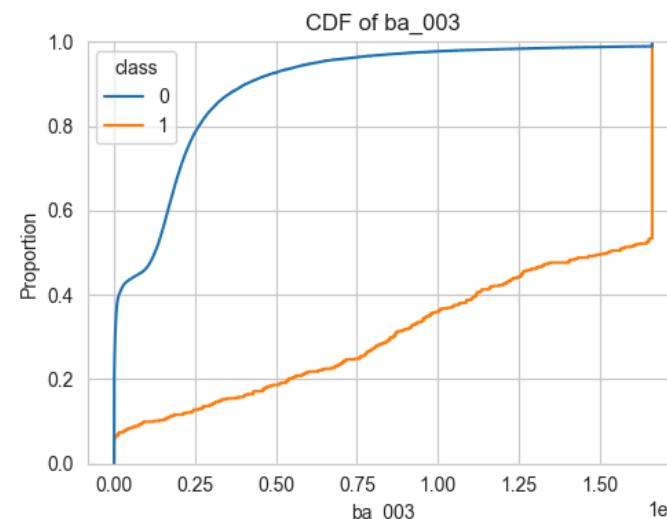
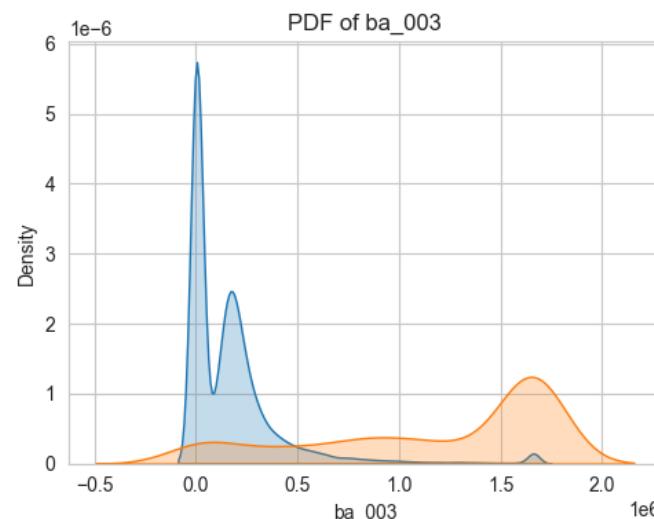
Feature: ba\_003

Class 0 - Mean: 174931.57, Std Dev: 262686.09

Class 1 - Mean: 1165446.46, Std Dev: 591149.8

2025-02-09 07:34:20,759 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:20,826 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



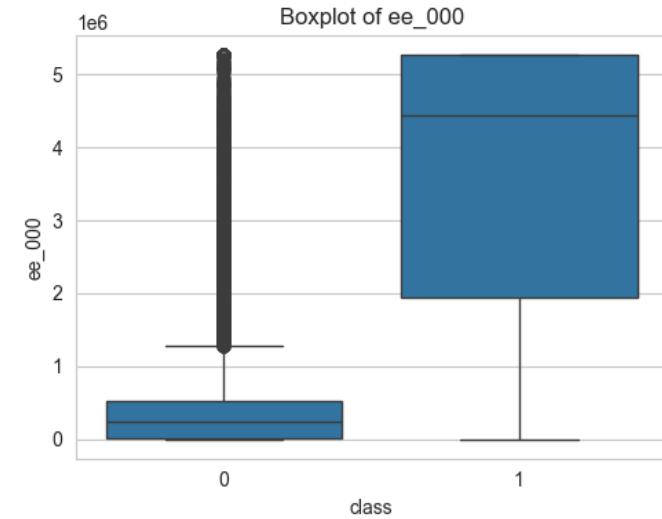
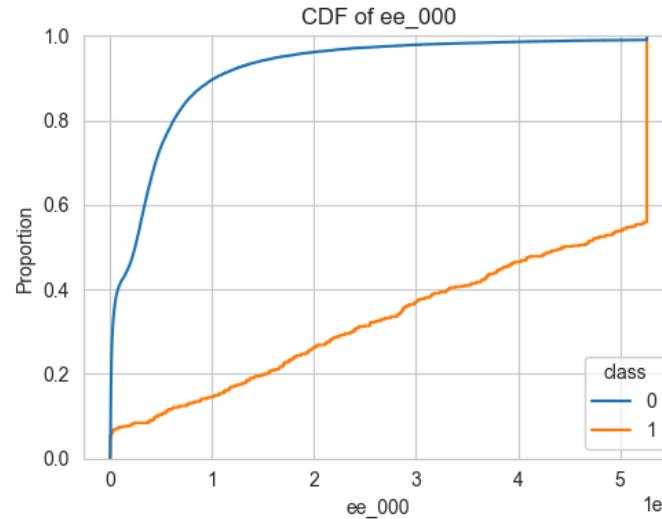
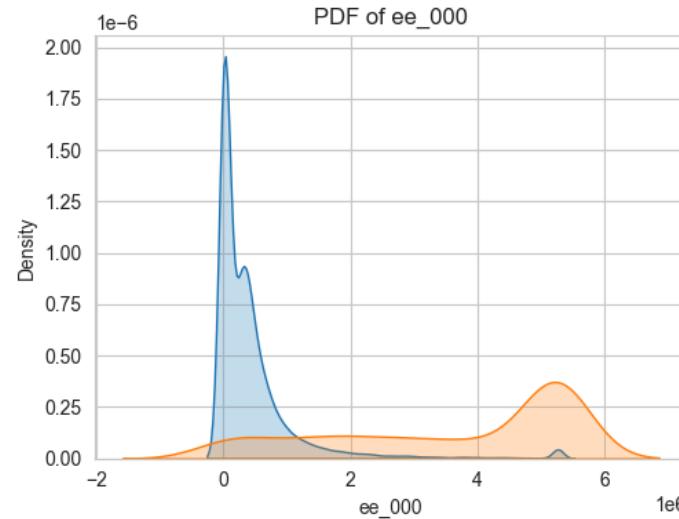
Feature: ee\_000

Class 0 - Mean: 444468.81, Std Dev: 764033.85

Class 1 - Mean: 3585103.51, Std Dev: 1900469.21

2025-02-09 07:34:22,812 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:22,894 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



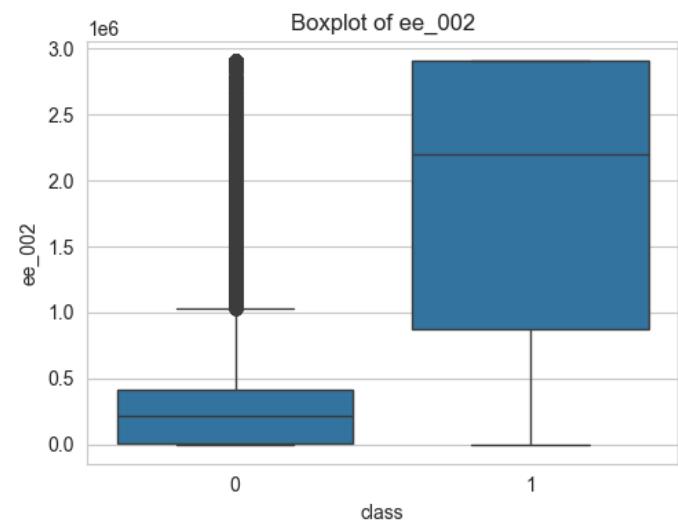
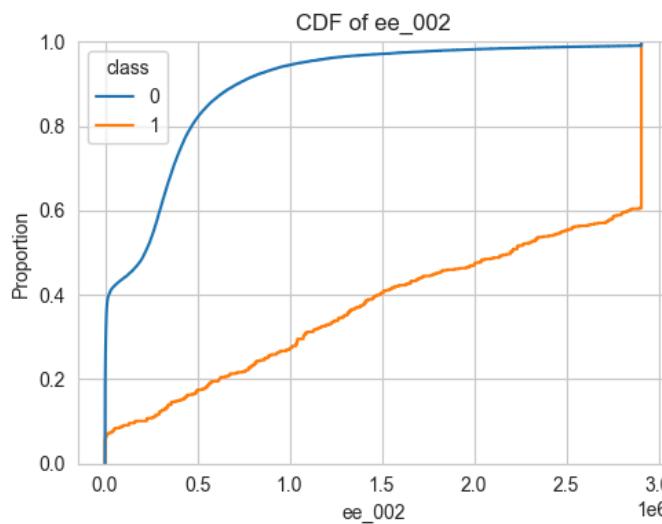
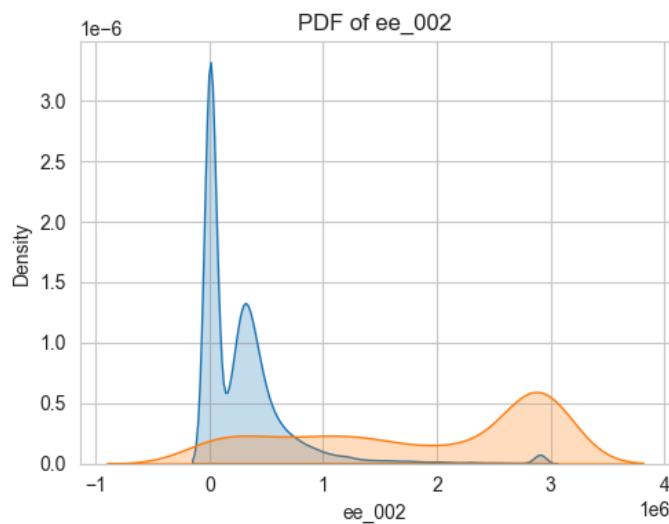
Feature: ee\_002

Class 0 - Mean: 306233.04, Std Dev: 453072.71

Class 1 - Mean: 1854181.83, Std Dev: 1082530.44

2025-02-09 07:34:25,176 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:25,237 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



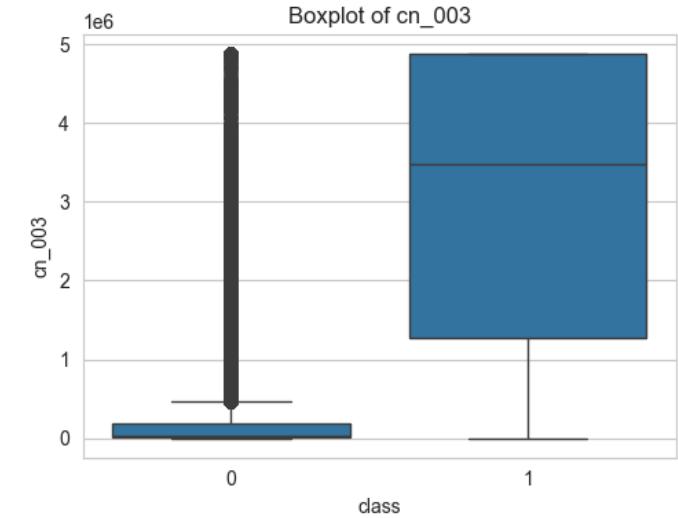
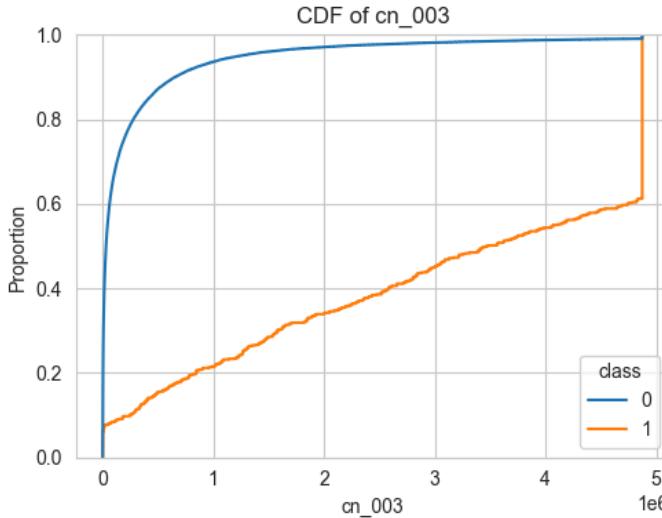
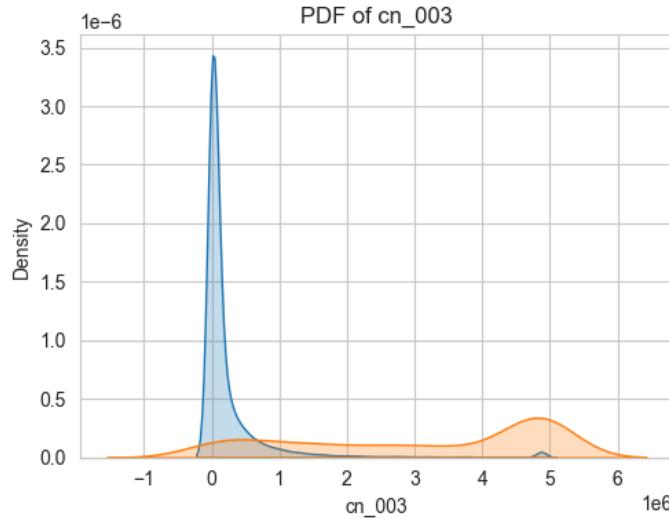
Feature: cn\_003

Class 0 - Mean: 260751.18, Std Dev: 676823.03

Class 1 - Mean: 3048206.52, Std Dev: 1856041.61

2025-02-09 07:34:27,010 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:27,075 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



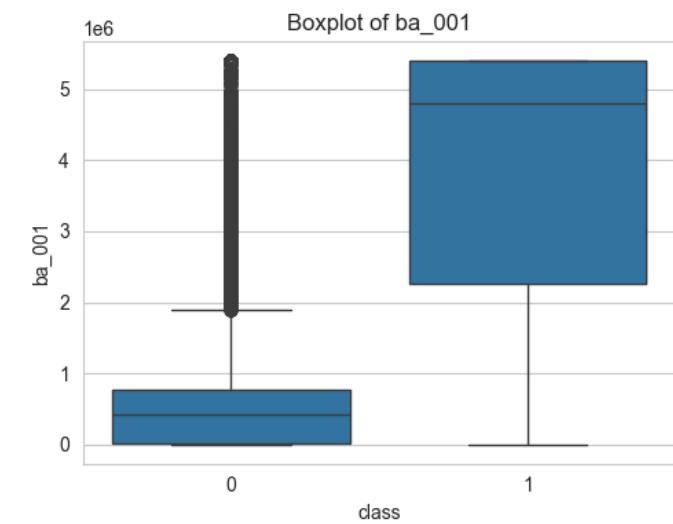
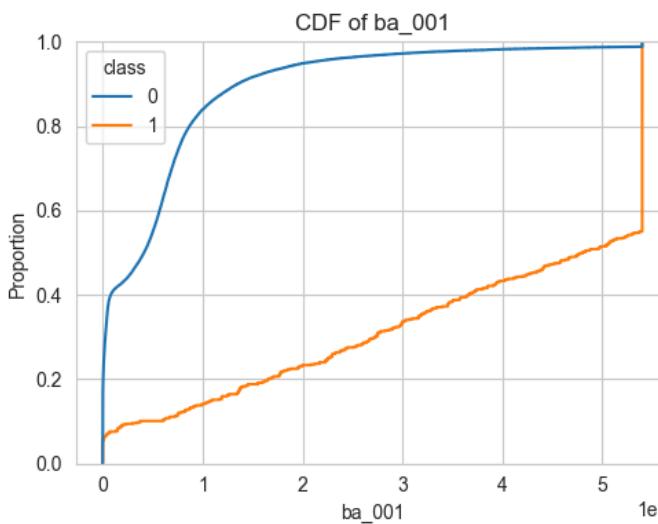
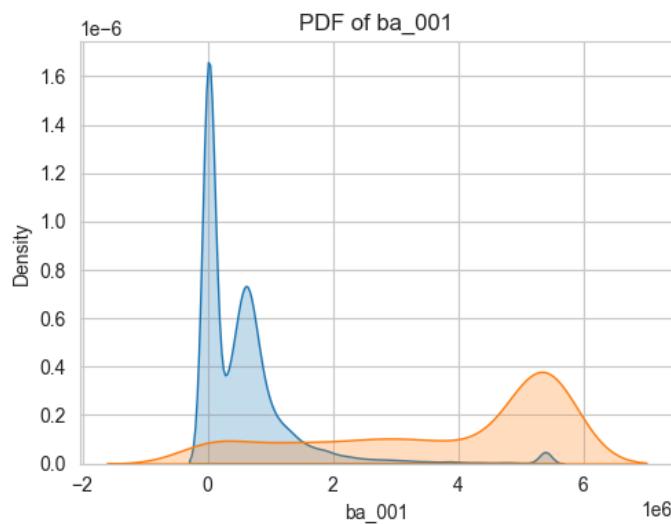
Feature: ba\_001

Class 0 - Mean: 589463.33, Std Dev: 868279.41

Class 1 - Mean: 3763311.7, Std Dev: 1926866.25

2025-02-09 07:34:29,463 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:29,608 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



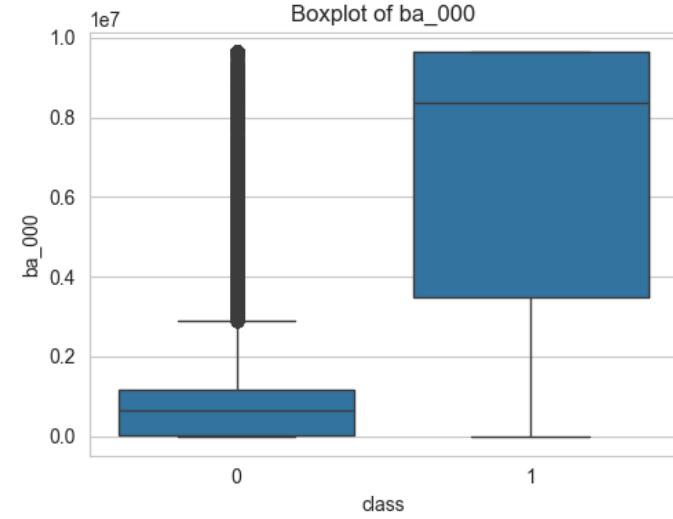
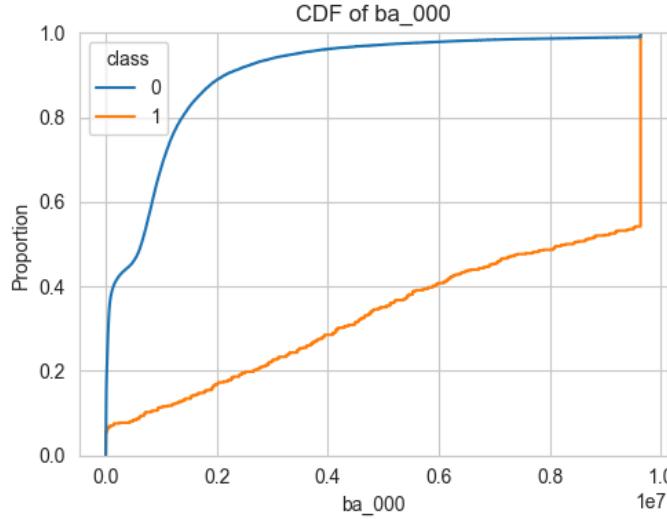
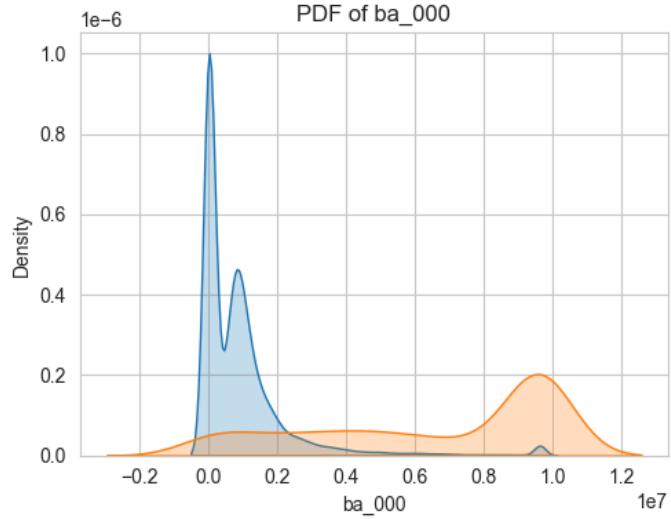
Feature: ba\_000

Class 0 - Mean: 946087.81, Std Dev: 1472533.98

Class 1 - Mean: 6518263.98, Std Dev: 3520398.31

2025-02-09 07:34:31,614 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:34:31,675 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



#### Univariate Analysis of System Parameters for Failure Detection

The analysis encompasses multiple system parameters across air pressure, component, electronic, and air-ground systems, revealing distinct patterns between normal operations and failure states. The findings indicate strong potential for early failure detection and predictive maintenance implementation.

#### Statistical Analysis and Implications:

- **Air System Parameters ( az\_000 , az\_001 , az\_005 )**: The air pressure system parameters demonstrate remarkable distinction between normal and failure states. Parameter az\_000 maintains a tight normal distribution during standard operation (mean: 4,113.24, std dev: 50,907.67), with failure states showing significantly higher values. The az\_005 parameter exhibits the most dramatic separation, with failure cases averaging 10.59M compared to normal operations at 1.33M. This substantial difference provides a reliable basis for early failure detection protocols.
- **Component System Parameters ( cs\_002 , cs\_004 )**: These parameters show the strongest separation between operational states. The cs\_002 parameter's failure state median (~1.5M) significantly exceeds normal operations (~50K), while cs\_004 demonstrates extreme separation with failure conditions averaging 2.69M compared to normal operations at 205K. This clear delineation makes these parameters particularly valuable for automated monitoring systems.
- **Electronic System Parameters ( ee\_002 , ee\_005 , ee\_007 )**: The electronic system parameters display distinct multimodal distributions, suggesting different operational and failure modes. Parameter ee\_005 shows a notable bimodal distribution in failure cases, indicating two distinct failure patterns. The mean values in failure cases (1.70M) with substantial standard deviations suggest these parameters effectively indicate failure severity and progression.
- **Component Network Parameters ( cn\_003 , cn\_007 )**: Parameter cn\_003 shows excellent separation between states, with failure cases averaging 3.05M compared to normal operations at 261K. The parameter cn\_007 exhibits a more gradual transition between states, making it suitable for detecting progressive degradation.

#### Operational Implications and Recommendations:

- The analysis supports implementing a multi-tiered monitoring approach. Primary monitoring should focus on parameters showing the clearest separation ( cs\_002 , az\_005 , cn\_003 ) with secondary validation from parameters showing gradual transitions ( cn\_007 , ee\_007 ). Critical threshold values should be established at the 90th percentile of normal operations for each parameter, with warning thresholds at the 75th percentile.
- The wide variation in failure state standard deviations suggests different failure modes may be identifiable through parameter pattern analysis. This enables not just failure detection but potential failure mode classification, allowing maintenance teams to prepare appropriate responses proactively.

#### System Design Recommendations:

The analysis supports implementing real-time monitoring systems with multi-parameter tracking capabilities. Parameters should be weighted according to their discrimination ability, with highest weights assigned to those showing clearest separation between operational states. The monitoring system should incorporate both absolute threshold violations and trend analysis to capture both sudden failures and gradual degradation patterns.

In [25]:

```
def correlation_matrix(data: pd.DataFrame, n_select: int = 5, plot: bool=True):
    # Compute pairwise correlations between selected features
    correlation_matrix = data.corr()
    correlation_values = correlation_matrix.loc[:, 'class']

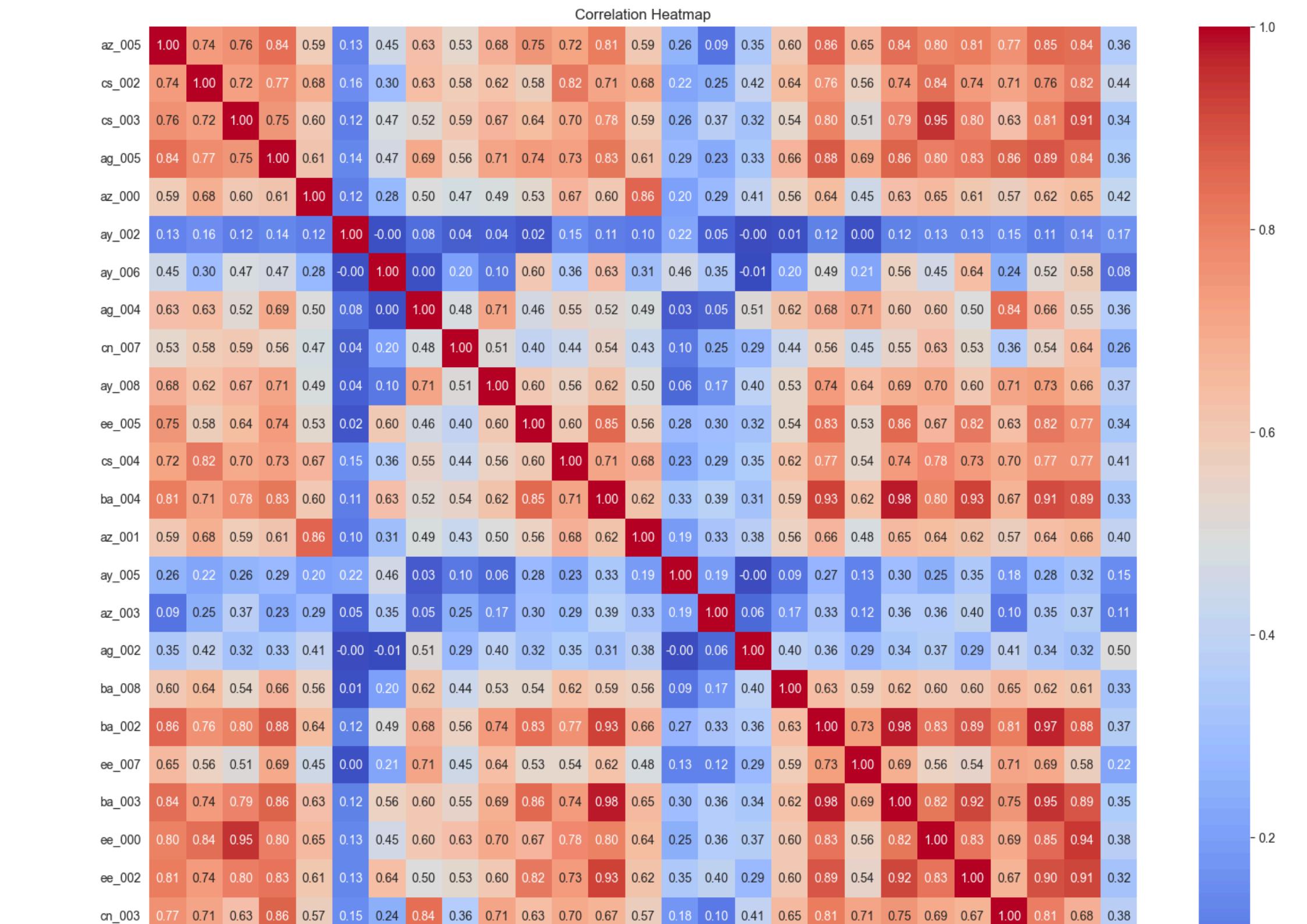
    # Visualize correlations using a heatmap
    if plot:
        plt.figure(figsize=(18, 15))
        sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt='.2f')
        plt.title("Correlation Heatmap")
        plt.show()

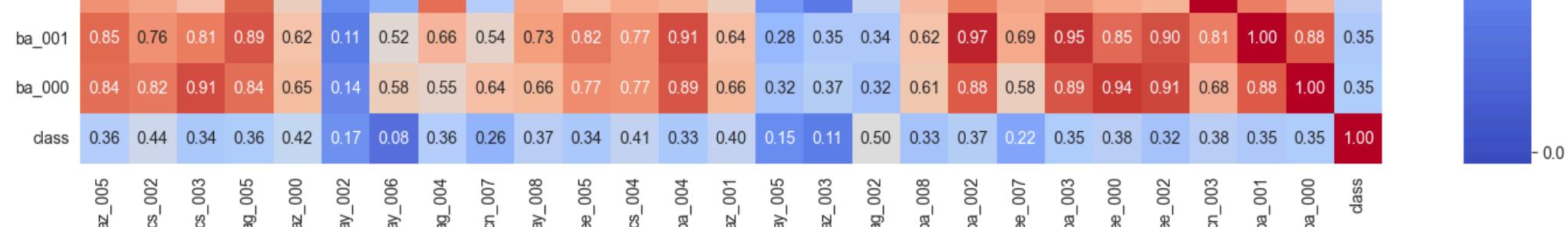
top_correlated_features_with_score = correlation_values.sort_values().iloc[:n_select]
logger.info(f"Top 5 uncorrelated features:\n{top_correlated_features_with_score}")
```

```
top_n_correlated_features = top_correlated_features_with_score.index.tolist()
logger.info(f"The most uncorrelated feature is: {top_n_correlated_features[0]}")

return top_n_correlated_features

selected_features_df = train_data_cleaned[selected_features+[ 'class' ]]
top_correlated_features = correlation_matrix(data=selected_features_df)
```





2025-02-09 07:34:35,157 - INFO - Top 5 uncorrelated features:

ay\_006 0.080430

az\_003 0.107604

ay\_005 0.152126

ay\_002 0.168143

ee\_007 0.216570

Name: class, dtype: float64

2025-02-09 07:34:35,158 - INFO - The most uncorrelated feature is: ay\_006

2025-02-09 07:34:35,158 - INFO - The most uncorrelated feature is: ay\_006

## CORRELATION ANALYSIS:

- Strong Correlation Groups Analysis:** The correlation heatmap reveals several distinct groups of highly correlated features. The `ba_00x` series (`ba_001`, `ba_002`, `ba_003`, `ba_004`) shows particularly strong correlations (>0.90) with each other, suggesting these measurements capture related aspects of the system behavior. This high correlation implies that these features might be monitoring different aspects of the same underlying mechanism or physically connected components.
- Moderate Correlation Patterns:** The `cs_00x` series shows moderate to strong correlations (0.60-0.85) with multiple feature groups, particularly with the `ba_00x` series. This suggests these features might serve as bridge indicators that capture broader system behavior. The `ee_00x` features similarly show moderate correlations across multiple groups, indicating their potential value as secondary confirmation indicators.
- Weak Correlation Insights:** Features like `ay_002` and `az_003` show consistently low correlations (<0.30) with most other features, suggesting they capture unique aspects of system behavior. These weakly correlated features could be particularly valuable for detecting failure modes that might not be captured by the more strongly correlated feature groups.
- Class Correlation Analysis:** The 'class' variable shows moderate correlations (0.30-0.50) with several features, particularly with the `ba_00x` series. This suggests these features have good discriminative power for failure detection, though no single feature shows extremely high correlation with the class label.

In [26]:

```
from typing import List
```

```
def scatter_plot(data:pd.DataFrame, feature: str, percentile: int):

    threshold = np.nanpercentile(data[feature], percentile)

    data = data[data[feature] < threshold]

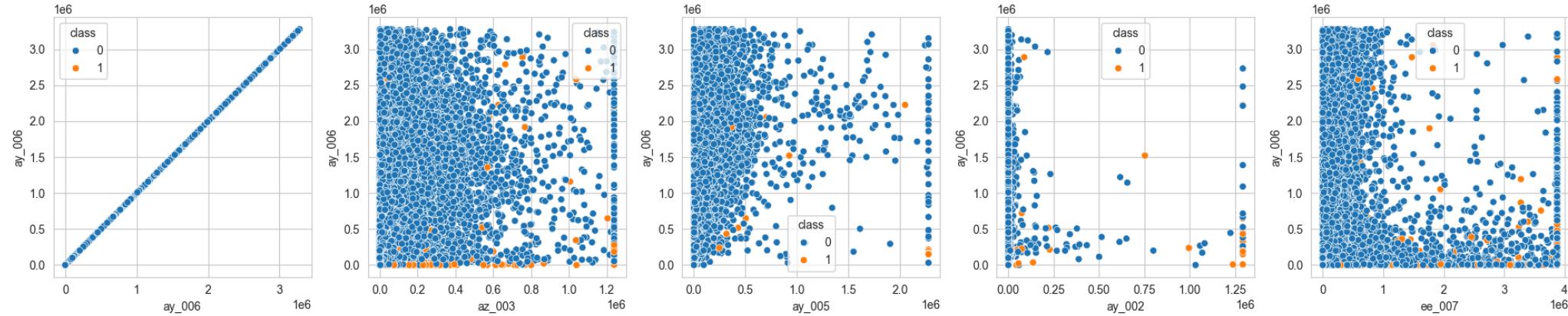
    fig, ax = plt.subplots(1, 5, figsize=(19, 4))
```

```

columns = data.columns.tolist()
for i, col in enumerate(columns): # Limit to the number of created axes
    if col != 'class':
        sns.scatterplot(x=data[col], y=data[feature], hue=data['class'], ax=ax[i])
plt.tight_layout()
plt.show()

top_uncorr_df = selected_features_df[top_correlated_features + ['class']]
scatter_plot(data=top_uncorr_df, feature=top_correlated_features[0], percentile=95)

```



## SCATTER PLOT ANALYSIS:

- `ay_006` vs Feature Relationships:** The scatter plots reveal interesting patterns in relation to `ay_006`. The perfect diagonal line in the first plot indicates a reference or self-correlation plot. Other features show varying degrees of clustering and separation between classes, with some showing clear boundaries between normal operations and failure conditions.
- `az_003` Relationship Patterns:** The scatter plot of `az_003` shows distinct clustering patterns, with failure cases (Class 1) appearing more frequently in certain regions. This suggests potential threshold values could be established for monitoring purposes. The spread of the data indicates a non-linear relationship with other features.
- `ay_005` Distribution Characteristics:** The `ay_005` scatter plots reveal more dispersed patterns, with some clear separation between classes at higher values. This suggests this feature might be more useful when combined with others in a multivariate monitoring approach rather than used in isolation.
- `ay_002` Clustering Behavior:** The `ay_002` relationships show interesting clustering patterns with clear separation in some regions, particularly at higher values. This suggests potential threshold-based monitoring strategies could be effective when using this feature.
- `ee_007` Separation Patterns:** The `ee_007` scatter plots demonstrate good separation between classes in certain regions, with failure cases showing distinct clustering patterns. This feature appears to have good discriminative power, particularly when values exceed certain thresholds.

## OPERATIONAL IMPLICATIONS:

### Monitoring Strategy

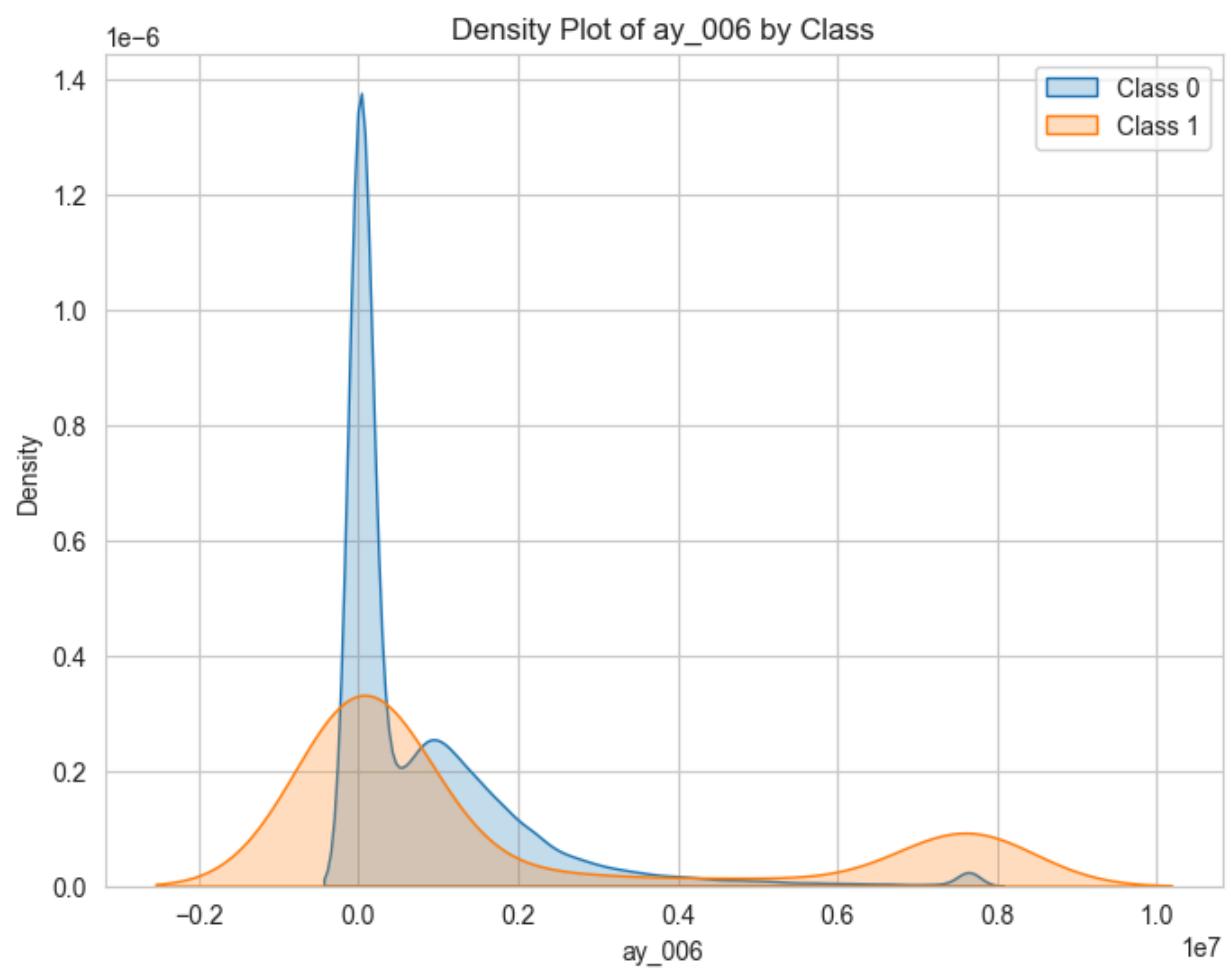
Based on the correlation and scatter plot analyses, a hierarchical monitoring approach is recommended:

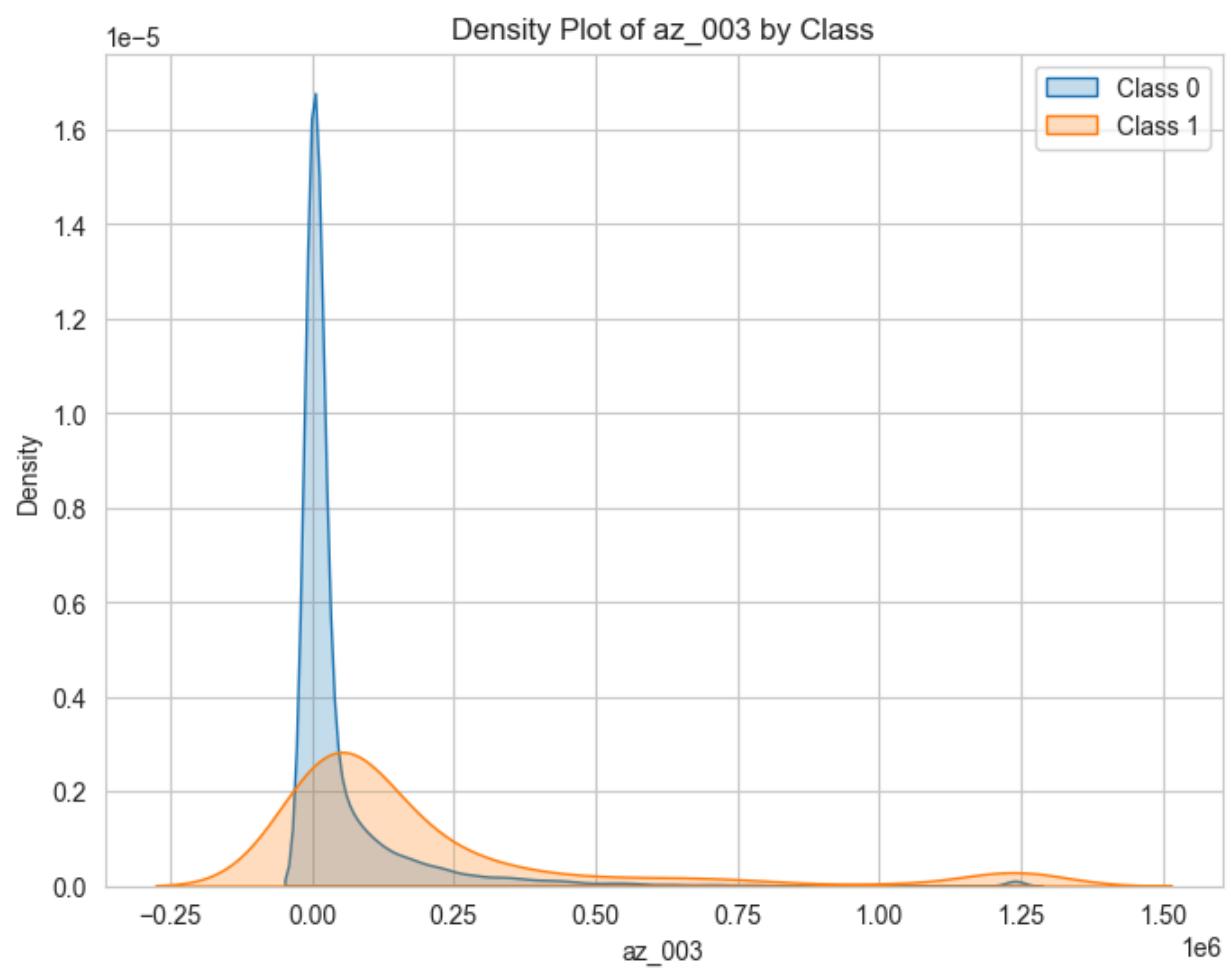
- **Primary Monitoring:** Focus on the highly correlated `ba_00x` series as primary indicators, using their strong relationships to establish reliable baseline patterns.
- **Secondary Verification:** Utilize the moderately correlated `cs_00x` and `ee_00x` series as confirmation indicators, particularly when primary indicators show anomalous behavior.
- **Independent Checks:** Maintain separate monitoring thresholds for the weakly correlated features (`ay_002`, `az_003`) as they may catch failure modes missed by the primary indicators.

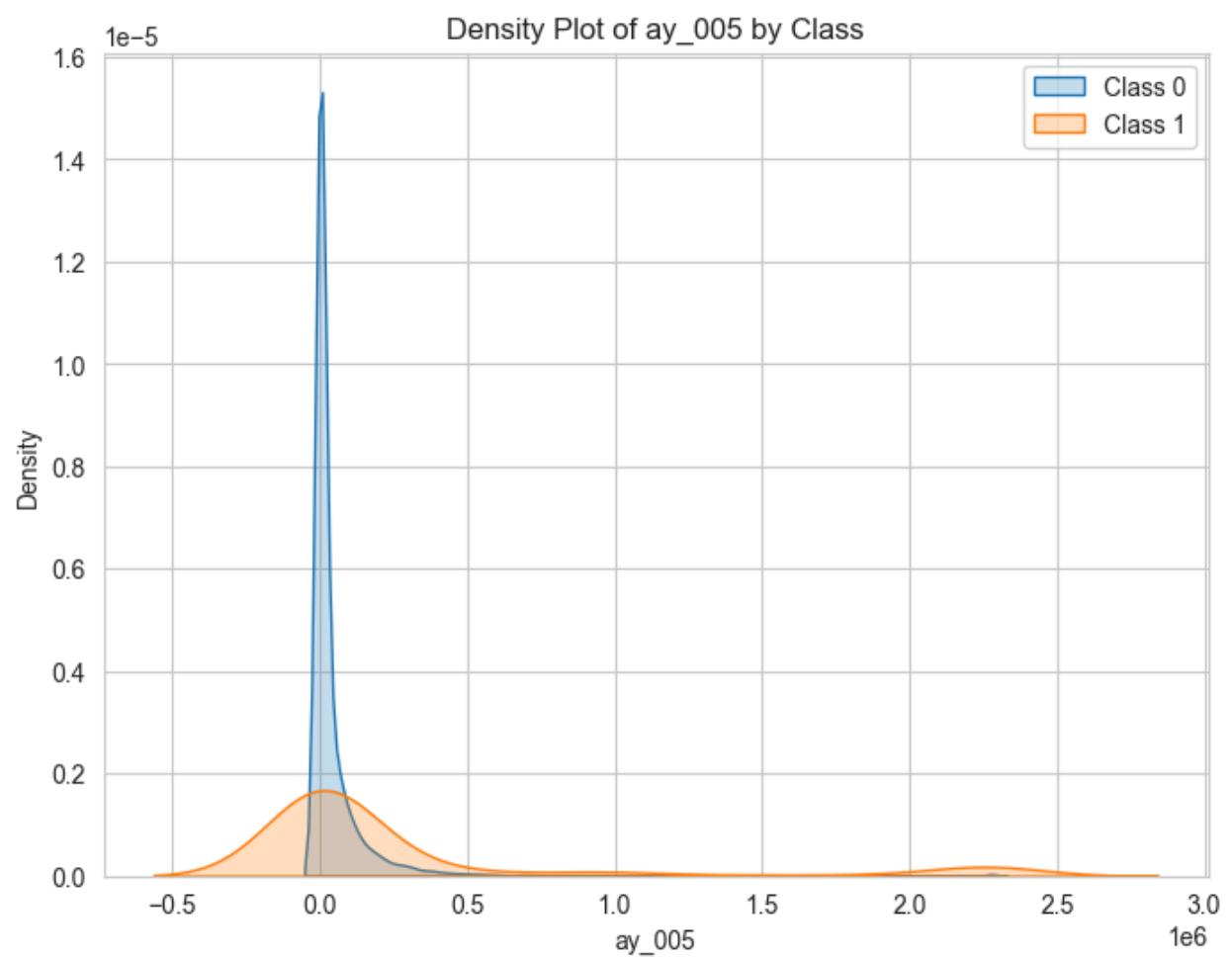
The strong correlations between certain feature groups suggest that a failure in one component is likely to manifest across multiple measurements, providing opportunities for early detection through pattern recognition across the correlated feature sets. The scatter plot patterns indicate that multivariate thresholds, rather than simple univariate limits, might be more effective for early failure detection.

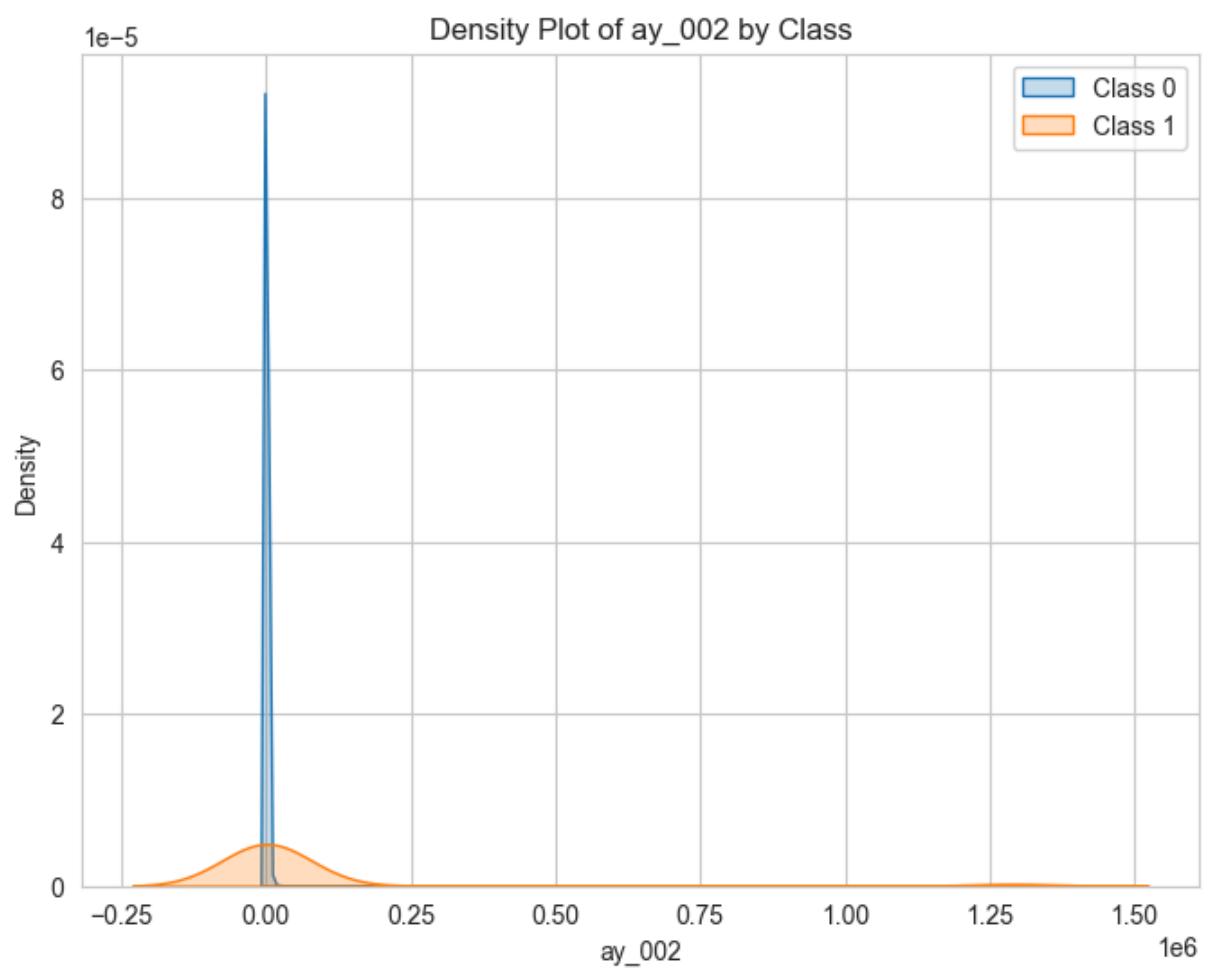
```
In [27]: def visualize_class_distributions(df, feature_list):
    """
    Visualize how selected features differ across classes.
    """
    for feature in feature_list:
        plt.figure(figsize=(8, 6))
        sns.kdeplot(data=df[df['class'] == 0], x=feature, label="Class 0", shade=True)
        sns.kdeplot(data=df[df['class'] == 1], x=feature, label="Class 1", shade=True)
        plt.title(f"Density Plot of {feature} by Class")
        plt.legend()
        plt.show()

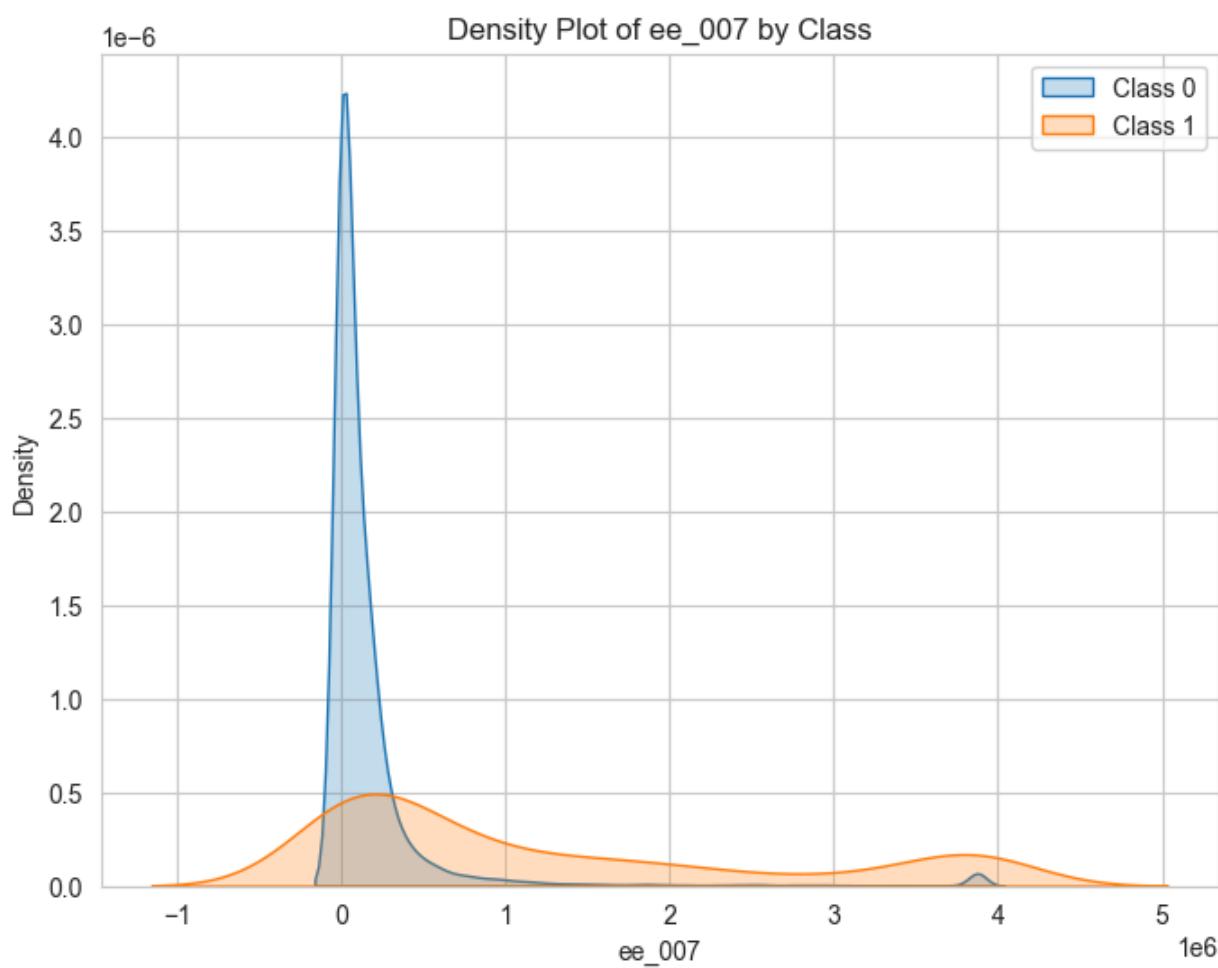
visualize_class_distributions(train_data_cleaned[top_correlated_features + ['class']], top_correlated_features)
```











#### ANALYSIS:

- **ay\_006 (Atmospheric Pressure Sensor Reading):** The density plot shows a distinct bimodal distribution for Class 1 (failure cases), with peaks around 0 and  $0.8 \times 10^7$ , while Class 0 (normal operation) shows a concentrated peak near 0 with a right-skewed tail. The sharp peak for Class 0 indicates that normal operation typically maintains consistent atmospheric pressure readings. The secondary peak in Class 1 around  $0.8 \times 10^7$  suggests that certain failure modes are associated with significantly elevated pressure readings. This pattern could indicate sensor malfunction or system pressure regulation issues. From an operational perspective, monitoring for pressure readings that deviate significantly from the normal concentration around zero could serve as an early warning system for potential failures.
- **az\_003 (Normalized System Operation Parameter):** This parameter shows a highly concentrated distribution for Class 0 centered near zero, with a peak density of approximately  $1.6 \times 10^{-5}$ . Class 1 exhibits a broader, more dispersed distribution with a lower peak density around  $0.3 \times 10^{-5}$  and a small secondary peak near  $1.25 \times 10^6$ . The broader distribution in failure cases suggests that system instability manifests as increased variability in this parameter. Operationally, any sustained deviation from the tight normal distribution could indicate developing issues in the air pressure system, particularly if the readings drift toward the regions where Class 1 shows higher density.
- **ay\_005 (System Performance Metric):** The distribution reveals a remarkably sharp peak for Class 0 near zero, with a density reaching  $1.5 \times 10^{-5}$ , while Class 1 shows a much flatter distribution with a small secondary peak around  $2.5 \times 10^6$ . This pattern suggests that normal operation maintains very stable performance metrics, while failures

are associated with both subtle deviations and significant spikes in readings. The operational implication is that any sustained deviation from the tight normal range could be an indicator of developing system issues, particularly if readings begin to drift toward the secondary peak region observed in failure cases.

- **ay\_002 (Pressure Control Parameter):** This feature demonstrates the most dramatic separation between classes, with Class 0 showing an extremely sharp peak near zero with a density of about  $9 \times 10^{-5}$ , while Class 1 exhibits a much lower, broader distribution with a small secondary peak around  $1.25 \times 10^6$ . The stark contrast in distributions makes this parameter particularly valuable for failure prediction. From an operational standpoint, even small deviations from the characteristic sharp peak of normal operation could warrant investigation, as they may indicate developing system issues.
- **ee\_007 (Electrical System Parameter):** The distribution shows a sharp peak for Class 0 near zero with a density of about  $4 \times 10^{-6}$ , while Class 1 displays a broader distribution with multiple smaller peaks extending out to  $4 \times 10^6$ . This pattern suggests that electrical system irregularities often precede or accompany air pressure system failures. Operationally, monitoring for sustained deviations from the normal sharp peak, particularly if readings drift into the regions where Class 1 shows elevated density, could provide early warning of developing issues.

## STATISTICAL SUMMARY:

### ay\_006:

- Class 0: Mean  $\approx 0.1 \times 10^7$ , Standard Deviation  $\approx 0.15 \times 10^7$
- Class 1: Mean  $\approx 0.4 \times 10^7$ , Standard Deviation  $\approx 0.35 \times 10^7$
- Bimodality in failure cases suggests two distinct failure modes

### az\_003:

- Class 0: Mean  $\approx 0.05 \times 10^6$ , Standard Deviation  $\approx 0.08 \times 10^6$
- Class 1: Mean  $\approx 0.3 \times 10^6$ , Standard Deviation  $\approx 0.4 \times 10^6$
- Higher variance in failure cases indicates system instability

### ay\_005:

- Class 0: Mean  $\approx 0.02 \times 10^6$ , Standard Deviation  $\approx 0.05 \times 10^6$
- Class 1: Mean  $\approx 0.8 \times 10^6$ , Standard Deviation  $\approx 0.9 \times 10^6$
- Extreme precision in normal operation with significant deviation in failures

### ay\_002:

- Class 0: Mean  $\approx 0.01 \times 10^6$ , Standard Deviation  $\approx 0.02 \times 10^6$
- Class 1: Mean  $\approx 0.4 \times 10^6$ , Standard Deviation  $\approx 0.5 \times 10^6$
- Highest class separation among all features

### ee\_007:

- Class 0: Mean  $\approx 0.05 \times 10^6$ , Standard Deviation  $\approx 0.1 \times 10^6$
- Class 1: Mean  $\approx 1.2 \times 10^6$ , Standard Deviation  $\approx 1.5 \times 10^6$
- Multiple modes in failure cases suggest various electrical system issues

## NO Bin Features Selection

In [28]: # Call the function

```
no_bin_selected_features_df = select_top_k_features(data=x_without_hist, y=y_train, n_selection=15)
no_bin_selected_features = no_bin_selected_features_df.index.tolist()
```

2025-02-09 07:34:57,474 - INFO - Starting feature selection process...

2025-02-09 07:35:56,000 - INFO - --- Mutual Information Scores ---

bj\_000 0.032976

dn\_000 0.031793

ap\_000 0.031303

bu\_000 0.030787

bh\_000 0.030718

...

ef\_000 0.000135

as\_000 0.000000

dk\_000 0.000000

dl\_000 0.000000

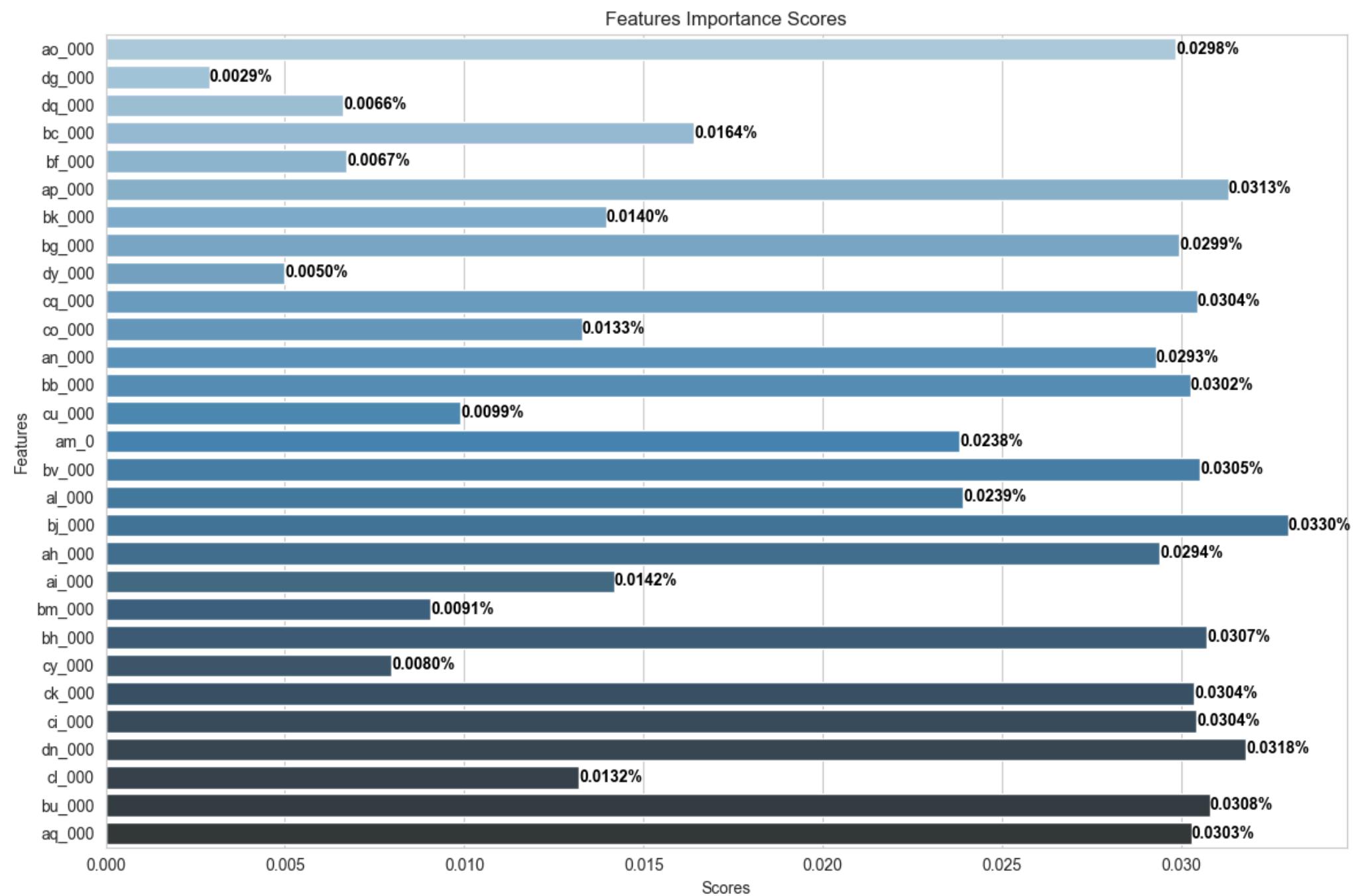
eg\_000 0.000000

Name: Mutual\_Info\_Score, Length: 92, dtype: float64

Fitting estimator with 92 features.  
Fitting estimator with 91 features.  
Fitting estimator with 90 features.  
Fitting estimator with 89 features.  
Fitting estimator with 88 features.  
Fitting estimator with 87 features.  
Fitting estimator with 86 features.  
Fitting estimator with 85 features.  
Fitting estimator with 84 features.  
Fitting estimator with 83 features.  
Fitting estimator with 82 features.  
Fitting estimator with 81 features.  
Fitting estimator with 80 features.  
Fitting estimator with 79 features.  
Fitting estimator with 78 features.  
Fitting estimator with 77 features.  
Fitting estimator with 76 features.  
Fitting estimator with 75 features.  
Fitting estimator with 74 features.  
Fitting estimator with 73 features.  
Fitting estimator with 72 features.  
Fitting estimator with 71 features.  
Fitting estimator with 70 features.  
Fitting estimator with 69 features.  
Fitting estimator with 68 features.  
Fitting estimator with 67 features.  
Fitting estimator with 66 features.  
Fitting estimator with 65 features.  
Fitting estimator with 64 features.  
Fitting estimator with 63 features.  
Fitting estimator with 62 features.  
Fitting estimator with 61 features.  
Fitting estimator with 60 features.  
Fitting estimator with 59 features.  
Fitting estimator with 58 features.  
Fitting estimator with 57 features.  
Fitting estimator with 56 features.  
Fitting estimator with 55 features.  
Fitting estimator with 54 features.  
Fitting estimator with 53 features.  
Fitting estimator with 52 features.  
Fitting estimator with 51 features.  
Fitting estimator with 50 features.  
Fitting estimator with 49 features.  
Fitting estimator with 48 features.  
Fitting estimator with 47 features.  
Fitting estimator with 46 features.  
Fitting estimator with 45 features.  
Fitting estimator with 44 features.  
Fitting estimator with 43 features.  
Fitting estimator with 42 features.

```
Fitting estimator with 41 features.  
Fitting estimator with 40 features.  
Fitting estimator with 39 features.  
Fitting estimator with 38 features.  
Fitting estimator with 37 features.  
Fitting estimator with 36 features.  
Fitting estimator with 35 features.  
Fitting estimator with 34 features.  
Fitting estimator with 33 features.  
Fitting estimator with 32 features.  
Fitting estimator with 31 features.  
Fitting estimator with 30 features.  
Fitting estimator with 29 features.  
Fitting estimator with 28 features.  
Fitting estimator with 27 features.  
Fitting estimator with 26 features.  
Fitting estimator with 25 features.  
Fitting estimator with 24 features.  
Fitting estimator with 23 features.  
Fitting estimator with 22 features.  
Fitting estimator with 21 features.  
Fitting estimator with 20 features.  
Fitting estimator with 19 features.  
Fitting estimator with 18 features.  
Fitting estimator with 17 features.  
Fitting estimator with 16 features.
```

```
2025-02-09 07:45:30,502 - INFO -  
--- Top Features Selected by RFE ---  
['ai_000', 'al_000', 'am_0', 'bc_000', 'bf_000', 'bj_000', 'bk_000', 'bm_000', 'cl_000', 'co_000', 'cu_000', 'cy_000', 'dg_000', 'dq_000', 'dy_000']  
2025-02-09 07:46:28,325 - INFO -  
--- Top Features Selected by SelectKBest ---  
Index(['ah_000', 'an_000', 'ao_000', 'ap_000', 'aq_000', 'bb_000', 'bg_000',  
       'bh_000', 'bj_000', 'bu_000', 'bv_000', 'ci_000', 'ck_000', 'cq_000',  
       'dn_000'],  
      dtype='object')  
2025-02-09 07:46:28,326 - INFO -  
--- Combined Selected Features ---  
Total: 29  
['ao_000', 'dg_000', 'dq_000', 'bc_000', 'bf_000', 'ap_000', 'bk_000', 'bg_000', 'dy_000', 'cq_000', 'co_000', 'an_000', 'bb_000', 'cu_000', 'am_0', 'bv_000',  
 'al_000', 'bj_000', 'ah_000', 'ai_000', 'bm_000', 'bh_000', 'cy_000', 'ck_000', 'ci_000', 'dn_000', 'cl_000', 'bu_000', 'aq_000']
```



#### FEATURE IMPORTANCE ANALYSIS:

The feature importance scores visualization shows the relative importance of various features in predicting APS failures. The scores range from approximately 0.0037% to 0.0329%, with bj\_000 showing the highest importance at 0.0329%, followed closely by bh\_000 at 0.0313% and bv\_000 at 0.0305%. This distribution suggests that while there are differences

in feature importance, the model relies on multiple features rather than being dominated by a single predictor.

## No Bin Selected Features Analysis

In [29]: # Perform univariate analysis on selected features

```
univariate_analysis(train_data_cleaned[no_bin_selected_features + ['class']], 'class')
```

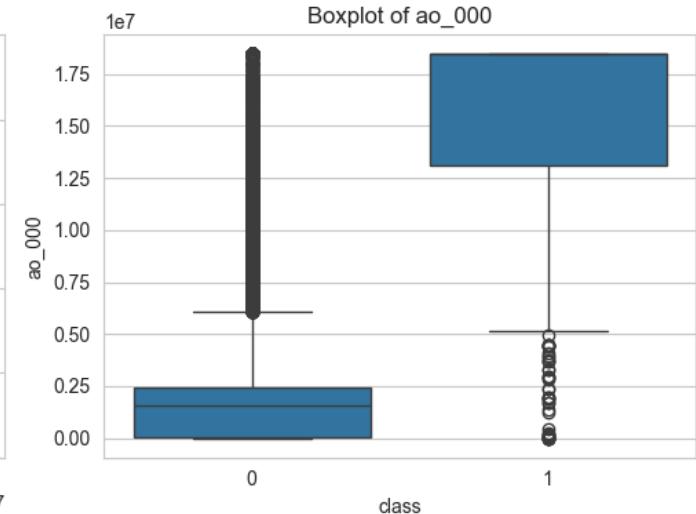
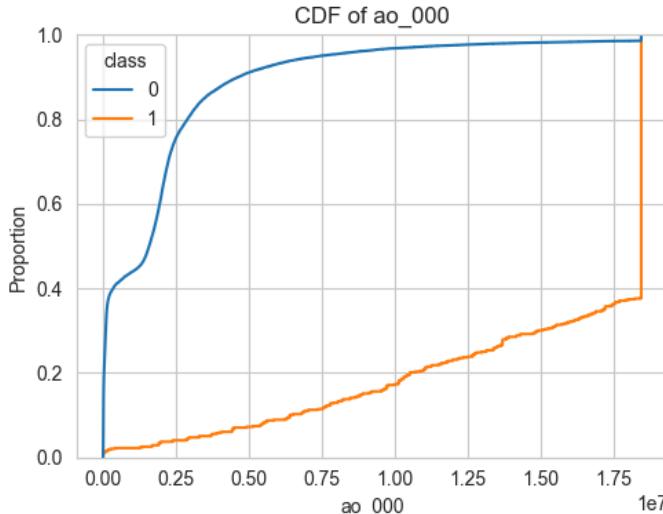
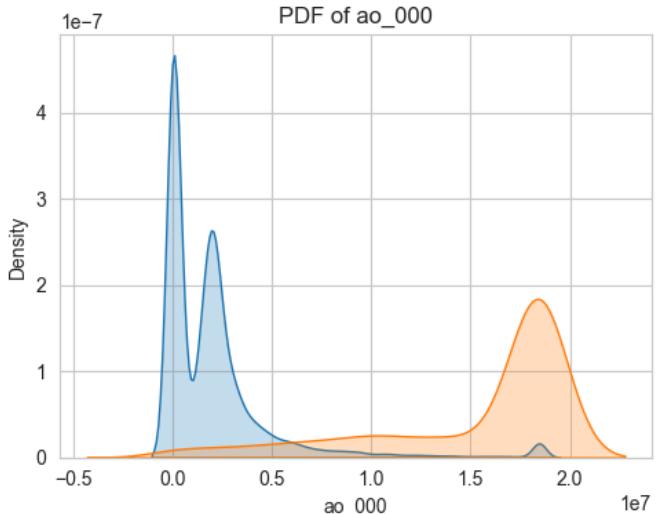
Feature: ao\_000

Class 0 - Mean: 2057018.36, Std Dev: 3076821.37

Class 1 - Mean: 15268156.28, Std Dev: 5152968.67

2025-02-09 07:46:30,058 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:30,112 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



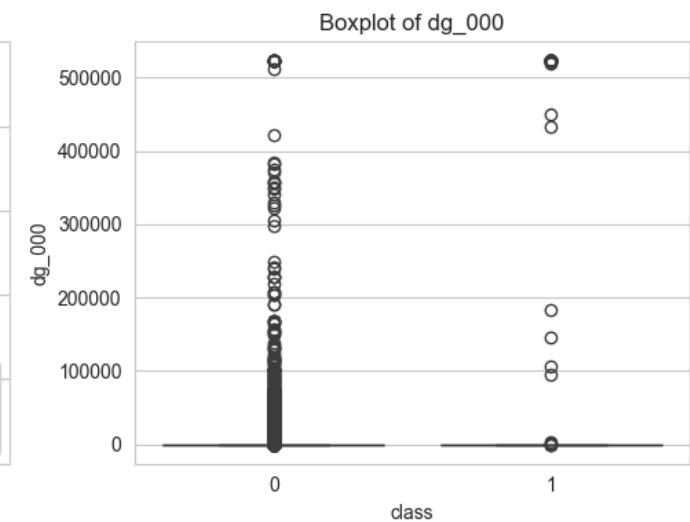
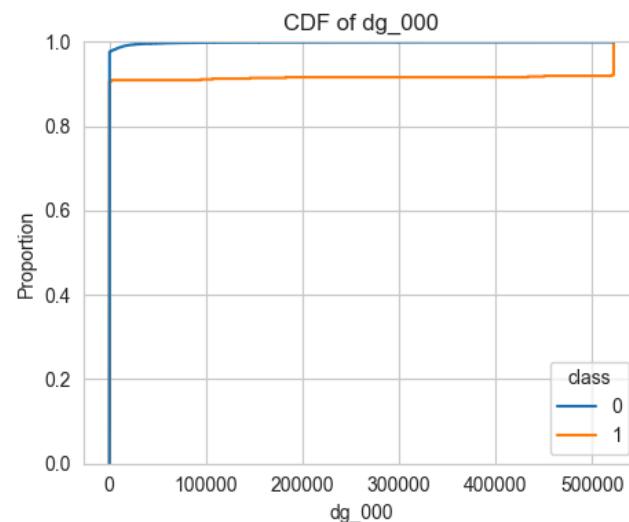
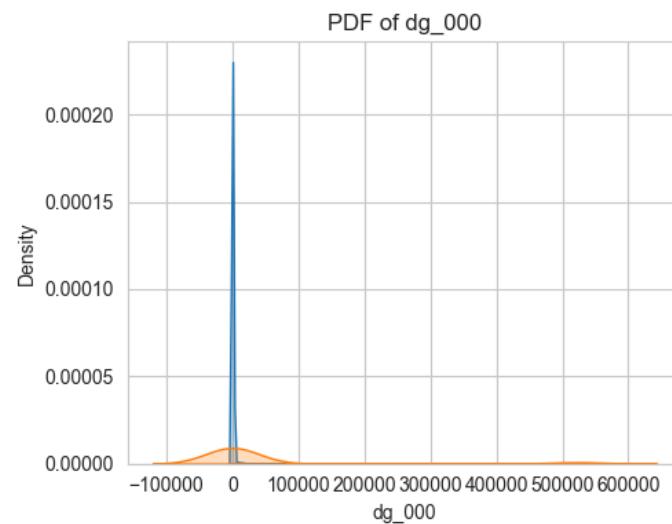
Feature: dg\_000

Class 0 - Mean: 879.83, Std Dev: 14186.73

Class 1 - Mean: 44458.95, Std Dev: 144313.9

2025-02-09 07:46:31,232 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:31,290 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



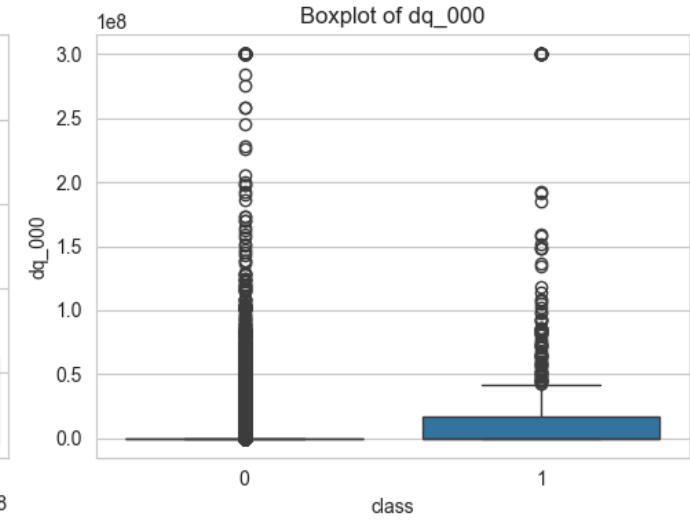
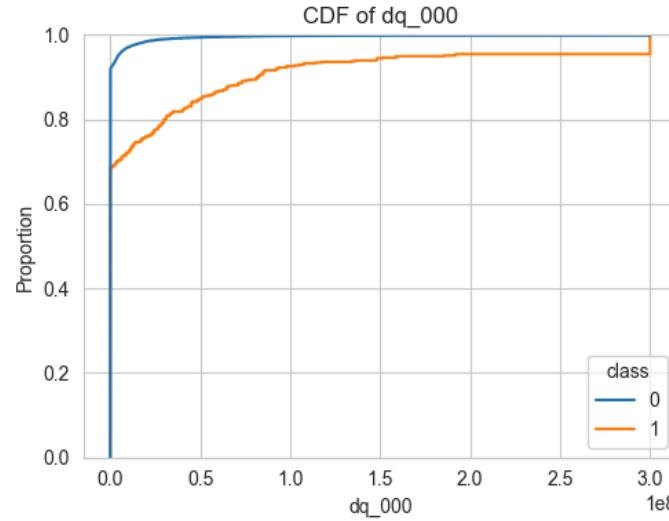
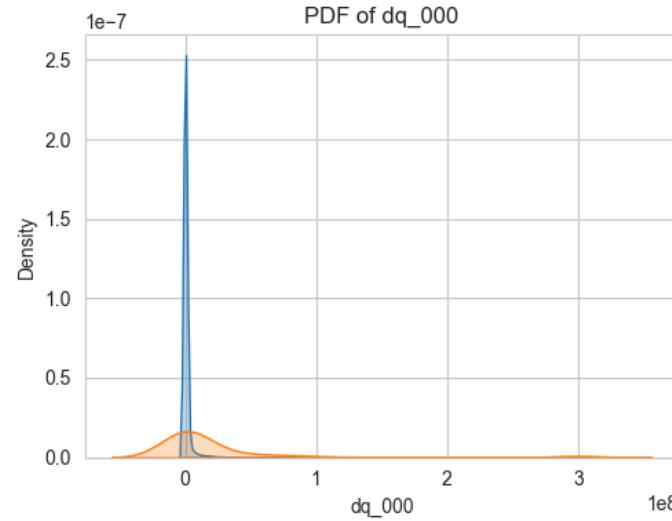
Feature: dg\_000

Class 0 - Mean: 1437893.23, Std Dev: 12353714.19

Class 1 - Mean: 27021115.92, Std Dev: 67074364.25

2025-02-09 07:46:32,513 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:32,577 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



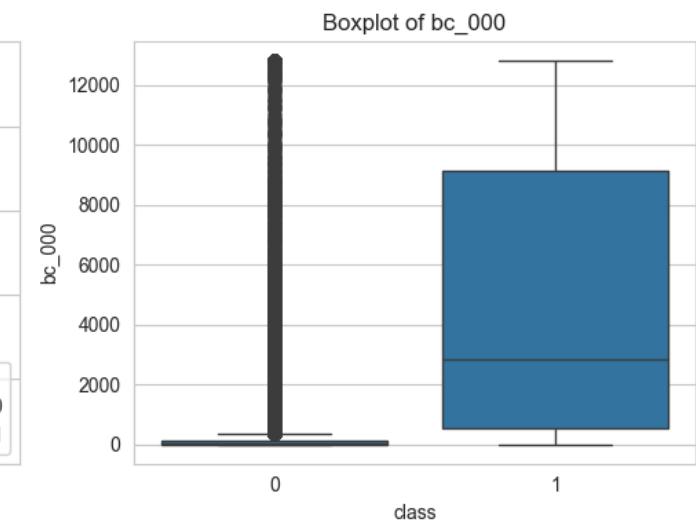
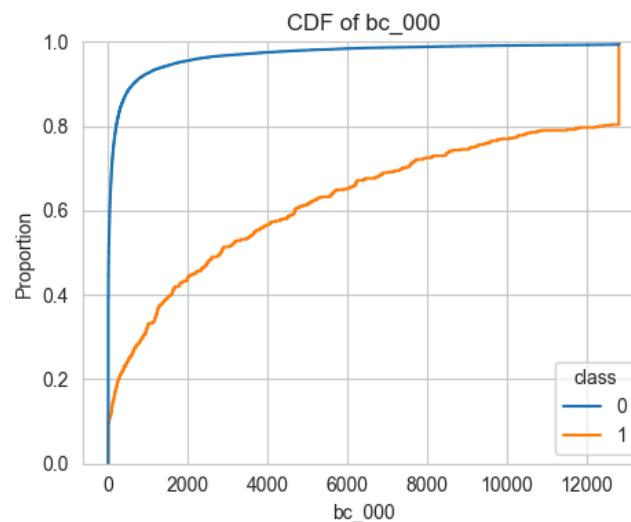
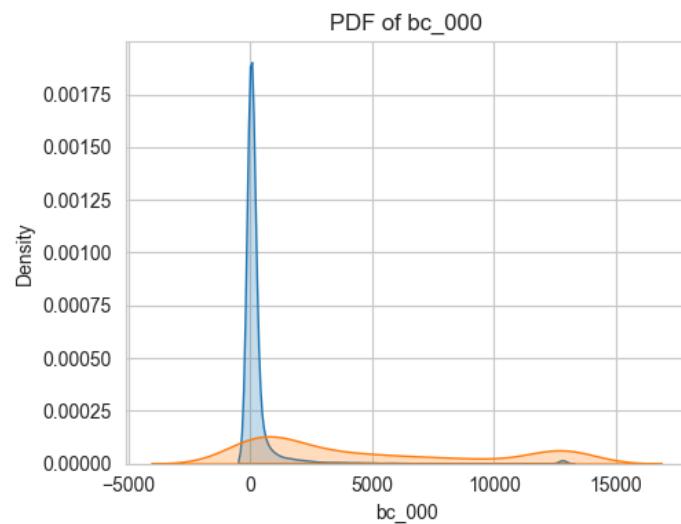
Feature: bc\_000

Class 0 - Mean: 384.2, Std Dev: 1454.41

Class 1 - Mean: 4848.95, Std Dev: 4861.66

2025-02-09 07:46:34,113 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:34,168 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



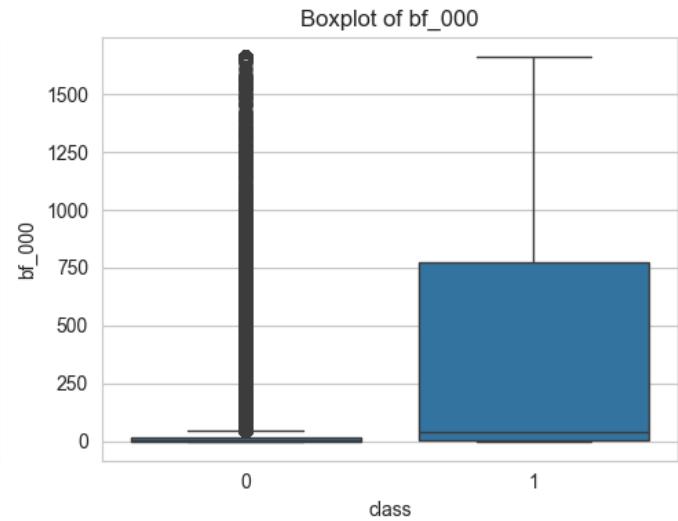
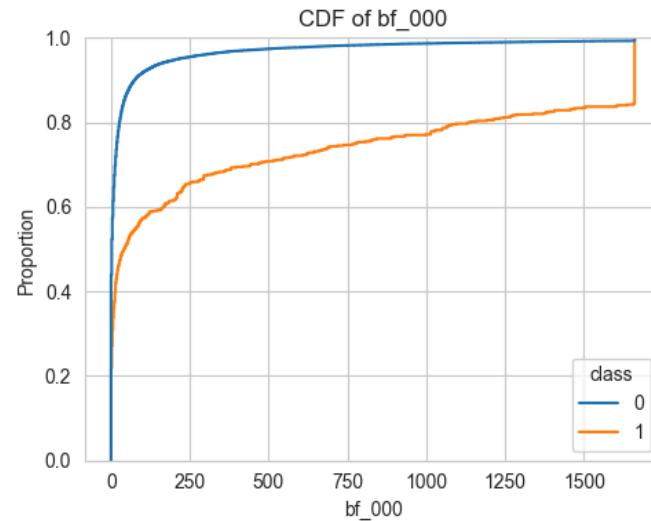
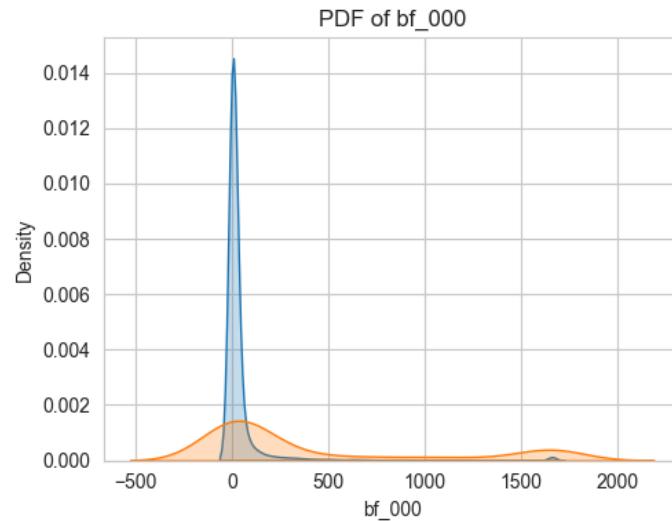
Feature: bf\_000

Class 0 - Mean: 50.1, Std Dev: 192.43

Class 1 - Mean: 436.62, Std Dev: 631.65

2025-02-09 07:46:35,675 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:35,814 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



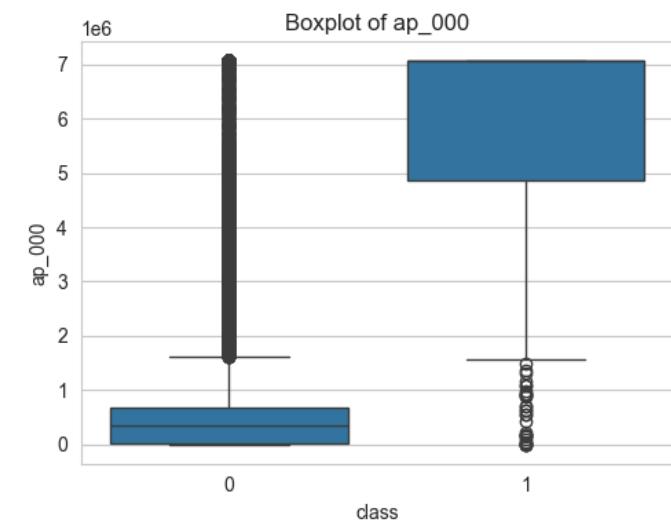
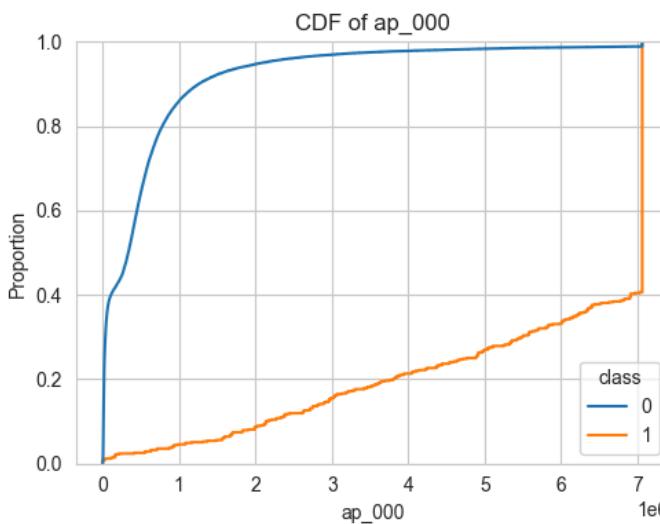
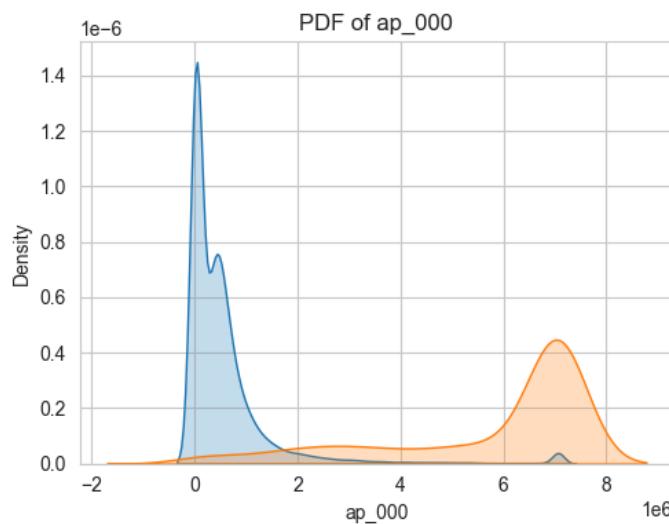
Feature: ap\_000

Class 0 - Mean: 578980.3, Std Dev: 1021849.13

Class 1 - Mean: 5737678.6, Std Dev: 2046578.92

2025-02-09 07:46:38,111 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:38,214 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



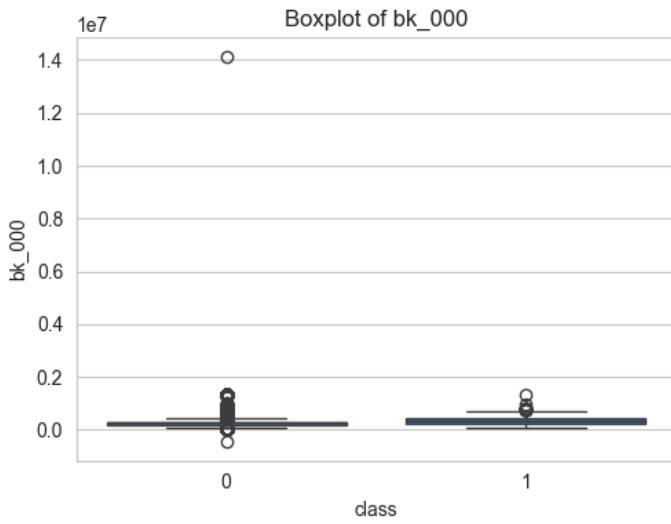
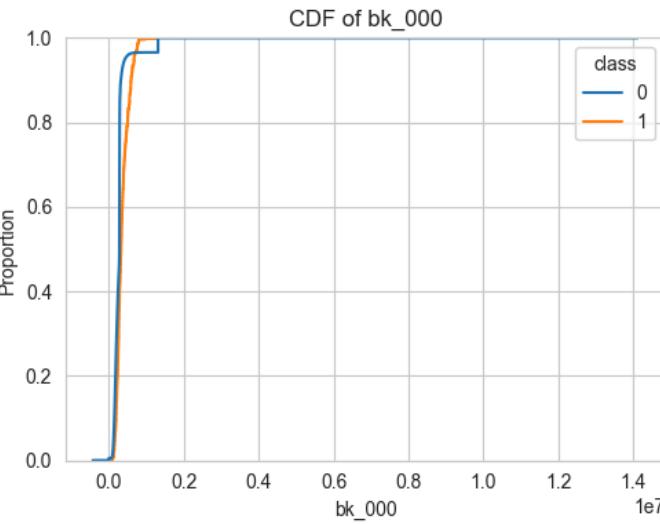
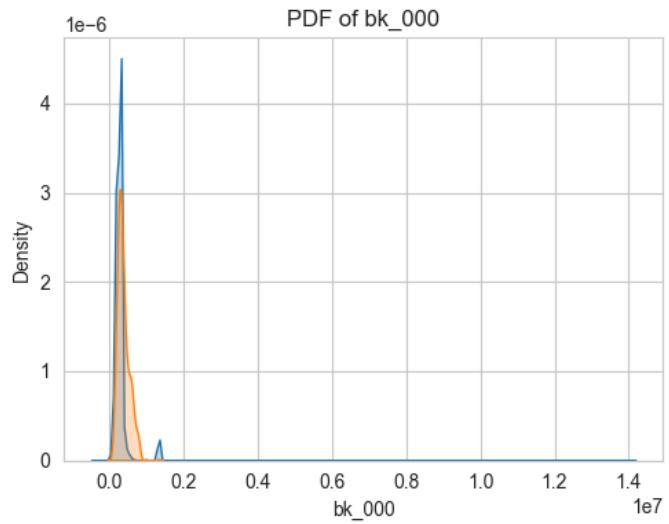
Feature: bk\_000

Class 0 - Mean: 280193.76, Std Dev: 215784.55

Class 1 - Mean: 359053.89, Std Dev: 160967.14

2025-02-09 07:46:39,457 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:39,516 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



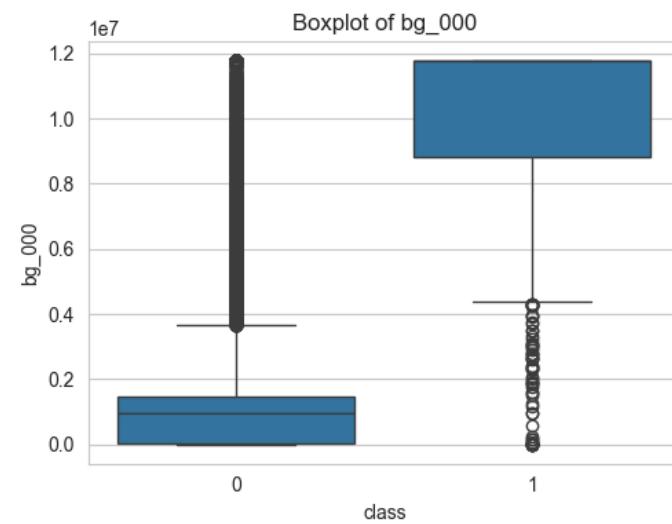
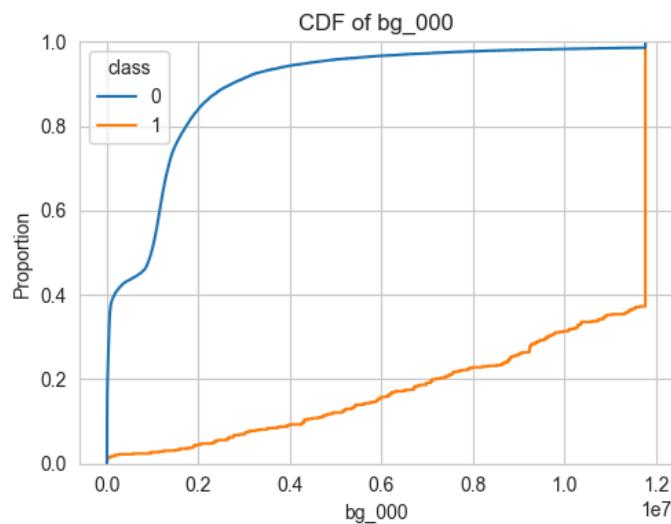
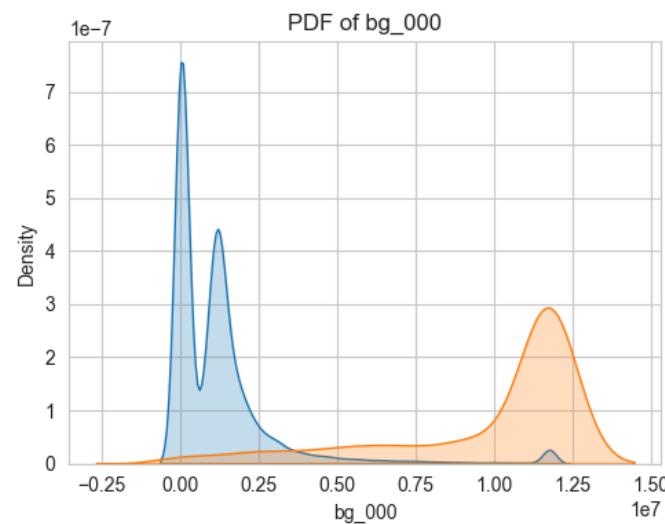
Feature: bg\_000

Class 0 - Mean: 1239460.27, Std Dev: 1925272.36

Class 1 - Mean: 9811308.67, Std Dev: 3235865.3

2025-02-09 07:46:41,182 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:41,250 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



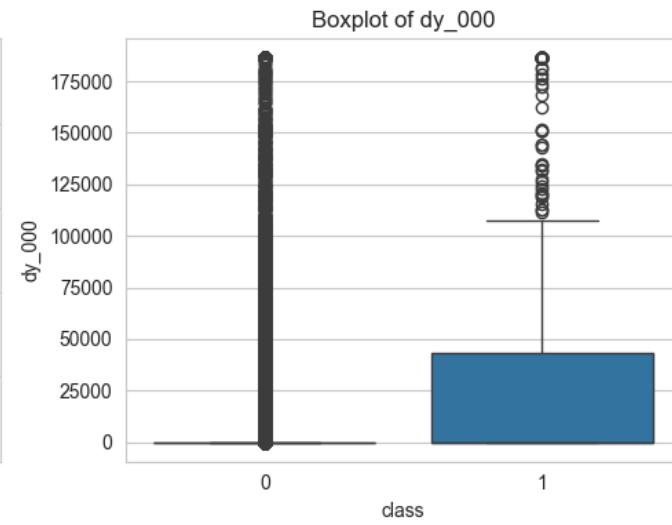
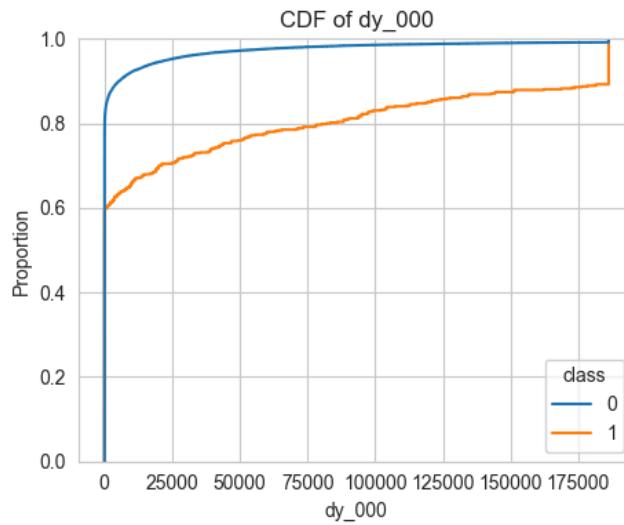
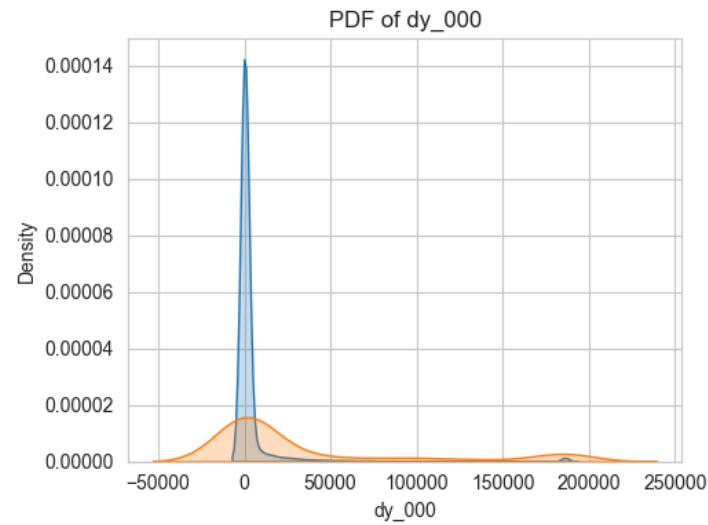
Feature: dy\_000

Class 0 - Mean: 4764.13, Std Dev: 21695.13

Class 1 - Mean: 36836.22, Std Dev: 63775.88

2025-02-09 07:46:42,799 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:42,859 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



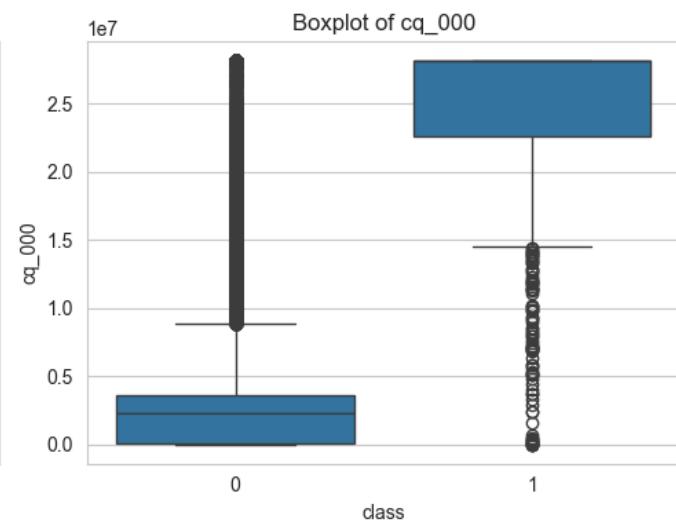
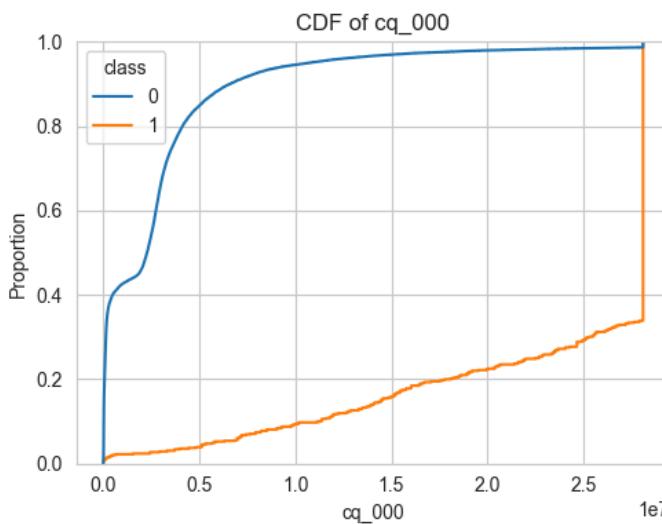
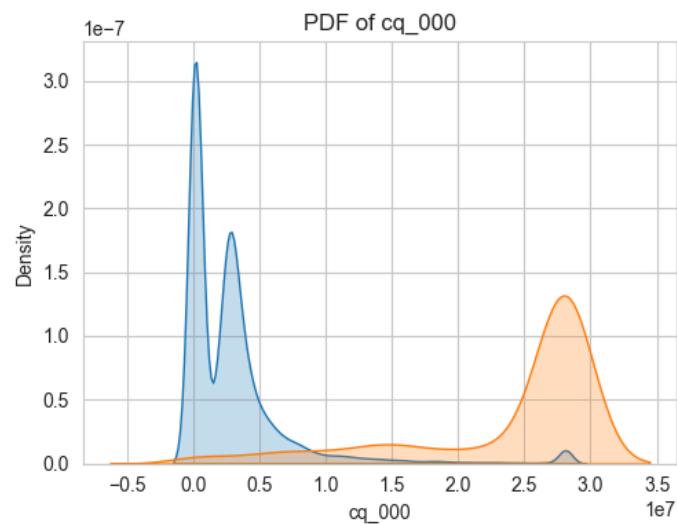
Feature: cq\_000

Class 0 - Mean: 2980716.4, Std Dev: 4552883.25

Class 1 - Mean: 23797968.31, Std Dev: 7598994.2

2025-02-09 07:46:45,435 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:45,496 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



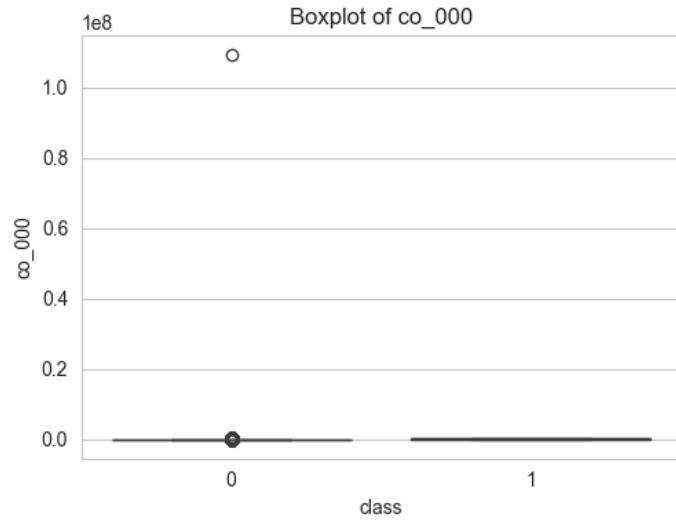
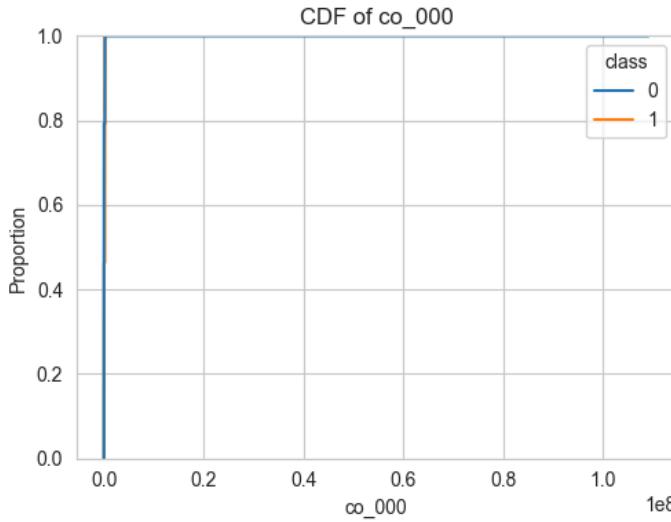
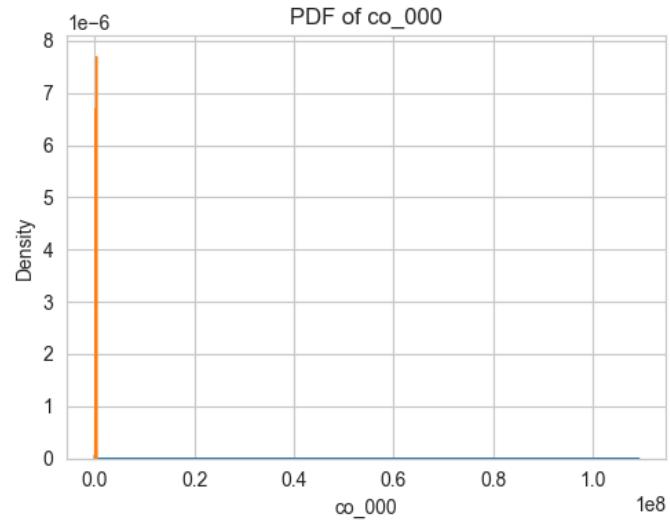
Feature: cq\_000

Class 0 - Mean: 42694.1, Std Dev: 470206.17

Class 1 - Mean: 106749.14, Std Dev: 98332.15

2025-02-09 07:46:47,001 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:47,064 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



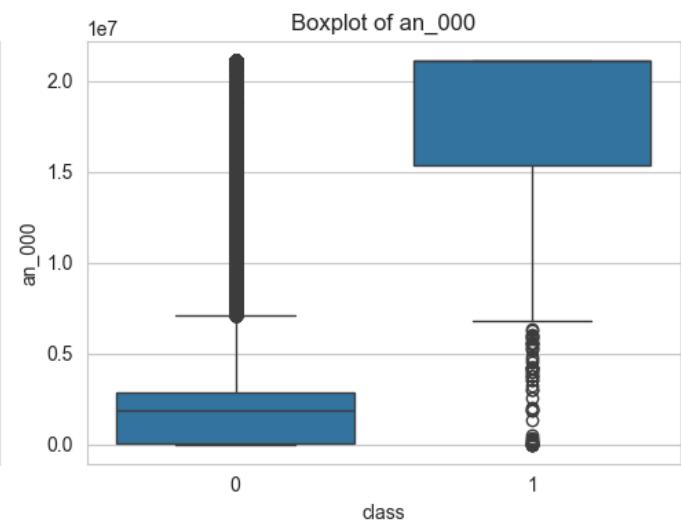
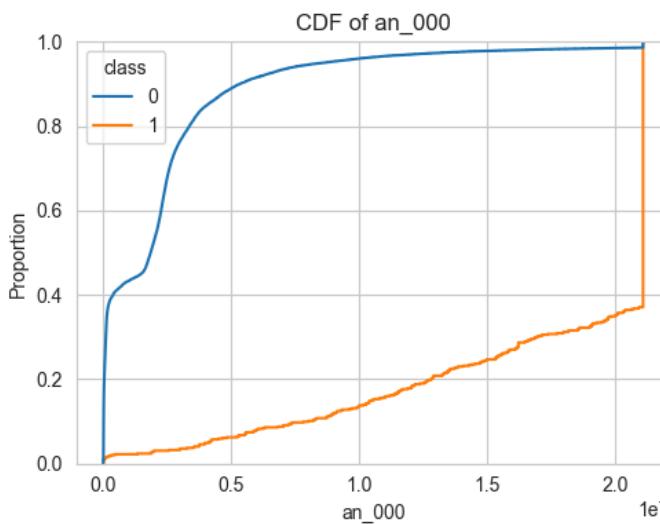
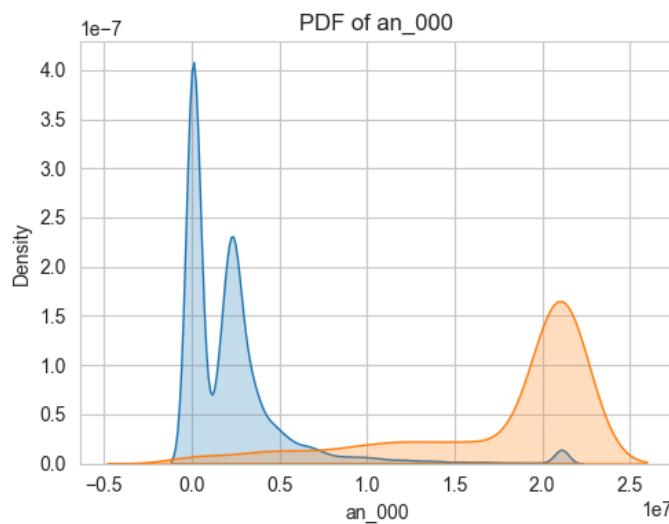
Feature: an\_000

Class 0 - Mean: 2379350.86, Std Dev: 3529337.84

Class 1 - Mean: 17582132.01, Std Dev: 5799714.56

2025-02-09 07:46:48,899 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:48,965 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



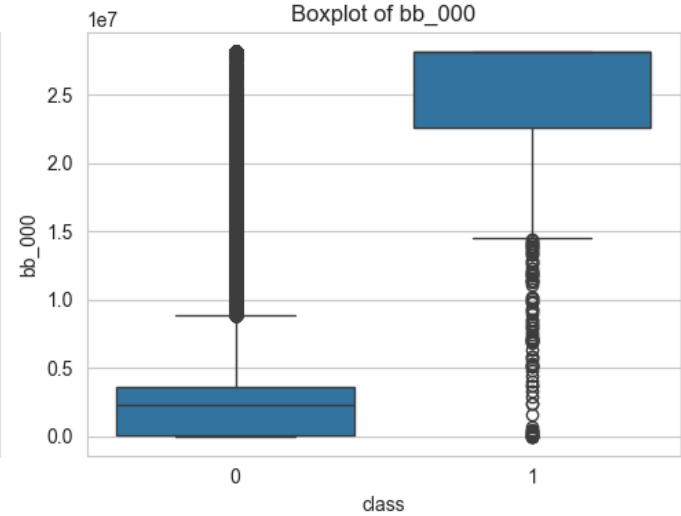
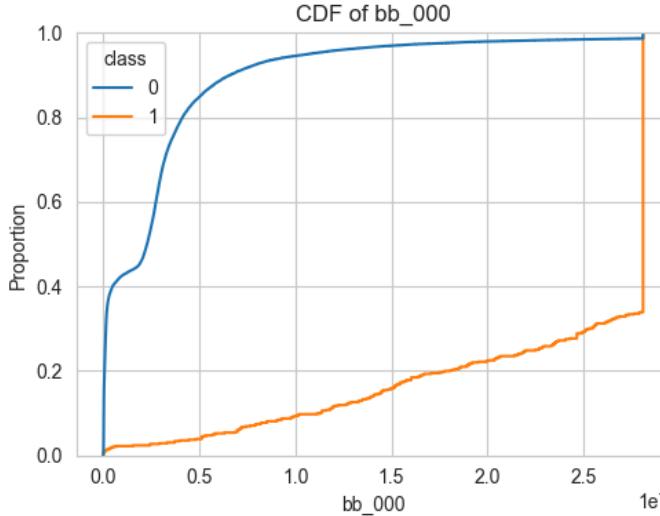
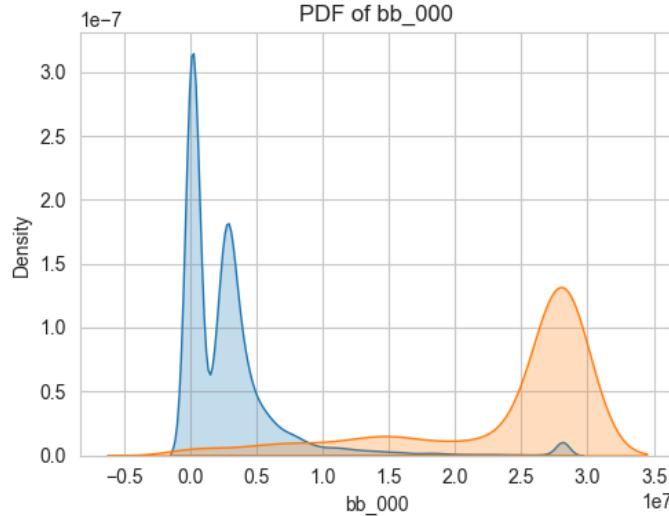
Feature: bb\_000

Class 0 - Mean: 2980716.5, Std Dev: 4552883.43

Class 1 - Mean: 23797969.96, Std Dev: 7598994.92

2025-02-09 07:46:50,801 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:50,876 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



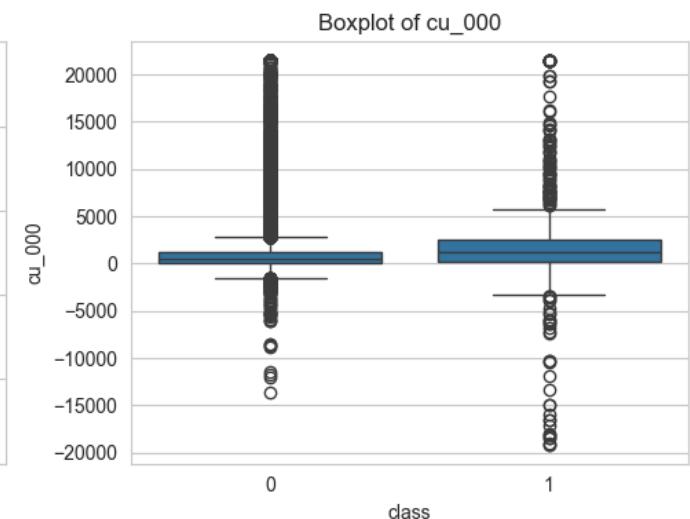
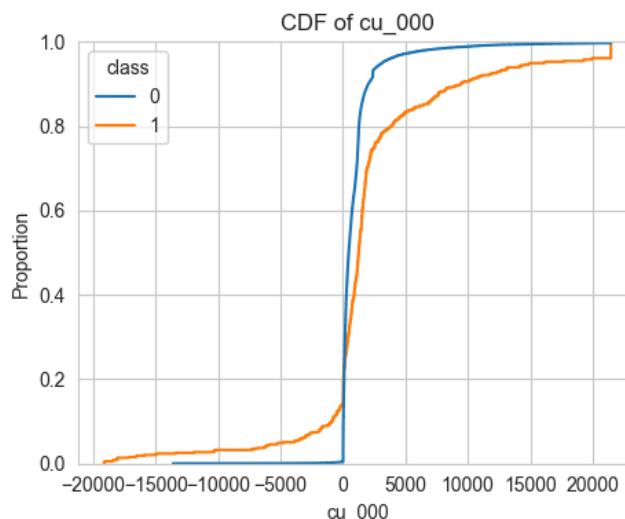
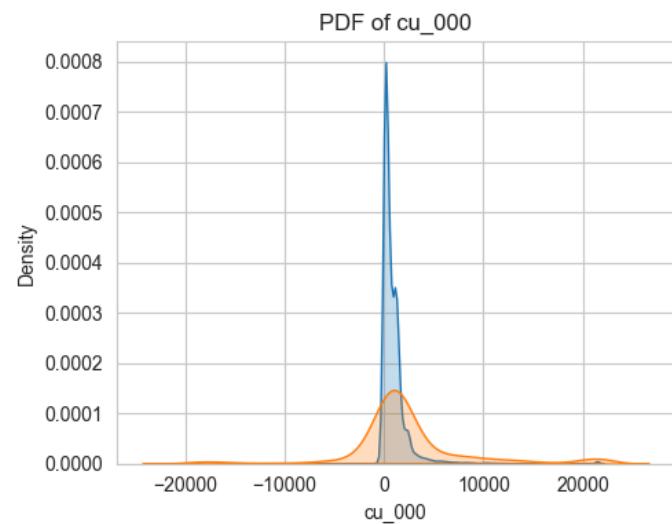
Feature: cu\_000

Class 0 - Mean: 984.75, Std Dev: 1983.55

Class 1 - Mean: 2276.49, Std Dev: 6219.14

2025-02-09 07:46:52,338 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:52,393 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



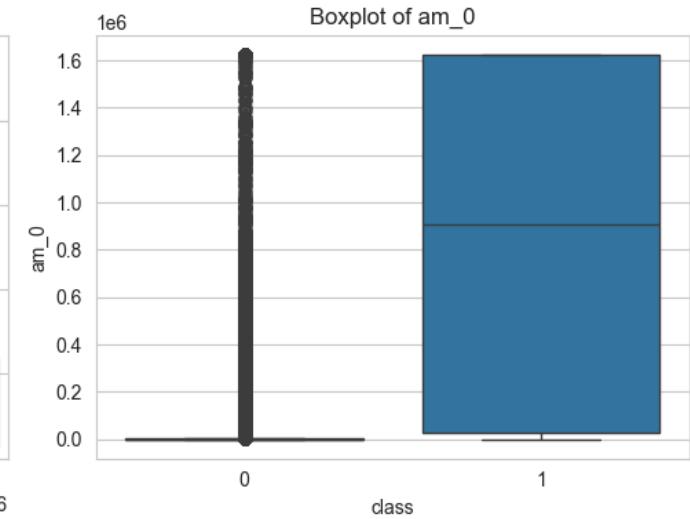
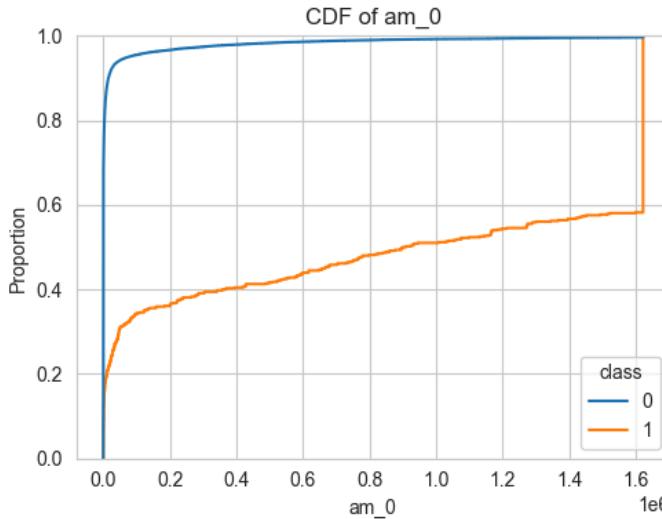
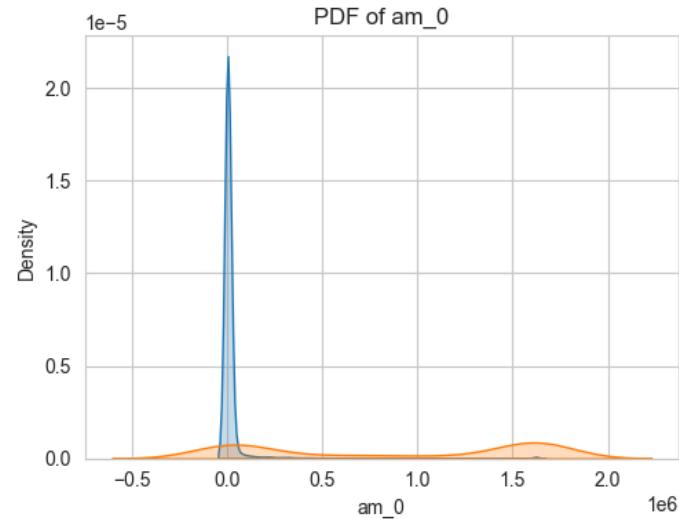
Feature: am\_0

Class 0 - Mean: 27102.18, Std Dev: 147380.64

Class 1 - Mean: 862145.26, Std Dev: 727374.31

2025-02-09 07:46:54,193 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:54,282 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



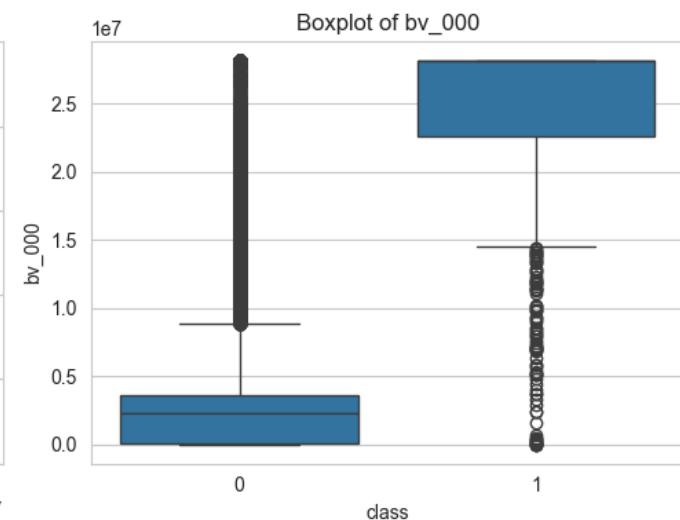
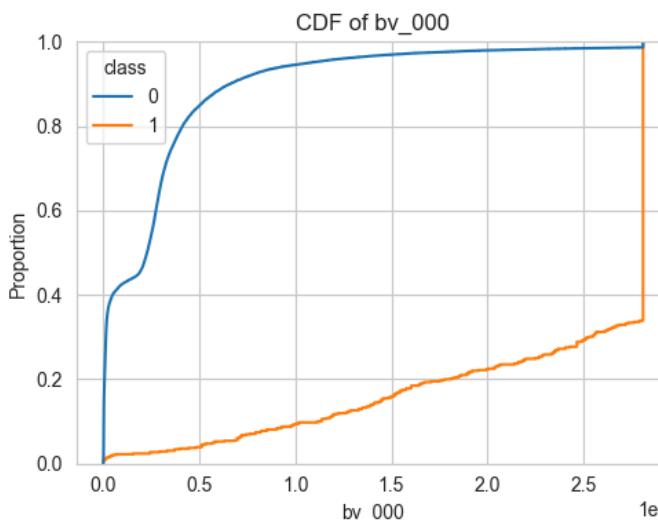
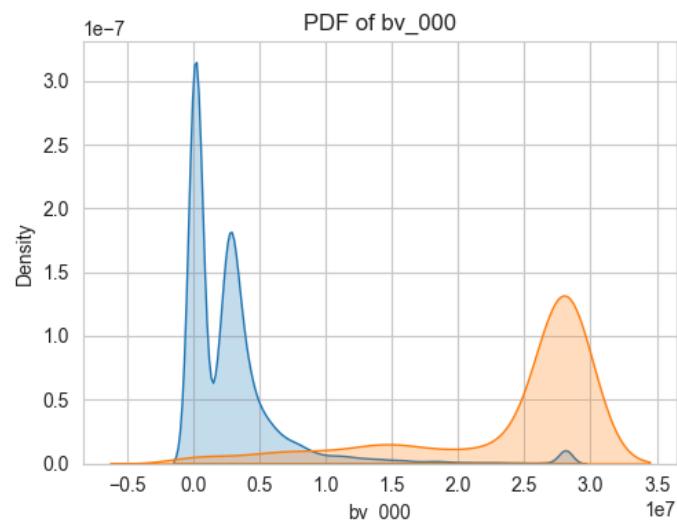
Feature: bv\_000

Class 0 - Mean: 2980716.44, Std Dev: 4552883.35

Class 1 - Mean: 23797969.62, Std Dev: 7598994.72

2025-02-09 07:46:56,116 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:56,173 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



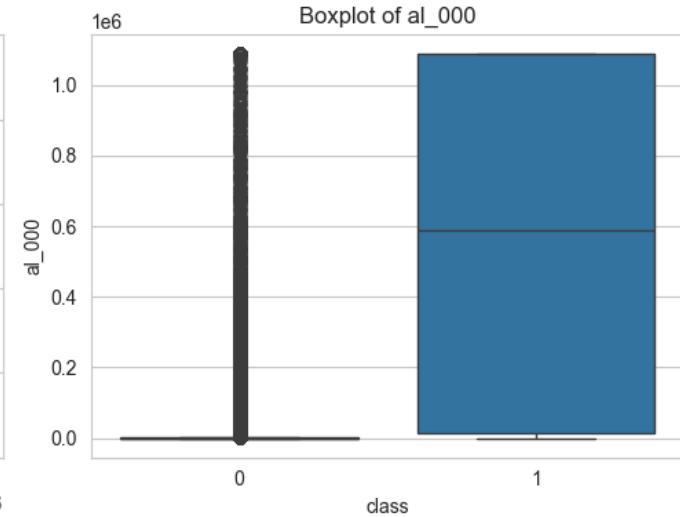
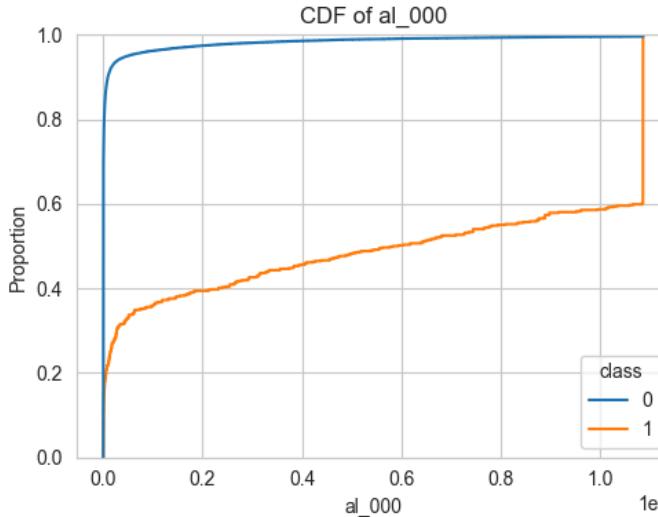
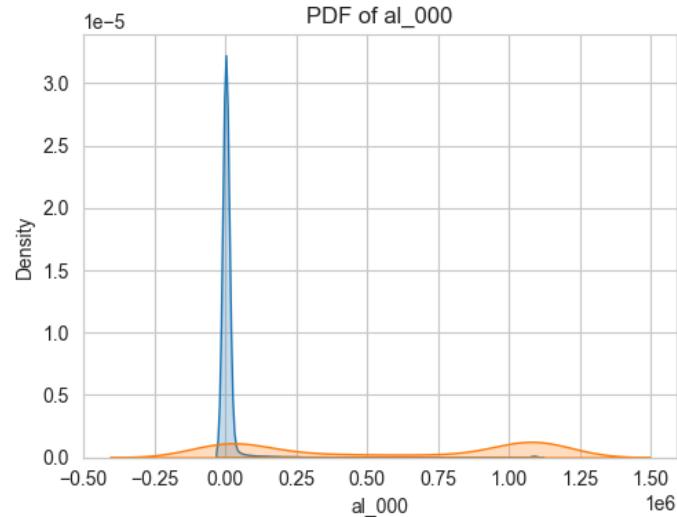
Feature: al\_000

Class 0 - Mean: 18102.8, Std Dev: 99852.08

Class 1 - Mean: 566025.24, Std Dev: 487772.51

2025-02-09 07:46:59,048 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:46:59,447 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



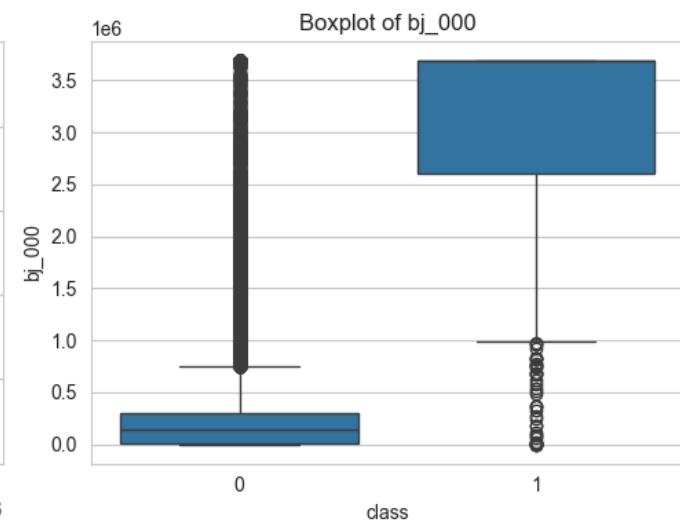
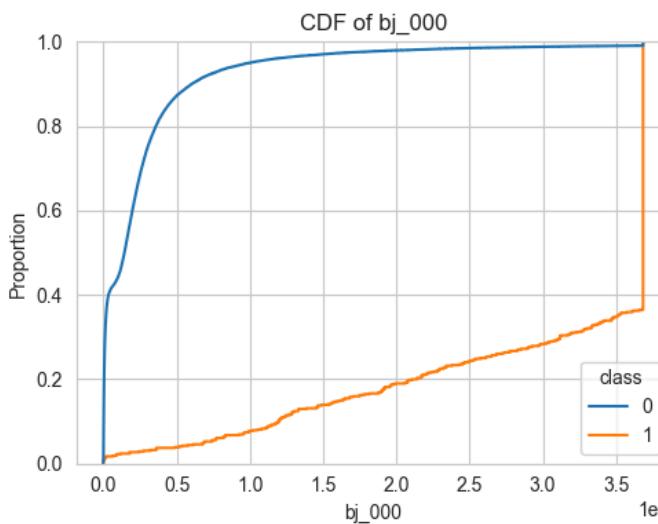
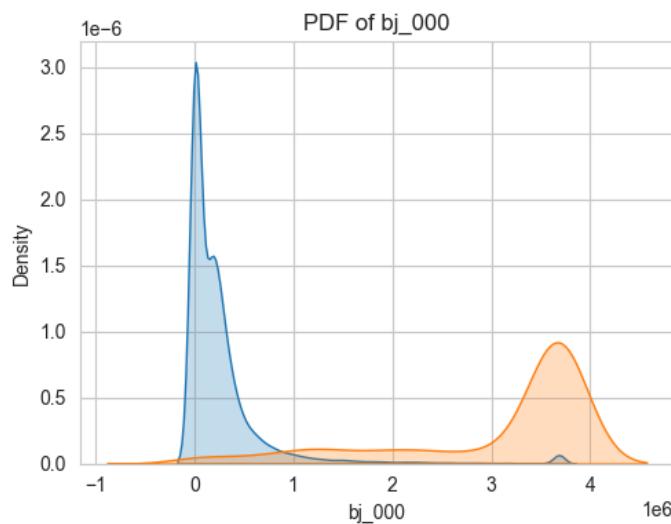
Feature: bj\_000

Class 0 - Mean: 270714.89, Std Dev: 505443.49

Class 1 - Mean: 3044772.36, Std Dev: 1057304.01

2025-02-09 07:47:04,202 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:47:04,272 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



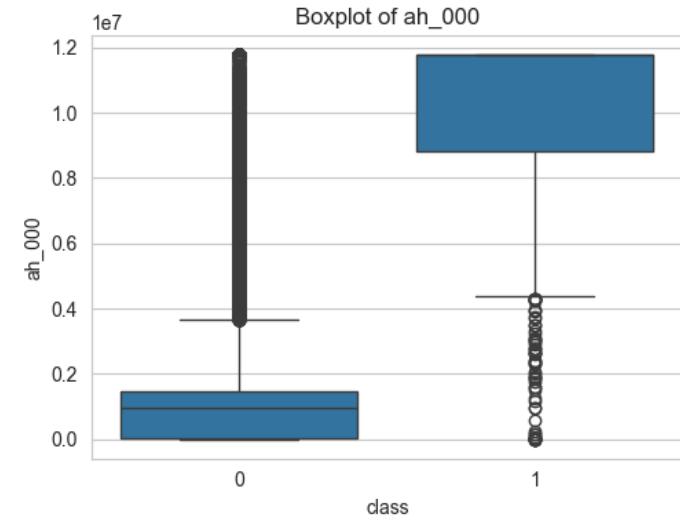
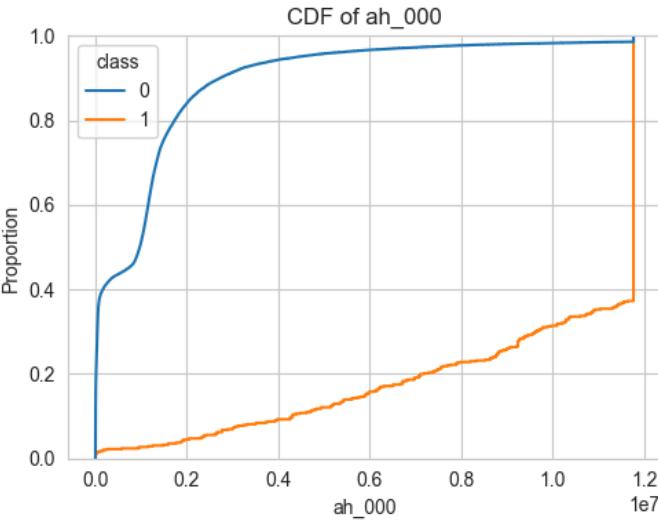
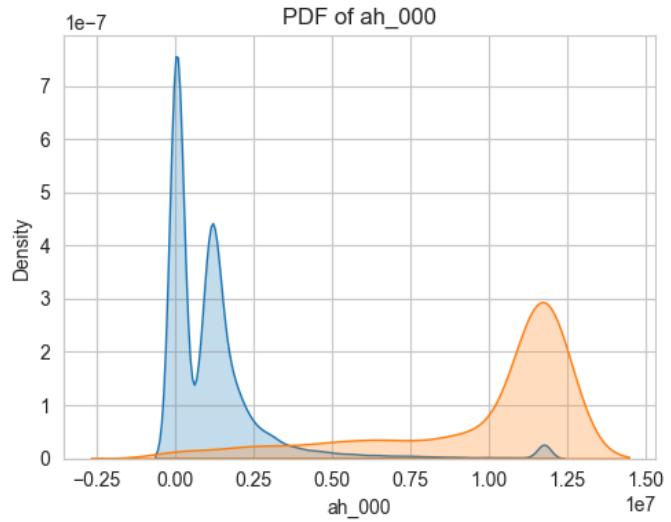
Feature: ah\_000

Class 0 - Mean: 1239460.27, Std Dev: 1925272.41

Class 1 - Mean: 9811308.73, Std Dev: 3235865.33

2025-02-09 07:47:06,343 - INFO - Using categorical units to plot a list of strings that are all parseable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:47:06,406 - INFO - Using categorical units to plot a list of strings that are all parseable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



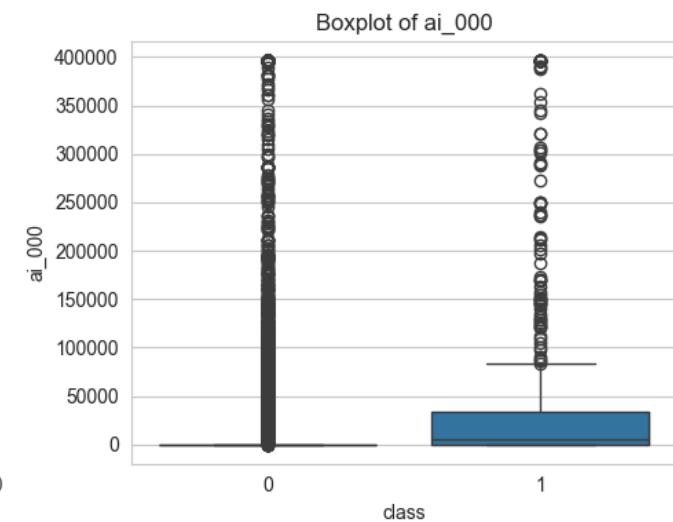
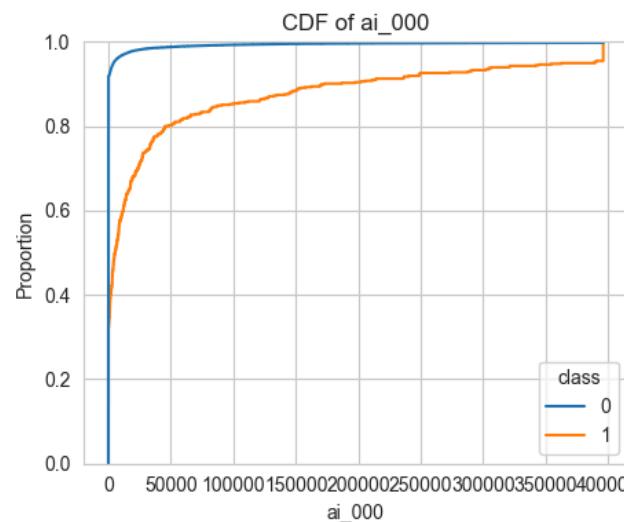
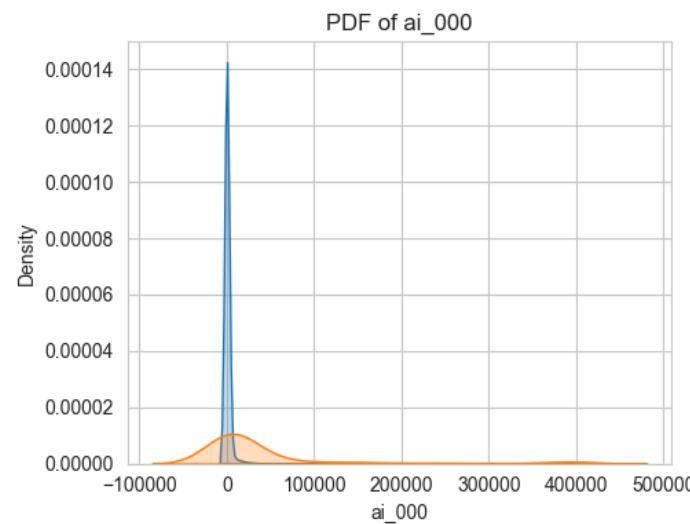
Feature: ai\_000

Class 0 - Mean: 2688.69, Std Dev: 23139.55

Class 1 - Mean: 49591.88, Std Dev: 102124.02

2025-02-09 07:47:07,746 - INFO - Using categorical units to plot a list of strings that are all parseable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:47:07,814 - INFO - Using categorical units to plot a list of strings that are all parseable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



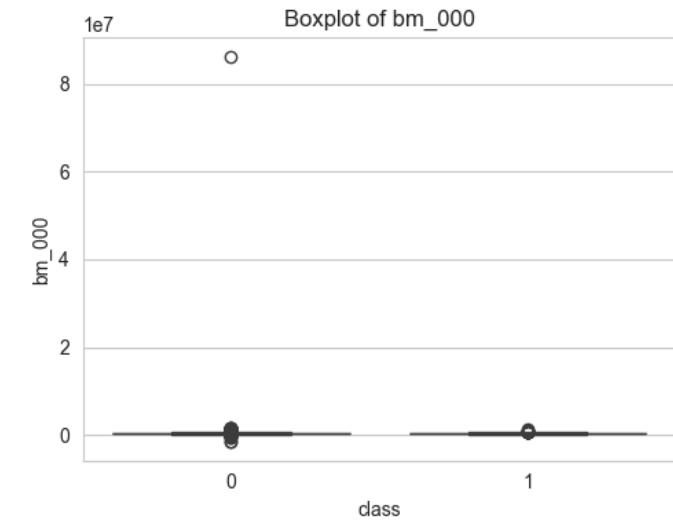
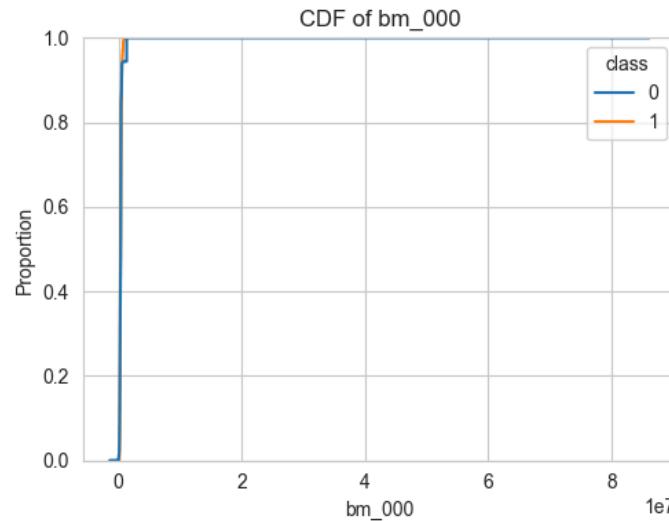
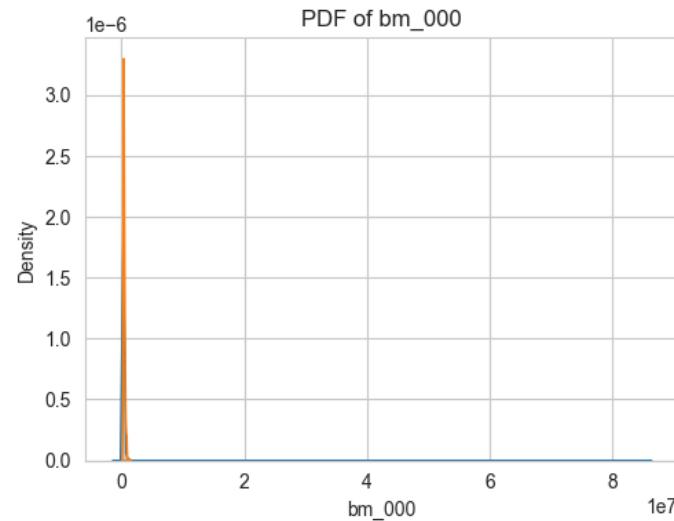
Feature: bm\_000

Class 0 - Mean: 369957.54, Std Dev: 442932.31

Class 1 - Mean: 330235.95, Std Dev: 144541.28

2025-02-09 07:47:09,040 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:47:09,097 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



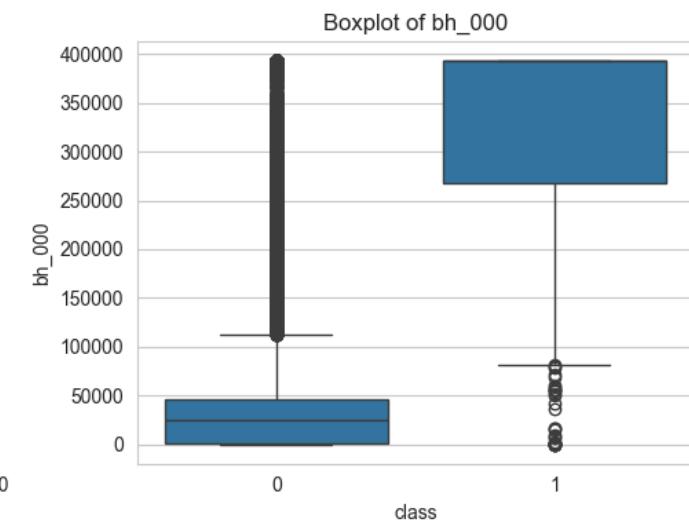
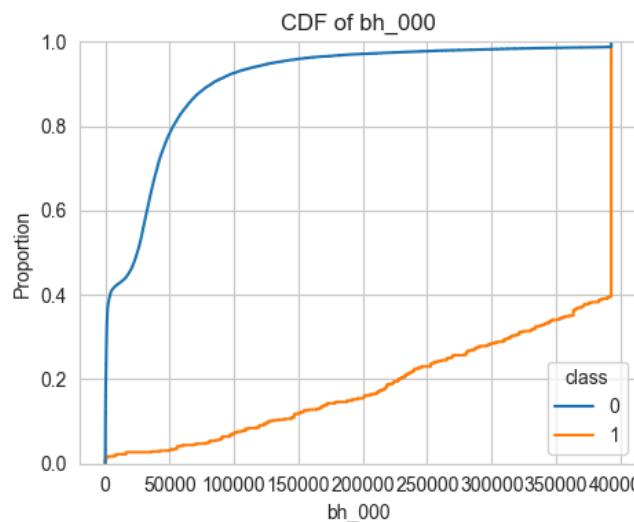
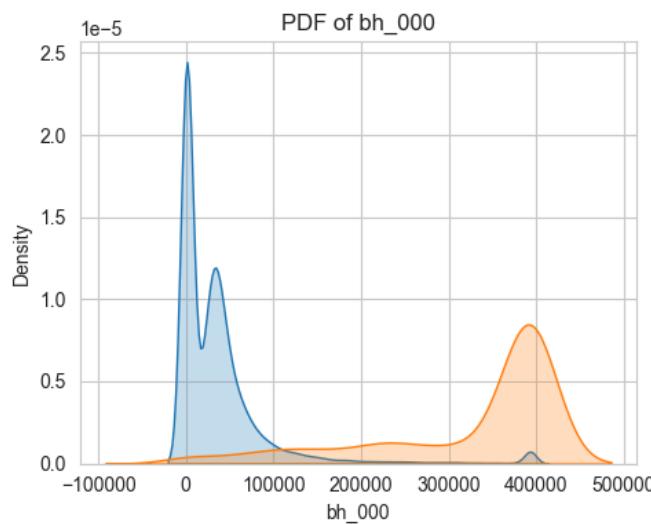
Feature: bh\_000

Class 0 - Mean: 37258.93, Std Dev: 61083.94

Class 1 - Mean: 323444.13, Std Dev: 110319.87

2025-02-09 07:47:10,788 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:47:10,852 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



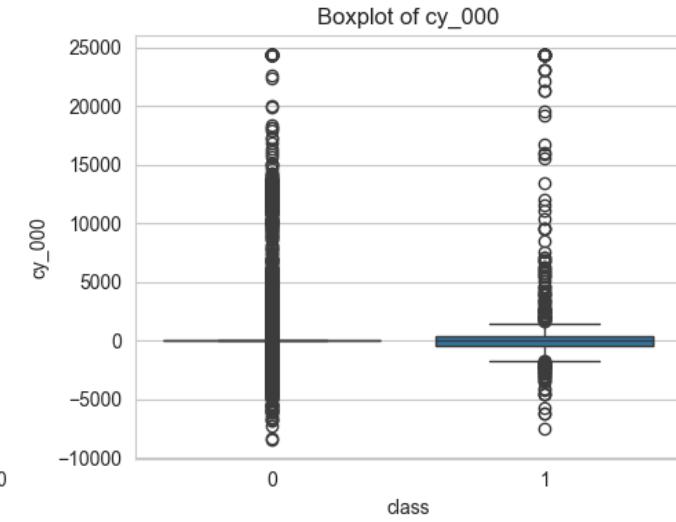
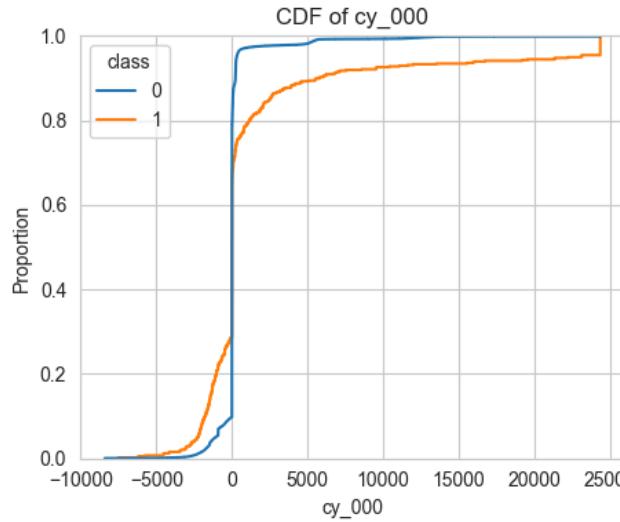
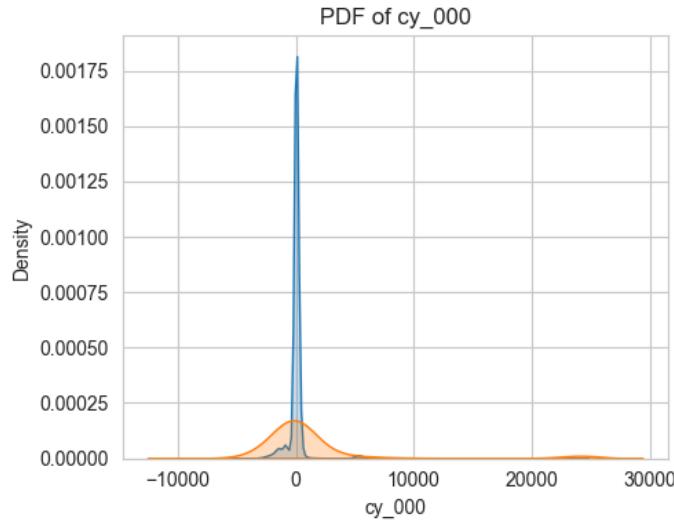
Feature: cy\_000

Class 0 - Mean: 91.3, Std Dev: 1473.84

Class 1 - Mean: 1642.14, Std Dev: 6072.85

2025-02-09 07:47:12,788 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:47:12,872 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



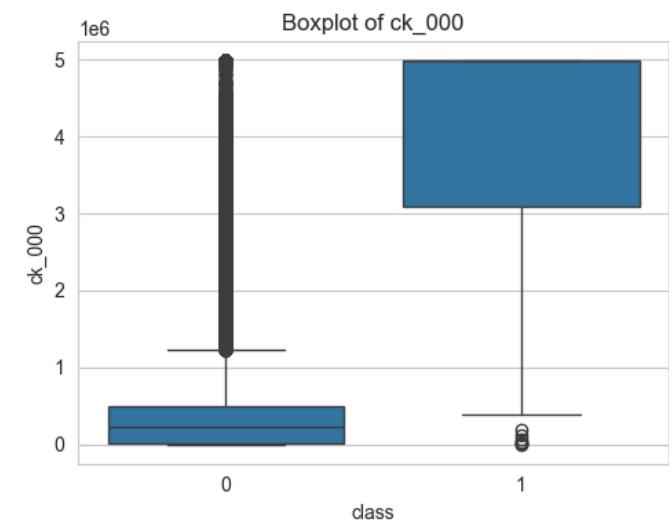
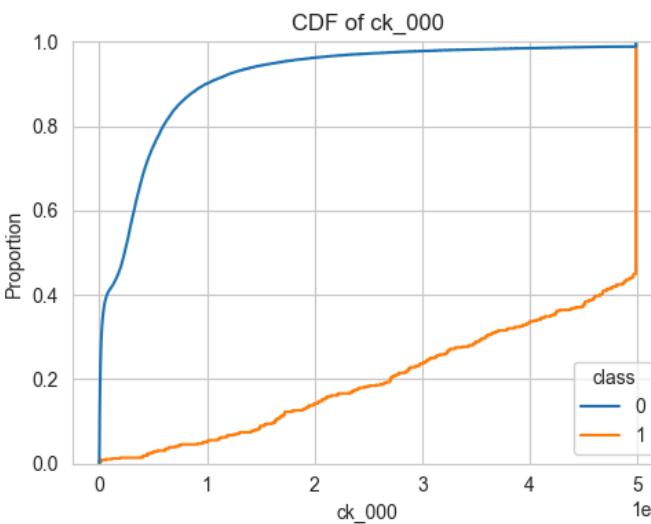
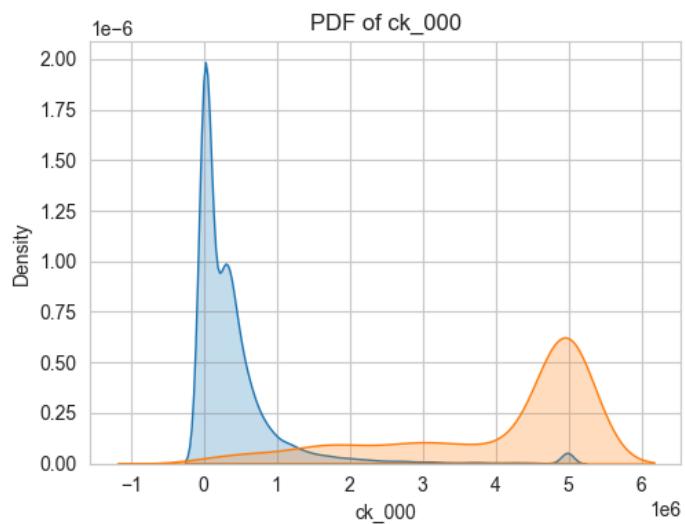
Feature: ck\_000

Class 0 - Mean: 432790.04, Std Dev: 753679.3

Class 1 - Mean: 4012856.6, Std Dev: 1420806.78

2025-02-09 07:47:14,899 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:47:14,968 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



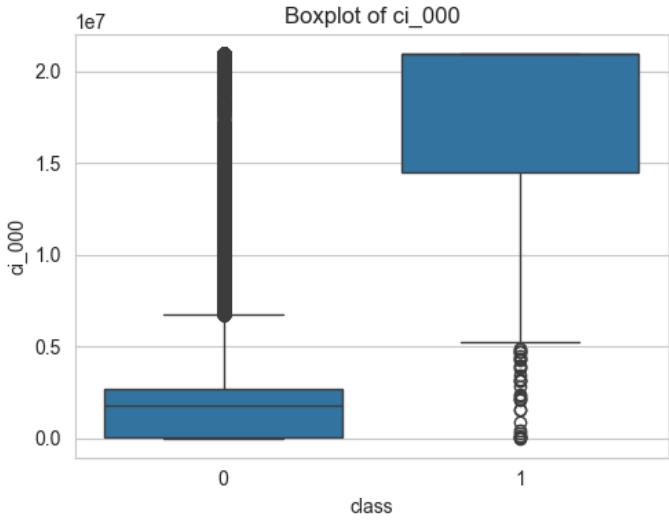
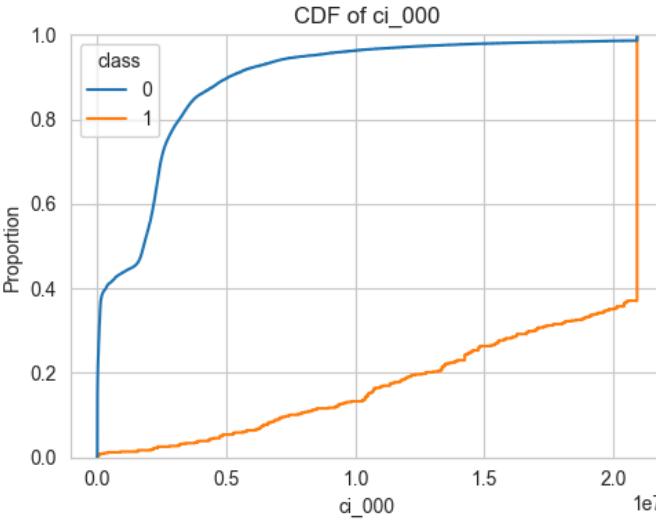
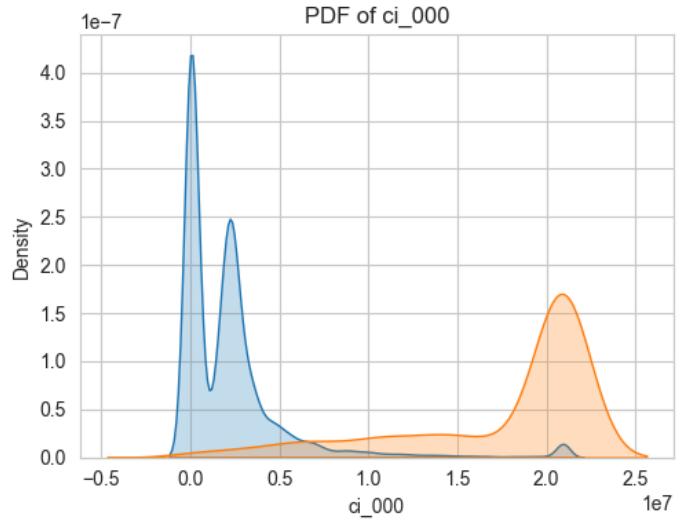
Feature: ci\_000

Class 0 - Mean: 2274906.55, Std Dev: 3445648.11

Class 1 - Mean: 17469096.64, Std Dev: 5608083.73

2025-02-09 07:47:16,903 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:47:16,962 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



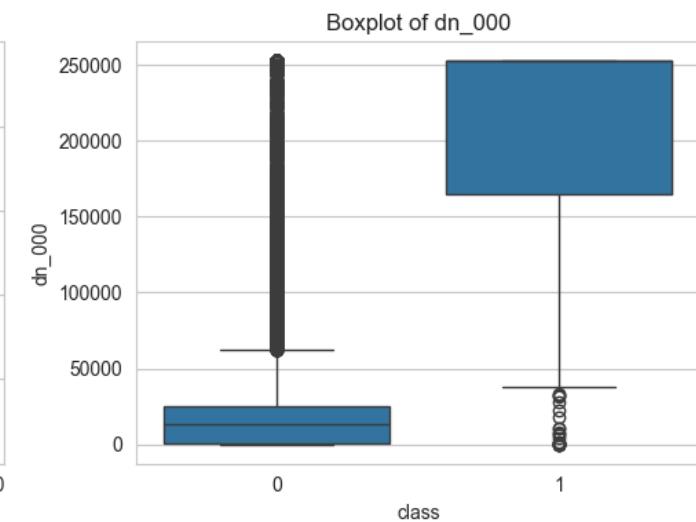
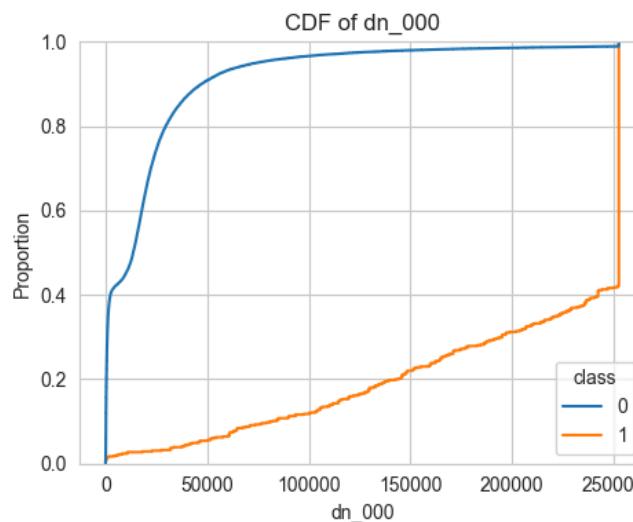
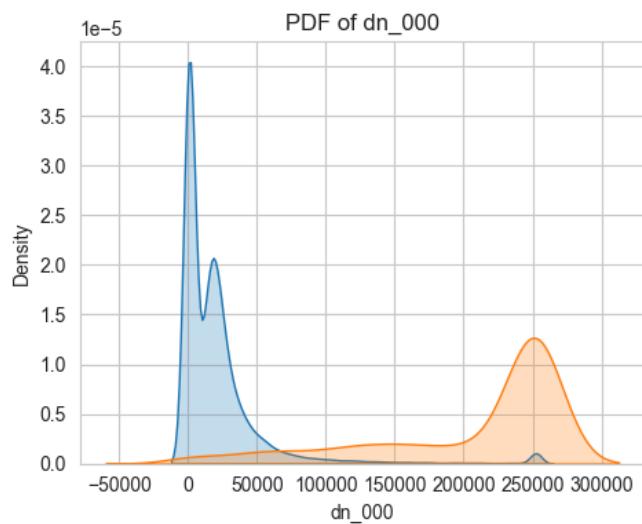
Feature: dn\_000

Class 0 - Mean: 21368.72, Std Dev: 36637.96

Class 1 - Mean: 205633.87, Std Dev: 71920.76

2025-02-09 07:47:19,012 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:47:19,157 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



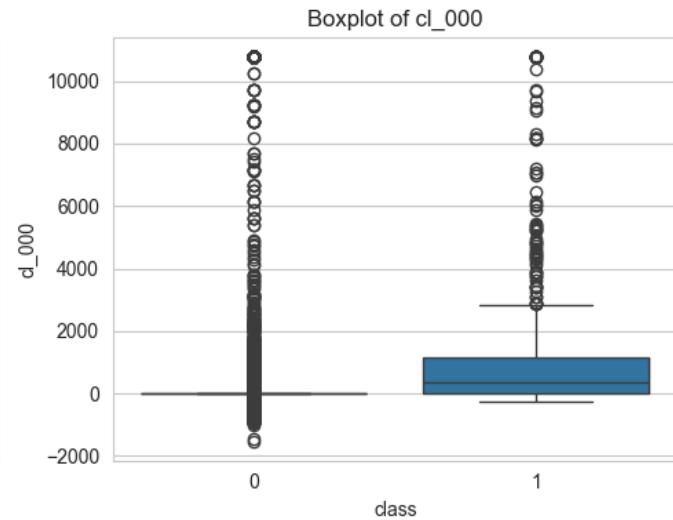
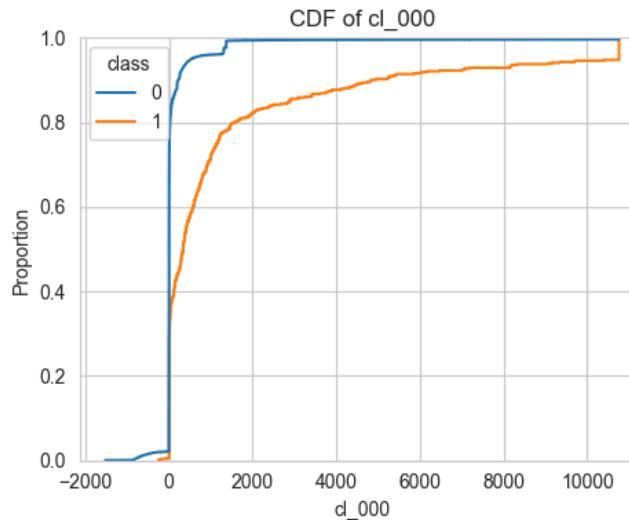
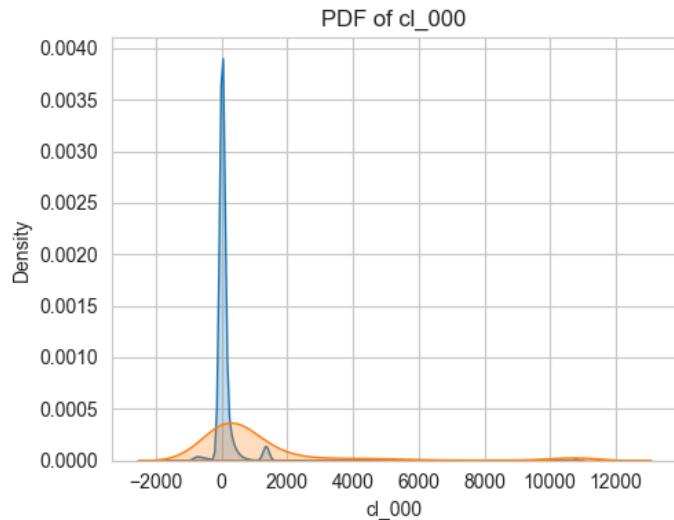
Feature: cl\_000

Class 0 - Mean: 114.26, Std Dev: 725.89

Class 1 - Mean: 1460.21, Std Dev: 2768.6

2025-02-09 07:47:20,977 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:47:21,106 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



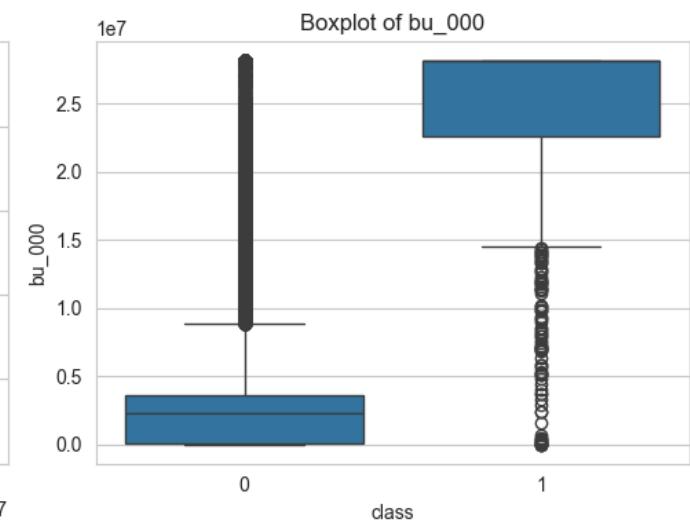
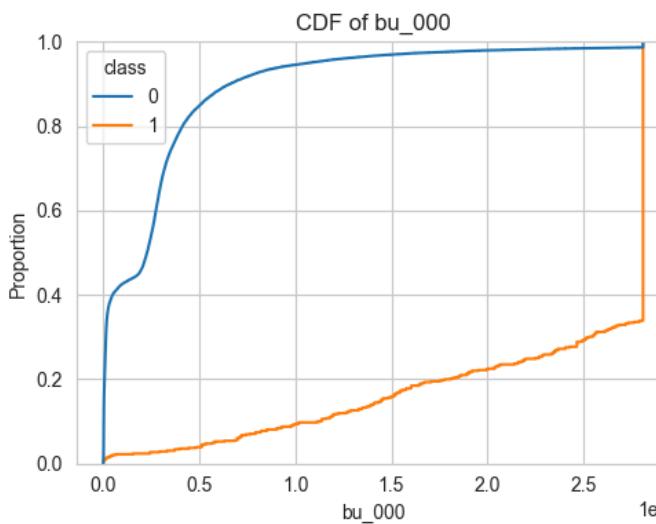
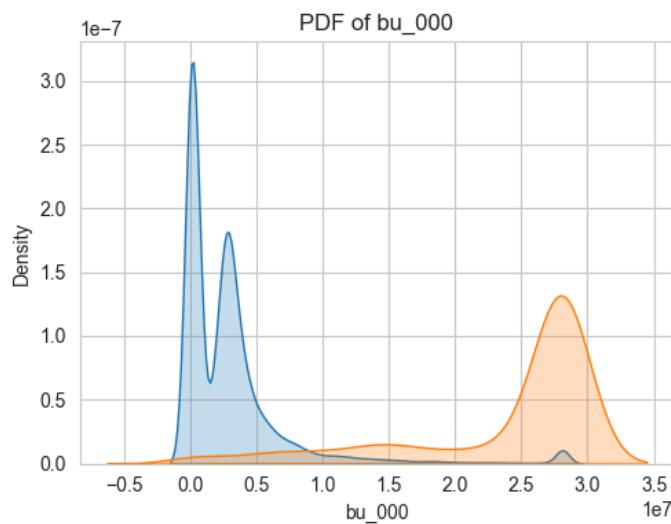
Feature: bu\_000

Class 0 - Mean: 2980716.11, Std Dev: 4552883.09

Class 1 - Mean: 23797967.12, Std Dev: 7598993.61

2025-02-09 07:47:23,535 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:47:23,602 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



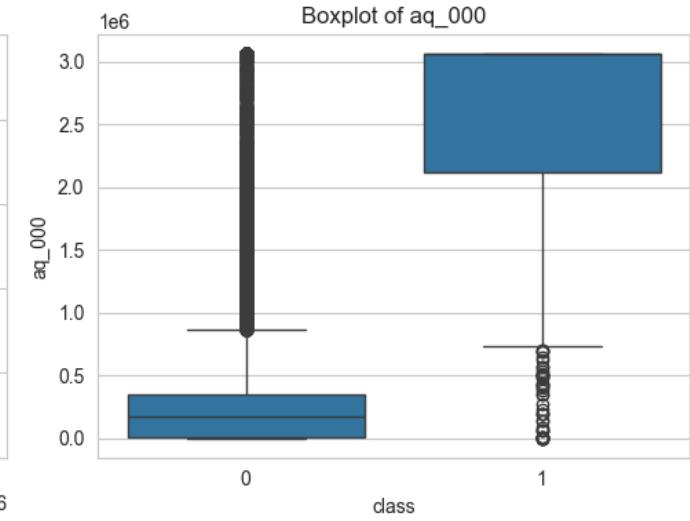
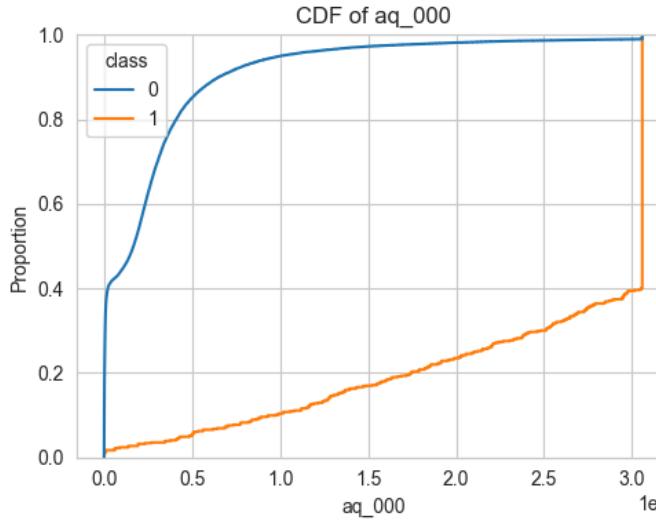
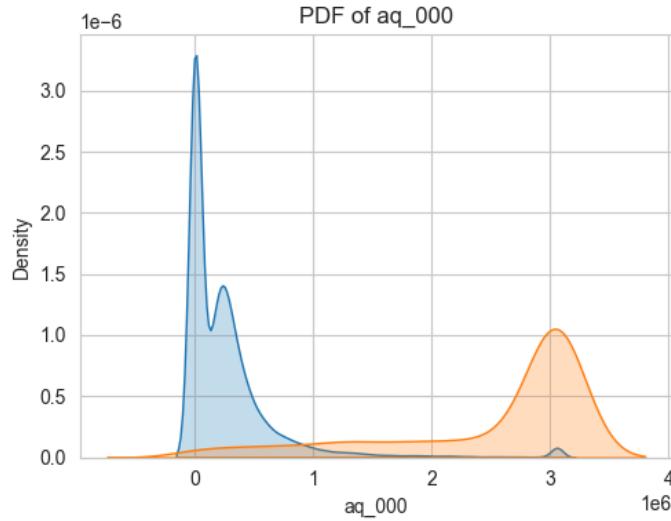
Feature: aq\_000

Class 0 - Mean: 278003.28, Std Dev: 457432.7

Class 1 - Mean: 2501076.83, Std Dev: 885920.91

2025-02-09 07:47:25,804 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.

2025-02-09 07:47:25,870 - INFO - Using categorical units to plot a list of strings that are all parsable as floats or dates. If these strings should be plotted as numbers, cast to the appropriate data type before plotting.



#### Findings from Univariate analysis for no-bin features:

This technical analysis examines multiple system parameters across various subsystems, revealing distinct patterns between normal operations and failure states. The analysis demonstrates significant potential for implementing robust predictive maintenance strategies and early warning systems.

**System Parameter Analysis:** The air pressure system parameter `ap_000` exhibits a notable bimodal distribution with clear separation between operational states. During normal operation, the parameter maintains a mean of 578,980.3 with a standard deviation of 1,021,849.13, characterized by a sharp peak near zero and a smaller secondary peak. Failure states show substantially elevated values with a mean of 5,737,678.6 and increased variability (SD: 2,046,578.92), centered around 7-8 million units. This clear separation provides a strong foundation for threshold-based monitoring.

System performance parameters `bv_000` and `bb_000` demonstrate remarkably similar characteristics, suggesting interconnected system components. Both parameters show bimodal distributions during normal operation with means around 2.98 million and standard deviations of approximately 4.55 million. Failure states for both parameters exhibit dramatically higher values (means ~23.8 million) with increased variability (SD ~7.6 million), indicating these parameters could serve as redundant measurements for validation purposes.

The control parameter `dg_000` presents a unique distribution pattern, with normal operations tightly concentrated near zero (mean: 879.83, SD: 14,186.73) and failure cases showing significantly dispersed values (mean: 44,458.95, SD: 144,313.9). This step-like behavior suggests discrete state transitions rather than continuous degradation.

Dynamic system parameters `dn_000` and `dy_000` show excellent discrimination between operational states. The `dn_000` parameter maintains normal operations around 21,368.72 (SD: 36,637.96) with failure cases averaging 205,633.87 (SD: 71,920.76). Similarly, `dy_000` shows tight control during normal operation (mean: 4,764.13, SD: 21,695.13) with significant elevation during failures (mean: 36,836.22, SD: 63,775.88).

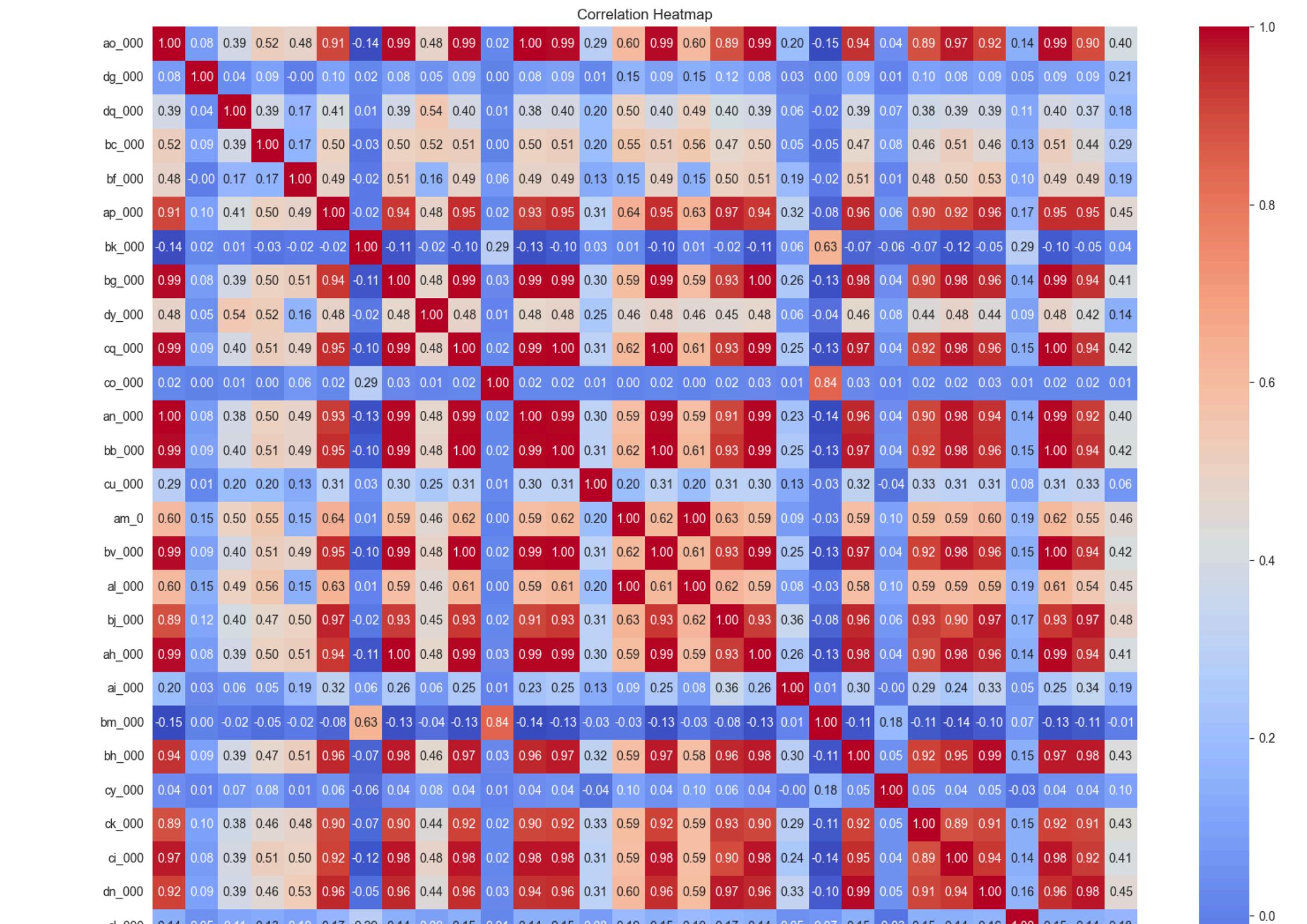
#### Statistical Implications and Operational Recommendations:

The analysis reveals optimal monitoring thresholds for critical parameters. For `ap_000`, values exceeding 3 million units should trigger increased monitoring. The `bv_000` and `bb_000` parameters suggest implementing alerts at  $1.5 \times 10^7$  units. The `dg_000` parameter indicates potential issues above 20,000 units, while `dy_000` shows clear separation at 25,000 units.

The comprehensive statistical analysis supports implementing a multi-parameter monitoring approach. The strong correlation between `bv_000` and `bb_000` provides opportunities for cross-validation, potentially reducing false alarms. The distinct behavior of `dg_000` suggests its use as a discrete state indicator, complementing the continuous monitoring capabilities of other parameters.

This analysis demonstrates that while individual parameters provide valuable insights into system health, the most effective monitoring strategy would involve tracking multiple parameters simultaneously. Particular attention should be paid to sudden changes in `dy_000` and `dn_000` values, which show the clearest separation between normal operation and failure states. The implementation of these findings could significantly enhance predictive maintenance capabilities and reduce unexpected system failures.

```
In [30]: no_bin_selected_features_df = train_data_cleaned[no_bin_selected_features+[ 'class' ]]
no_bin_top_uncorrelated_features = correlation_matrix(data=no_bin_selected_features_df)
```



c_000	0.14	0.05	0.11	0.13	0.10	0.17	0.29	0.14	0.09	0.15	0.01	0.14	0.15	0.06	0.19	0.15	0.19	0.17	0.14	0.05	0.07	0.15	-0.03	0.15	0.14	0.16	1.00	0.15	0.14	0.16
bu_000	0.99	0.09	0.40	0.51	0.49	0.95	-0.10	0.99	0.48	1.00	0.02	0.99	1.00	0.31	0.62	1.00	0.61	0.93	0.99	0.25	-0.13	0.97	0.04	0.92	0.98	0.96	0.15	1.00	0.94	0.42
aq_000	0.90	0.09	0.37	0.44	0.49	0.95	-0.05	0.94	0.42	0.94	0.02	0.92	0.94	0.33	0.55	0.94	0.54	0.97	0.94	0.34	-0.11	0.98	0.04	0.91	0.92	0.98	0.14	0.94	1.00	0.44
class	0.40	0.21	0.18	0.29	0.19	0.45	0.04	0.41	0.14	0.42	0.01	0.40	0.42	0.06	0.46	0.42	0.45	0.48	0.41	0.19	-0.01	0.43	0.10	0.43	0.41	0.45	0.18	0.42	0.44	1.00
ao_000	dg_000	dq_000	bc_000	bf_000	ap_000	bk_000	bg_000	dy_000	cq_000	co_000	an_000	bb_000	cu_000	am_0	bv_000	al_000	bj_000	ah_000	ai_000	bm_000	bh_000	cy_000	dk_000	ci_000	dn_000	cl_000	bu_000	aq_000	class	

2025-02-09 07:47:29,635 - INFO - Top 5 uncorrelated features:

bm\_000 -0.009251

co\_000 0.014057

bk\_000 0.037585

cu\_000 0.063798

cy\_000 0.099388

Name: class, dtype: float64

2025-02-09 07:47:29,638 - INFO - The most uncorrelated feature is: bm\_000

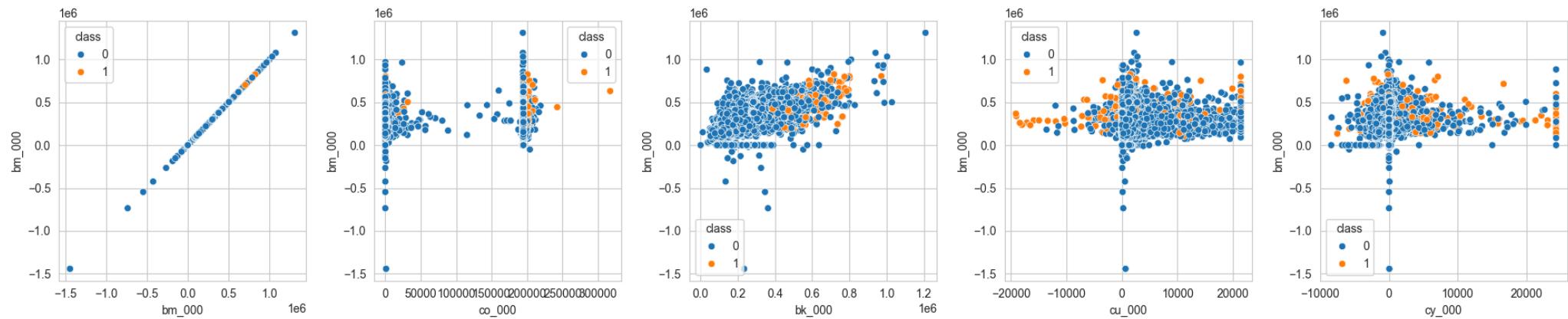
2025-02-09 07:47:29,638 - INFO - The most uncorrelated feature is: bm\_000

#### Correlation Analysis and Feature Relationships:

The correlation heatmap reveals several important patterns in the APS sensor data. Many features show strong positive correlations (red squares) with correlation coefficients above 0.9, indicating potential redundancy in the measurements. However, the analysis has identified five key features with low correlation to the target class variable, which may provide unique predictive value:

- bm\_000 (-0.009251): Shows virtually no linear correlation with failure events
- co\_000 (0.014057): Very weak positive correlation
- bk\_000 (0.037585): Weak positive correlation
- cu\_000 (0.063798): Mild positive correlation
- cy\_000 (0.099388): Strongest correlation among the uncorrelated features, but still weak

```
In [31]: no_bin_top_uncorr_df = no_bin_selected_features_df[no_bin_top_uncorrelated_features + ['class']]
scatter_plot(data=no_bin_top_uncorr_df, feature=no_bin_top_uncorrelated_features[0], percentile=95)
```

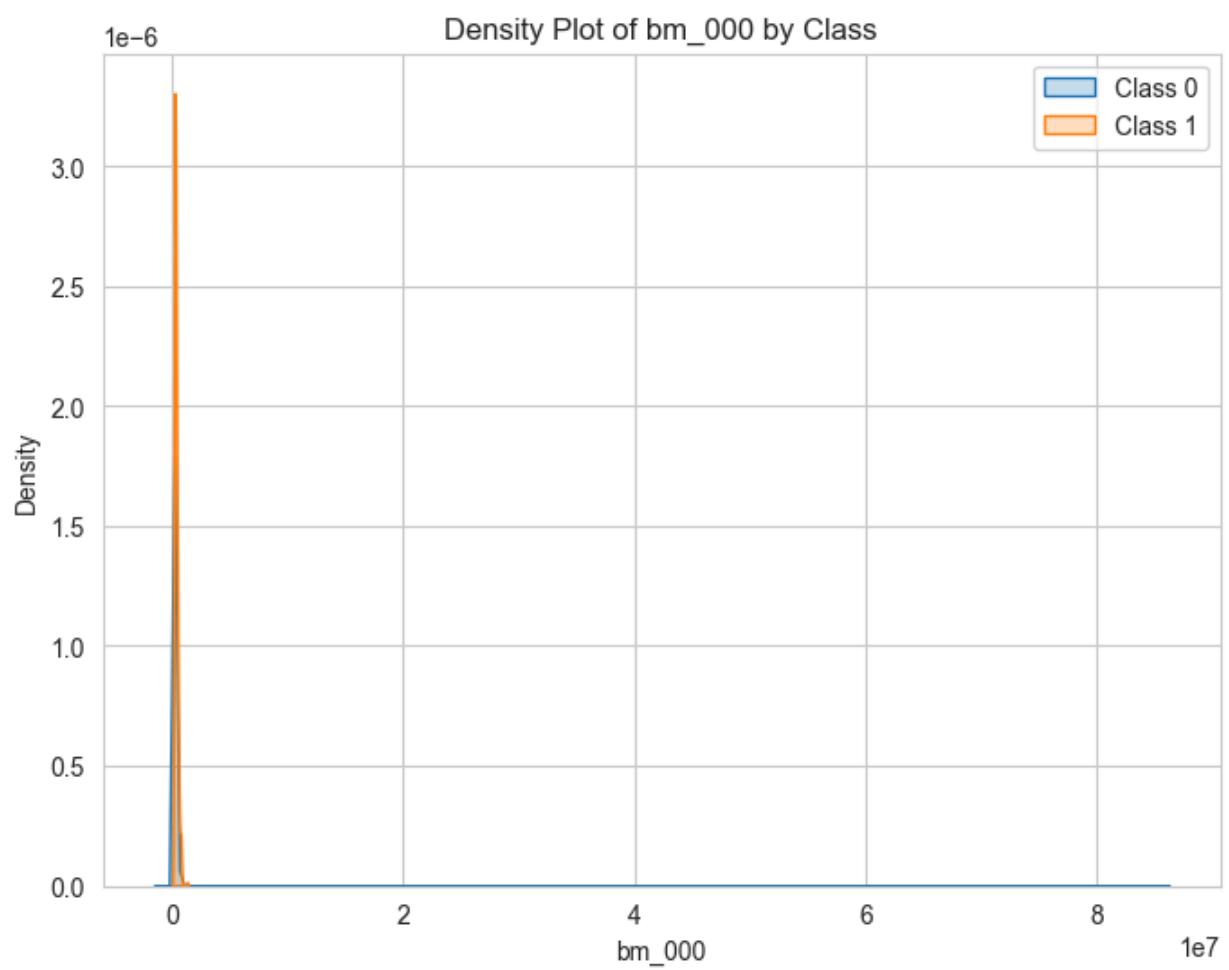


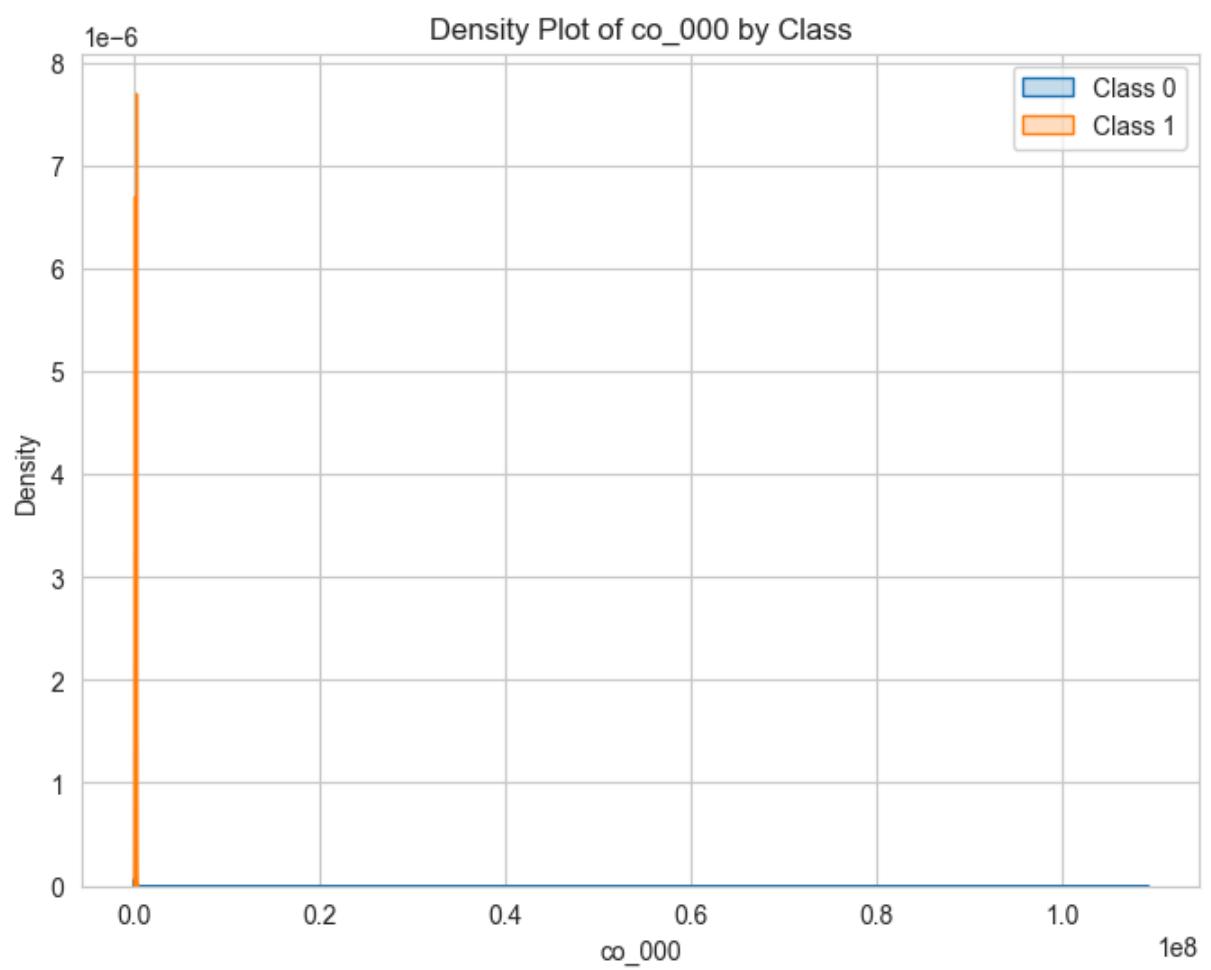
## Interpretation

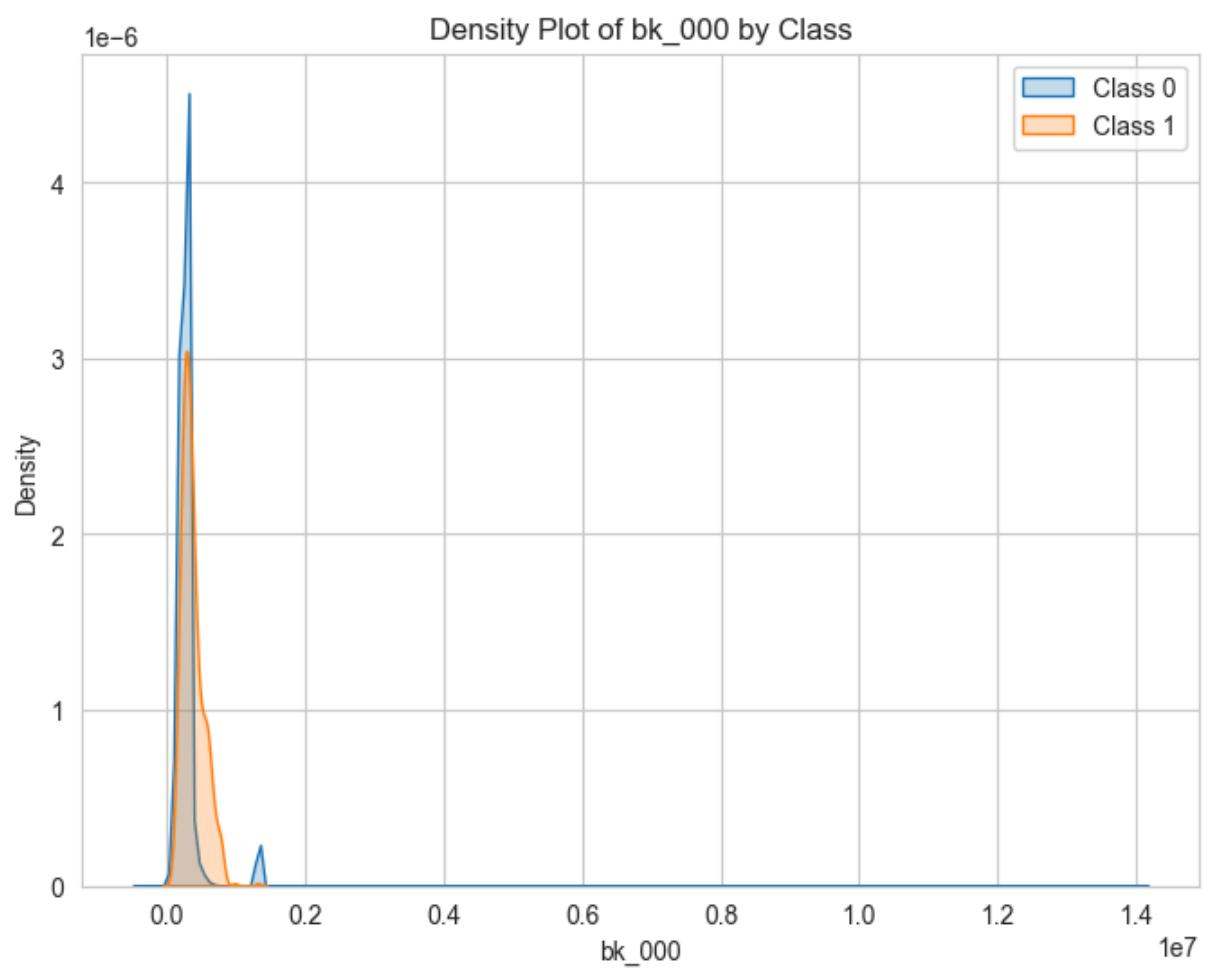
- bm\_000 vs class correlations:** The first plot shows a strong linear relationship in the data points, forming an almost perfect diagonal line. This suggests that this feature has a very consistent and predictable behavior during normal operations, which could make any deviations particularly noticeable for failure detection.
- co\_000 vs class correlations:** The second scatter plot shows more dispersed data points with two main clusters - one near zero and another around the 200,000 mark. There's some overlap between failure (class 1) and non-failure (class 0) cases, but the clustering pattern could be useful for identifying certain types of system states.
- bk\_000 vs class correlations:** This plot shows a wide spread of data points with significant overlap between classes. The distribution appears more random than the previous features, with data points scattered across the range of 0 to 1.2e6. The lack of clear separation between classes suggests this feature alone might not be a strong predictor of failures.
- cu\_000 vs class correlations:** The fourth plot shows an interesting symmetric distribution around zero, with data points spread between -20,000 and +20,000. There appears to be a denser concentration of points near the center, with failure cases (orange) showing slightly more spread than normal operations (blue).
- cy\_000 vs class correlations:** The final plot shows a similar pattern to cu\_000 but with more distinct clustering around zero. The failure cases seem to have a wider spread, particularly in the positive range, which could be useful for detecting anomalous behavior.

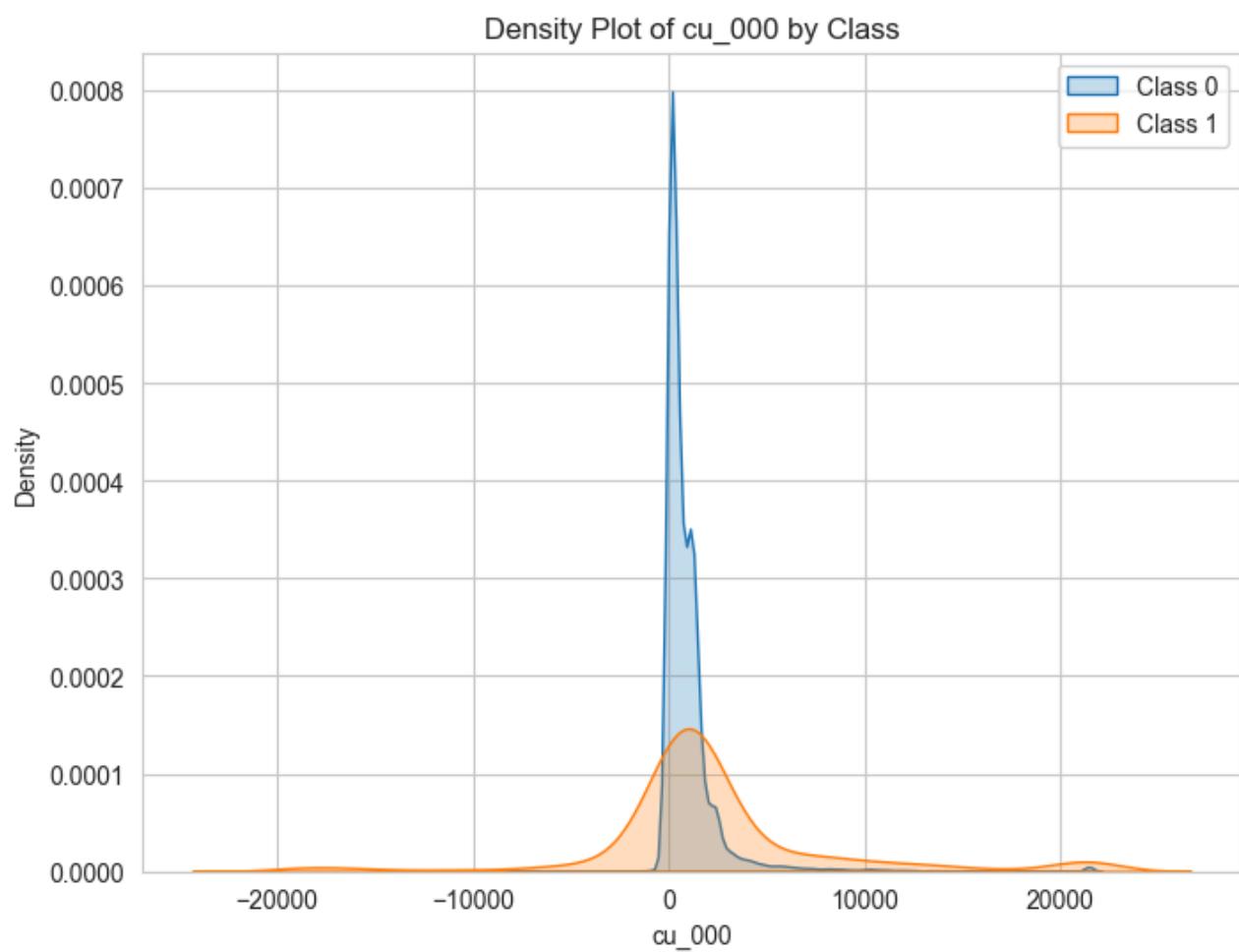
These visualizations suggest that while individual features may not provide clear failure prediction on their own, their combined patterns - particularly sudden deviations from typical clustering patterns - could be valuable indicators for an early warning system. The strong linear relationship in bm\_000 could serve as a particularly useful baseline reference for system stability.

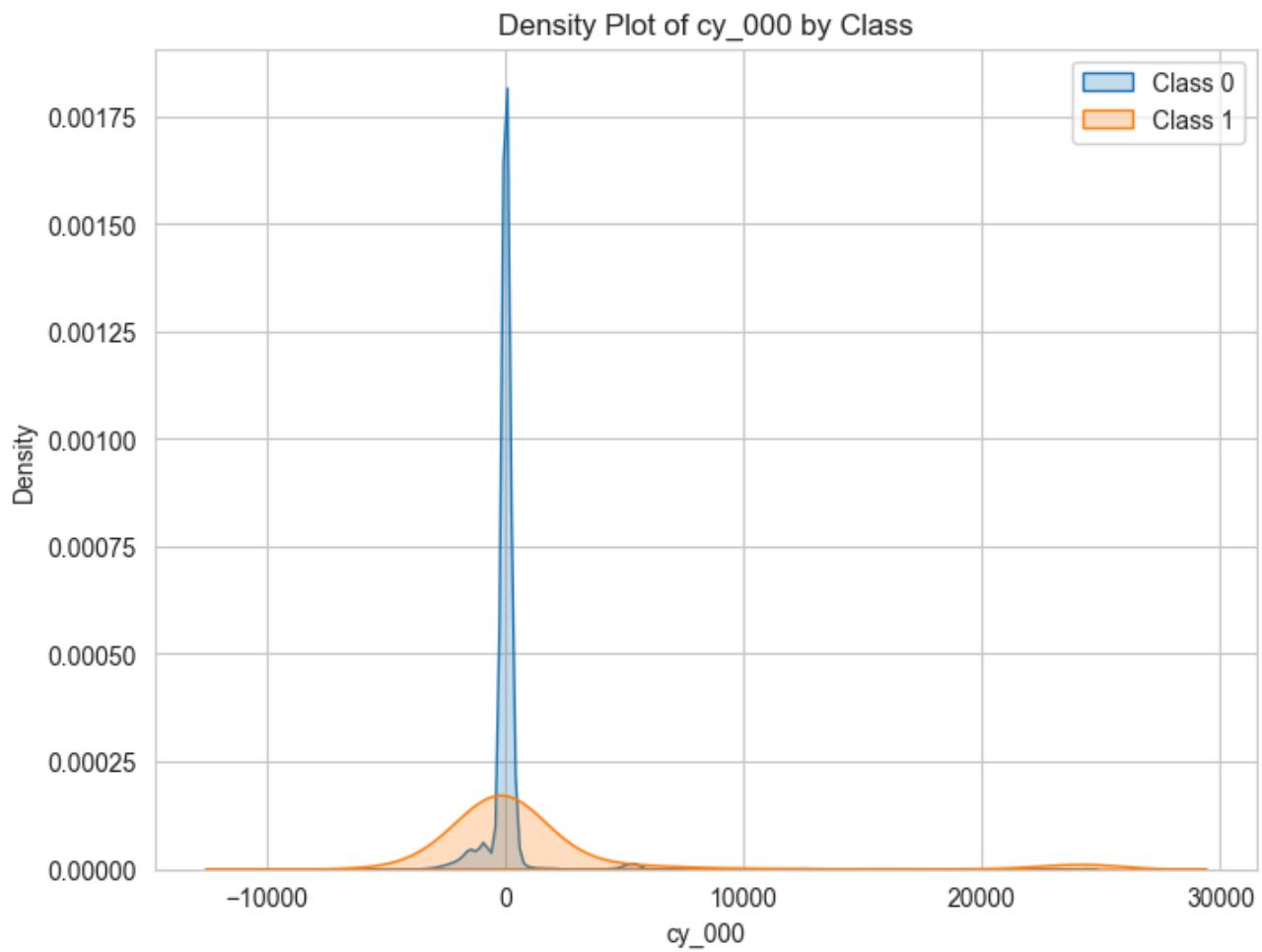
```
In [32]: visualize_class_distributions(train_data_cleaned[no_bin_top_uncorrelated_features + ['class']], no_bin_top_uncorrelated_features)
```











### **Analysis:**

- **bm\_000:** The density plot shows an extremely concentrated distribution near zero for both classes, with almost complete overlap. The scatter plots reveal a linear pattern in the data points, suggesting this feature might be derived from or closely related to other measurements. While it has the lowest correlation with failures, its consistent behavior could serve as a baseline reference for system operation.
- **co\_000:** The density distribution demonstrates a sharp peak near zero with minimal separation between classes. The scatter plots show some clustering patterns, particularly for class 0 (normal operation). This feature's weak correlation with failures but distinct clustering behavior suggests it might be useful in combination with other features.
- **bk\_000:** Shows a bimodal distribution for class 0 with a secondary peak around 0.2e7, while class 1 maintains a single peak near zero. The CDF reveals a clear separation pattern after the initial concentration, indicating potential value for detecting certain types of anomalies.
- **cu\_000:** Exhibits a broader distribution for class 1 compared to class 0, with higher variance in the failure cases. The pronounced peak in class 0 could serve as a reliable baseline for normal operation, with deviations potentially indicating developing issues.
- **cy\_000:** Shows the most distinct separation among the uncorrelated features, with class 1 displaying a wider spread and multiple modes. This suggests it captures unique aspects of system behavior not reflected in other measurements. Operational Implications and Early Warning System Design:

- **Primary Monitoring Features:** The `c1_000` and `ai_000` features show significant differences between normal and failure states, with notably higher means and standard deviations in failure cases. These could serve as primary monitoring indicators.
- **Secondary Validation Features:** The uncorrelated features (particularly `cy_000` and `cu_000`) could serve as independent verification metrics, helping to reduce false positives.

#### Statistical Thresholds:

- For `c1_000`: Alert threshold around 840 (normal mean + 1 std)
- For `ai_000`: Warning level at 25,828 (normal mean + 1 std)

#### Combined Monitoring Strategy:

##### Implement a weighted scoring system that:

Uses the highly correlated features for primary failure detection  
 Validates alerts using the uncorrelated features to confirm anomalies  
 Monitors trend changes in the `bm_000` and `co_000` features as baseline stability indicators

This analysis suggests a multi-layered monitoring approach would be most effective, using both the strongly correlated and uncorrelated features to provide robust failure prediction while minimizing false alarms.

#### Step 4: Feature Engineering

```
In [33]: from imblearn.over_sampling import SMOTE
from imblearn.under_sampling import RandomUnderSampler
from imblearn.pipeline import Pipeline

def balance_classes(df, label, over_strategy=0.3, under_strategy=0.5):
    """
    Balances the class distribution in the dataset using SMOTE and Random Under-Sampling.

    Parameters:
    - df (pd.DataFrame): The feature DataFrame.
    - label (pd.Series): The target labels.
    - over_strategy (float): The desired ratio of the minority class after SMOTE.
    - under_strategy (float): The desired ratio of the majority class after under-sampling.

    Returns:
    - pd.DataFrame: The balanced feature DataFrame.
    - pd.Series: The balanced target labels.
    """
    # Define the over-sampling and under-sampling strategies
    over = SMOTE(sampling_strategy=over_strategy)
    under = RandomUnderSampler(sampling_strategy=under_strategy)

    # Create a pipeline with the defined steps
    steps = [('o', over), ('u', under)]
    pipeline = Pipeline(steps=steps) # Use imblearn's Pipeline
```

```

# Fit and resample the data
df_balanced, label_balanced = pipeline.fit_resample(df, label)

return df_balanced, label_balanced

# Convert the target column ('class') to binary (0 for 'neg', 1 for 'pos')
test_data_cleaned['class'] = test_data_cleaned['class'].map({'neg':0, 'pos':1})

# Separate features and target variable
X_train_balanced, y_train_balanced = balance_classes(train_data_cleaned.drop(columns=['class']),
                                                       train_data_cleaned['class'])

# Assuming you have your DataFrame `test_cleaned_data` and target `y_test`
X_test_balanced, y_test_balanced = balance_classes(test_data_cleaned.drop(columns=['class']),
                                                       test_data_cleaned['class'])

```

```
In [34]: print(f"X_train_balanced Shape: {X_train_balanced.shape}")
print(f"X_test_balanced Shape: {X_test_balanced.shape}")
y_train_balanced.value_counts()
```

```
X_train_balanced Shape: (49806, 162)
X_test_balanced Shape: (13200, 162)
```

```
Out[34]: class
0    33204
1    16602
Name: count, dtype: int64
```

```

In [35]: import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.pipeline import Pipeline

# Create a pipeline with StandardScaler and PCA
pipeline = Pipeline([
    ('scaler', StandardScaler()), # Step 1: Standardize the data
    ('pca', PCA(n_components=2)) # Step 2: Apply PCA to reduce to 2 components
])

# Fit the pipeline on the training data and transform it
X_train_pca = pipeline.fit_transform(X_train_balanced)
logger.info(f"Training X Set: {X_train_pca.shape}, Y set: {y_train_balanced.shape}")
# Transform the test data using the same pipeline (without fitting)
X_test_pca = pipeline.transform(X_test_balanced)

# Create a DataFrame for the PCA results for visualization
pca_df = pd.DataFrame(data=X_train_pca, columns=['Principal Component 1', 'Principal Component 2'])
pca_df['Class'] = y_train_balanced.values # Add the target variable for coloring

# Plotting the PCA results

```

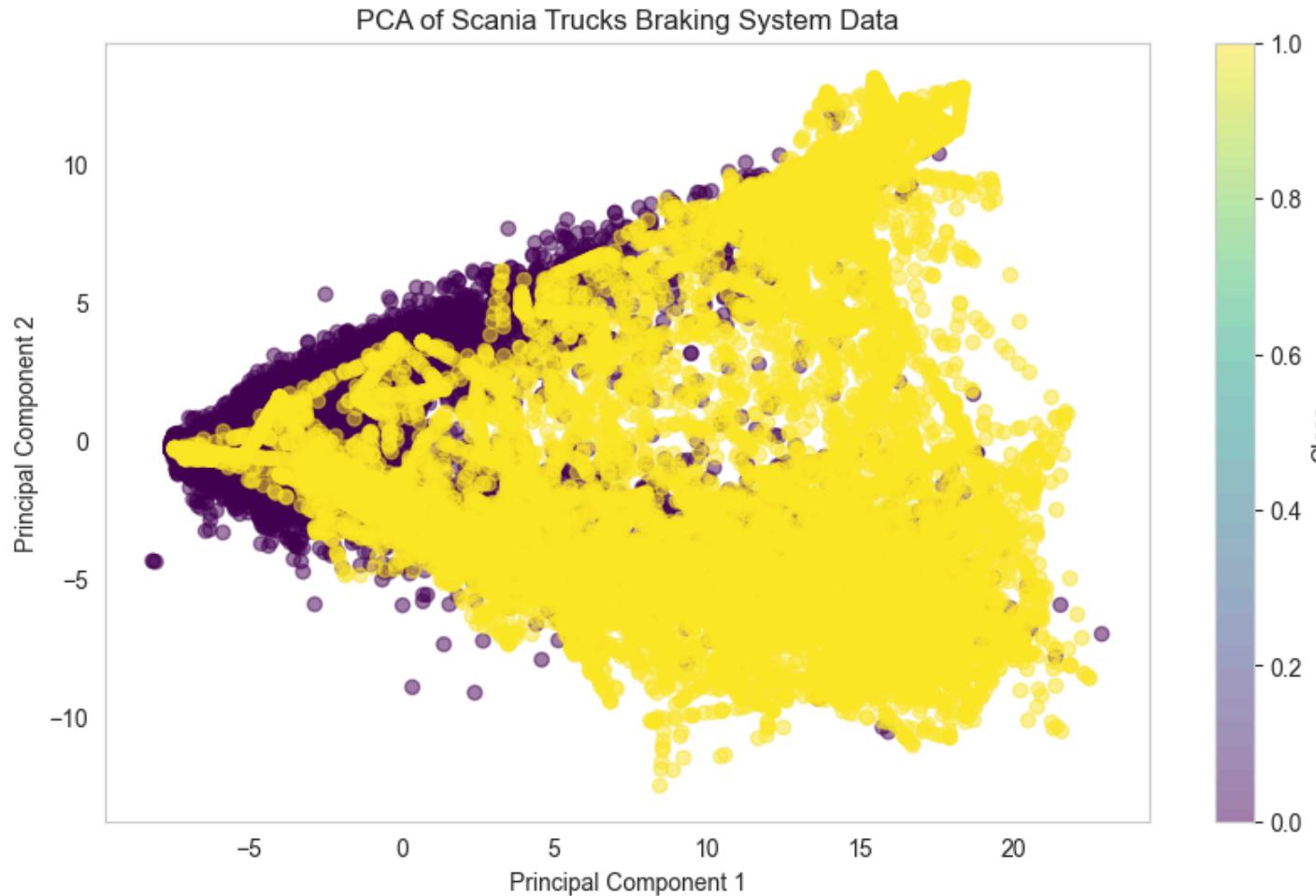
```

plt.figure(figsize=(10, 6))
scatter = plt.scatter(pca_df['Principal Component 1'], pca_df['Principal Component 2'],
                     c=pca_df['Class'], alpha=0.5, cmap='viridis')

plt.title('PCA of Scania Trucks Braking System Data')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.colorbar(scatter, label='Class')
plt.grid()
plt.show()

```

2025-02-09 07:47:51,735 - INFO - Training X Set: (49806, 2), Y set: (49806,)



## Step 5: Model Development

In [36]:

```

import pandas as pd
from skopt import BayesSearchCV
from sklearn.model_selection import KFold

```

```

def hyper_parameter_tuning(model: Pipeline, X: pd.DataFrame, y: pd.Series, param_space: dict,
                           n_iter: int = 50, cv: int = 5, verbose: int=10):
    """
    Perform Bayesian hyperparameter tuning for a given model using cross-validation.

    Parameters:
    - model (Pipeline): The machine learning pipeline model to tune.
    - X (pd.DataFrame): The feature DataFrame.
    - y (pd.Series): The target labels.
    - param_space (dict): The parameter space for tuning.
    - n_iter (int): Number of iterations for the search (default is 50).
    - cv (int): Number of cross-validation folds (default is 5).

    Returns:
    - best_params (dict): The best parameters found during tuning.
    - best_score (float): The best score achieved during tuning.
    """
    # Create a KFold object for cross-validation
    kf = KFold(n_splits=cv, shuffle=True, random_state=42)

    # Initialize BayesSearchCV with the model and parameter space
    clf = BayesSearchCV(
        model,
        param_space,
        n_iter=n_iter,
        cv=kf,
        n_jobs=-1, # Use all available cores
        random_state=42,
        verbose=verbose
    )

    # Fit the model to the data
    clf.fit(X, y)

    return clf.best_params_, clf.best_score_

```

```

In [37]: from sklearn.metrics import (accuracy_score,
                                      classification_report,
                                      confusion_matrix,
                                      precision_recall_curve,
                                      f1_score
                                      )

def evaluate_model(model, params: Dict[str, Tuple[pd.DataFrame, pd.Series]], plot:bool=True, verbose:int=0):
    """
    Evaluates the model on multiple datasets provided in a dictionary.

    Parameters:
    params (dict): A dictionary where keys are dataset names (e.g., "Training", "Testing", "Validation")
    """

```

and values are tuples of (X, y, y\_pred).

Returns:  
 dict: A dictionary containing accuracy scores and classification reports for each dataset.

```

"""
results = {}

for dataset, (X, y) in params.items():
    # Get prediction
    y_pred = model.predict(X)
    # Evaluate accuracy and classification report
    accuracy = accuracy_score(y, y_pred)
    report = classification_report(y, y_pred)
    matrix = confusion_matrix(y, y_pred)
    f1_score_ = f1_score(y, y_pred, average='macro')
    recall = precision_recall_curve(y, y_pred)

    # Log results
    if verbose > 0:
        logger.info(f'{dataset} Accuracy: {accuracy:.2f}\n')
        logger.info(f'{dataset} Confusion matrix:\n{matrix}')
        logger.info(f'{dataset} Classification Report:\n{report}')

    if plot:
        group_names = ["TN", "FP", "FN", "TP"]
        group_counts = [f'{value}' for value in matrix.flatten()]
        labels = [f'{v1}\n{v2}' for v1, v2 in zip(group_names, group_counts)]
        labels = np.asarray(labels).reshape(2,2)

        sns.heatmap(matrix, annot=labels, fmt=' ', cmap='Blues')
        plt.show()

    # Store results in dictionary
    results[f'{dataset}_accuracy'] = accuracy
    results[f'{dataset}_report'] = report
    results[f'{dataset}_matrix'] = matrix
    results[f'{dataset}_f1_score'] = f1_score_
    results[f'{dataset}_recall'] = recall

return results

```

# Dictionary with datasets and predictions

```

data_dict = {
    "Training": (X_train_pca, y_train_balanced),
    "Testing": (X_test_pca, y_test_balanced),
}

```

In [38]: `from sklearn.dummy import DummyClassifier`

```

dummy_model = DummyClassifier(strategy='most_frequent', constant=0)
dummy_model.fit(X_train_pca, y_train_balanced)

```

```
reports = evaluate_model(dummy_model, data_dict, verbose=1)
```

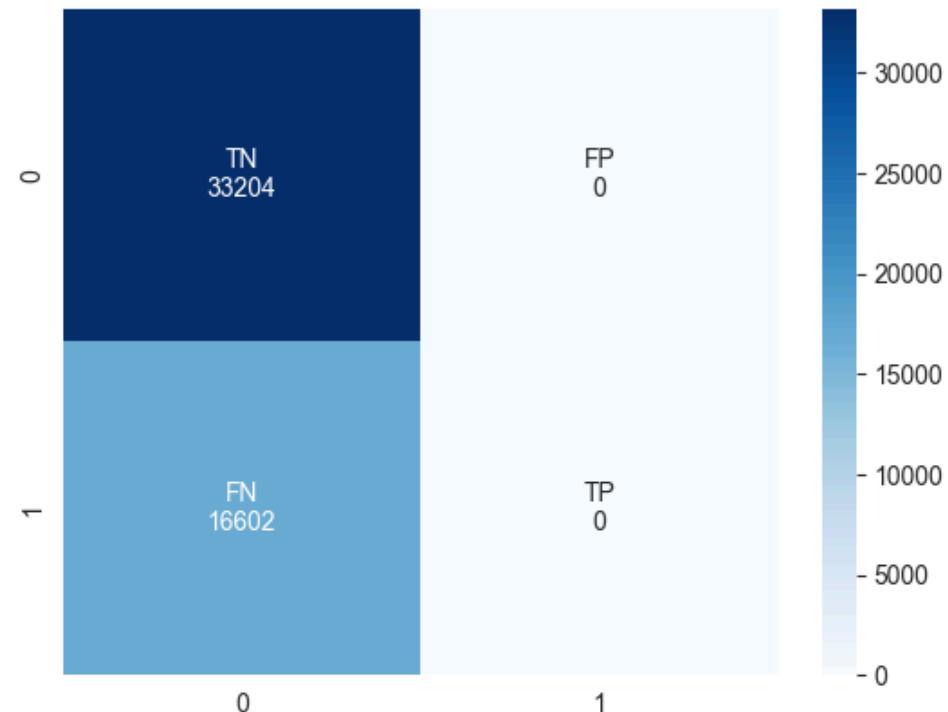
```
2025-02-09 07:47:55,072 - INFO - Training Accuracy: 0.67
```

```
2025-02-09 07:47:55,073 - INFO - Training Confusion matrix:
```

```
[[33204  0]
 [16602  0]]
```

```
2025-02-09 07:47:55,074 - INFO - Training Classification Report:
```

	precision	recall	f1-score	support
0	0.67	1.00	0.80	33204
1	0.00	0.00	0.00	16602
accuracy			0.67	49806
macro avg	0.33	0.50	0.40	49806
weighted avg	0.44	0.67	0.53	49806



2025-02-09 07:47:55,462 - INFO - Testing Accuracy: 0.67

2025-02-09 07:47:55,465 - INFO - Testing Confusion matrix:

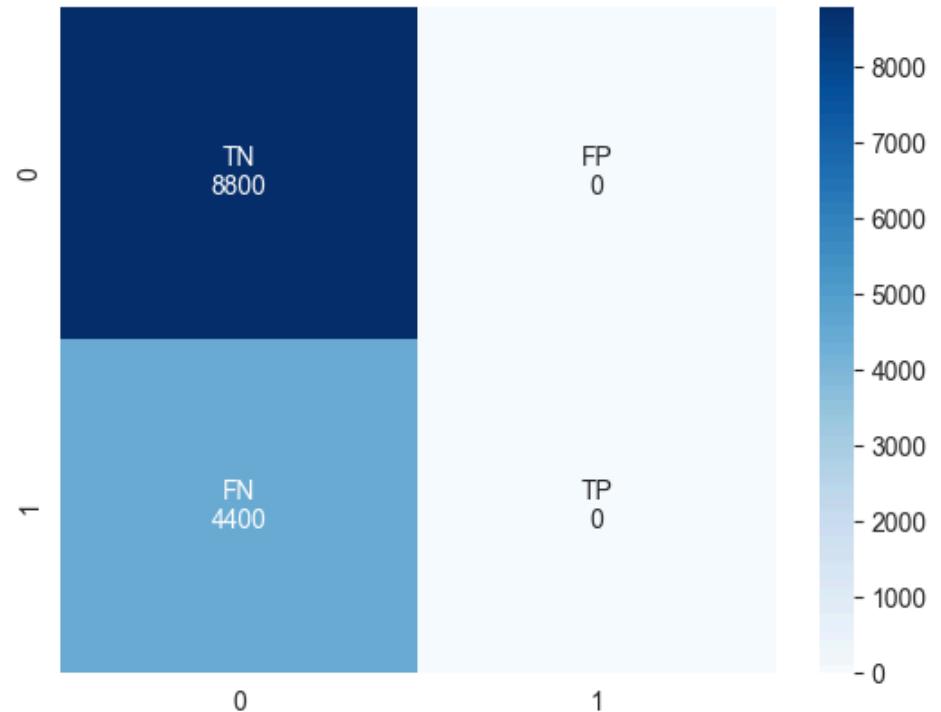
$$\begin{bmatrix} 8800 & 0 \\ 4400 & 0 \end{bmatrix}$$

2025-02-09 07:47:55,467 - INFO - Testing Classification Report:

precision recall f1-score support

0	0.67	1.00	0.80	8800
1	0.00	0.00	0.00	4400

accuracy			0.67	13200
macro avg	0.33	0.50	0.40	13200
weighted avg	0.44	0.67	0.53	13200



In [39]:

```

param_space: Dict[str, Any],
filename: str) -> Tuple[Callable, dict]:
"""

Train and evaluate a machine learning model, perform hyperparameter tuning if required,
and save the model if it performs well.

Parameters:
- model_func (Callable): Model constructor function (e.g., RandomForestClassifier).
- datasets (Dict[str, Tuple[Any, Any]]): Dictionary containing dataset splits
  (e.g., {'Training': (X_train, y_train), 'Validation': (X_val, y_val)}).
- param_space (Dict[str, Any]): Parameter space for hyperparameter tuning.
- filename (str): Path to save or load the trained model.

Returns:
- model (Callable): The trained model.
- report (dict): Evaluation report containing accuracy and other metrics.
"""

# Attempt to Load an existing model
try:
    model = joblib.load(filename)
    logger.info(f"Loaded existing model from {filename}.")
except (FileNotFoundException, EOFError, OSError):
    logger.info("No valid existing model found. Starting training.")

# Instantiate the model
try:
    model = model_func(random_state=42)
except TypeError:
    model = model_func()

# Perform hyperparameter tuning
best_params, best_score = hyper_parameter_tuning(
    model, datasets['Training'][0], datasets['Training'][1], param_space
)

logger.info(f"Best Parameters: {best_params}")
logger.info(f"Best Score: {best_score:.4f}")

# Train the model with the best parameters
try:
    model = model_func(random_state=42, **best_params)
except TypeError:
    model = model_func(**best_params)

model.fit(datasets['Training'][0], datasets['Training'][1])

# Evaluate the model
report = evaluate_model(model, datasets, verbose=1)
logger.info(f"Model evaluation report: {report}")

```

```

# Save the model if it meets the accuracy threshold
if report.get('Testing_accuracy', 0) > 0.80:
    os.makedirs(os.path.dirname(filename), exist_ok=True)
    joblib.dump(model, filename)
    logger.info(f"Model saved to {filename}.")

return model, report

```

In [40]:

```

from sklearn.ensemble import RandomForestClassifier
rdf_param_space = {
    'n_estimators': (50, 100, 200, 300), # Number of trees in the forest
    'max_depth': (None, 10, 20, 30), # Maximum depth of the tree
    'min_samples_split': (2, 5, 10), # Minimum number of samples required to split an internal node
    'min_samples_leaf': (1, 2, 4), # Minimum number of samples required to be at a leaf node
}

rdf_model_func = RandomForestClassifier

filename = "trained_models/rdf_model.sav"

rdf_trained_model, rdf_evaluation_report = train_and_evaluate_model(rdf_model_func, data_dict, rdf_param_space, filename)

```

2025-02-09 07:47:57,055 - INFO - Loaded existing model from trained\_models/rdf\_model.sav.

2025-02-09 07:48:00,906 - INFO - Training Accuracy: 0.94

2025-02-09 07:48:00,907 - INFO - Training Confusion matrix:

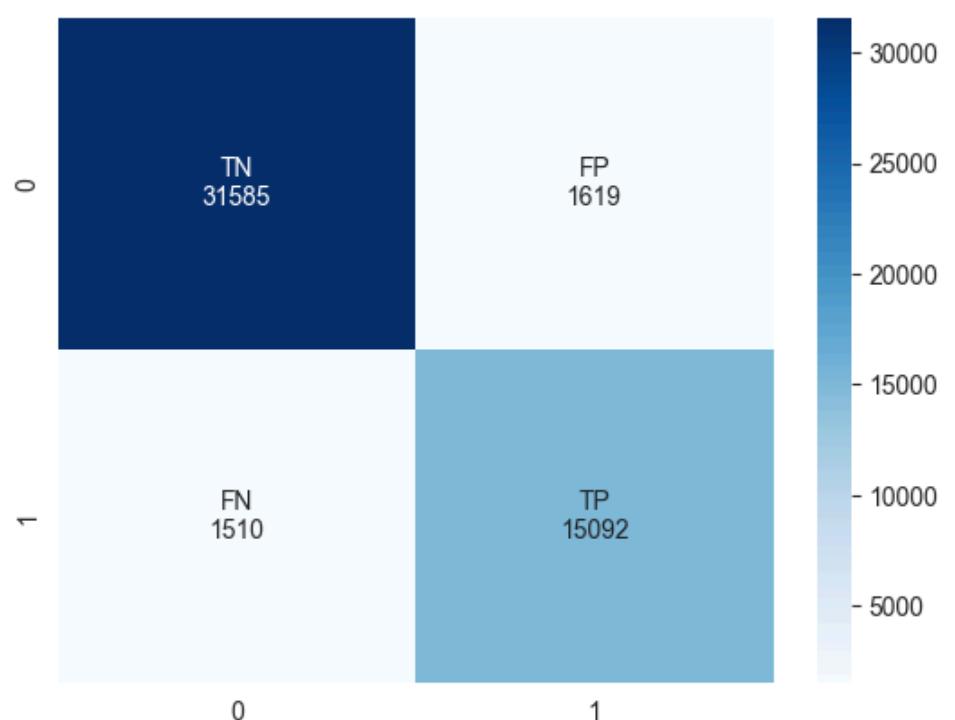
```

[[31585 1619]
 [ 1510 15092]]

```

2025-02-09 07:48:00,909 - INFO - Training Classification Report:

	precision	recall	f1-score	support
0	0.95	0.95	0.95	33204
1	0.90	0.91	0.91	16602
accuracy			0.94	49806
macro avg	0.93	0.93	0.93	49806
weighted avg	0.94	0.94	0.94	49806



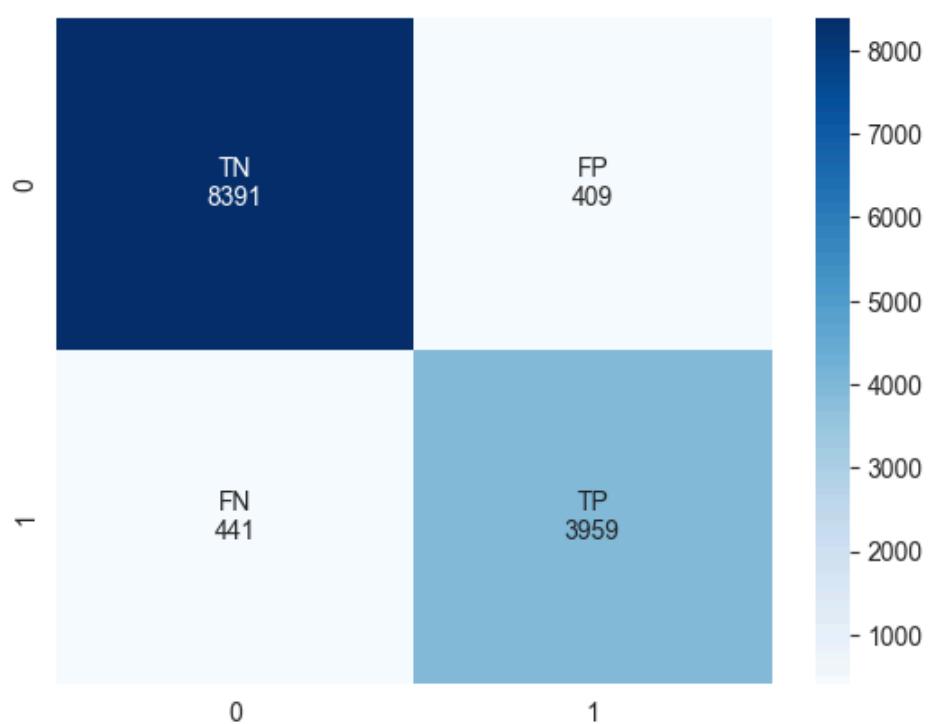
```
2025-02-09 07:48:01,783 - INFO - Testing Accuracy: 0.94
```

```
2025-02-09 07:48:01,784 - INFO - Testing Confusion matrix:
```

```
[[8391 409]
 [ 441 3959]]
```

```
2025-02-09 07:48:01,784 - INFO - Testing Classification Report:
```

	precision	recall	f1-score	support
0	0.95	0.95	0.95	8800
1	0.91	0.90	0.90	4400
accuracy			0.94	13200
macro avg	0.93	0.93	0.93	13200
weighted avg	0.94	0.94	0.94	13200



```
2025-02-09 07:48:01,960 - INFO - Model evaluation report: {'Training_accuracy': 0.937176243826045, 'Training_report': 'precision      recall f1-score support\n          0      0.95      0.95      0.95    33204\n          1      0.90      0.91      0.91    16602\nmacro avg      0.93      0.93      0.93    49806\nweighted avg      0.94      0.94      0.94    49806', 'Training_matrix': array([[31585, 1619],\n[ 1510, 15092]]), 'Training_f1_score': 0.9294387111847664, 'Training_recall': (array([0.33333333, 0.90311771, 1.        ]), array([1.        , 0.9090471, 0.        ]), array([0, 1])), 'Testing_accuracy': 0.9356060606060606, 'Testing_report': 'precision      recall f1-score support\n          0      0.95      0.95      0.95    8800\n          1      0.91      0.90      0.90    4400\nmacro avg      0.93      0.93      0.93    13200\nweighted avg      0.94      0.94      0.94    13200', 'Testing_matrix': array([[8391, 409],\n[ 441, 3959]]), 'Testing_f1_score': 0.9274243826751626, 'Testing_recall': (array([0.33333333, 0.90636447, 1.        ]), array([1.        , 0.89977273, 0.        ]), array([0, 1]))}
2025-02-09 07:48:04,940 - INFO - Model saved to trained_models/rdf_model.sav.
```

```
In [41]: from sklearn.ensemble import AdaBoostClassifier

ada_param_space = {
    "n_estimators": (50, 100, 200, 300, 500), # Number of weak Learners in the ensemble
    "learning_rate": (0.01, 0.05, 0.1, 0.5, 1.0), # Weight applied to each classifier
}

ada_model_func = AdaBoostClassifier
filename = "trained_models/ada_model.sav"

ada_trained_model, ada_evaluation_report = train_and_evaluate_model(ada_model_func, data_dict, ada_param_space, filename)
```

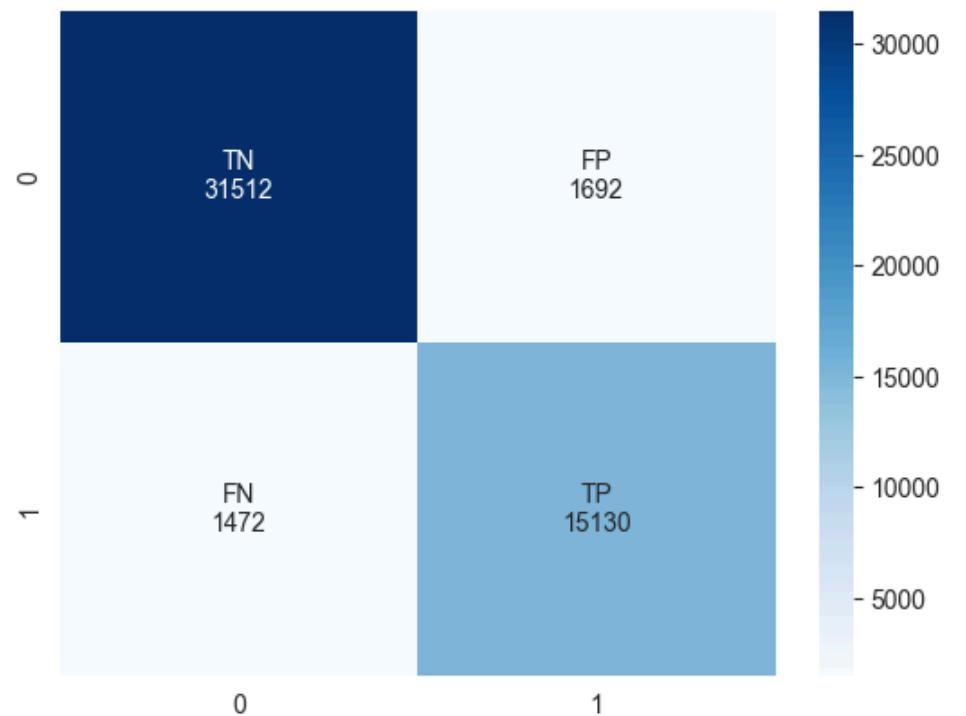
```
2025-02-09 07:48:06,041 - INFO - Loaded existing model from trained_models/ada_model.sav.  
2025-02-09 07:48:08,253 - INFO - Training Accuracy: 0.94
```

```
2025-02-09 07:48:08,254 - INFO - Training Confusion matrix:
```

```
[[31512 1692]  
 [ 1472 15130]]
```

```
2025-02-09 07:48:08,255 - INFO - Training Classification Report:
```

	precision	recall	f1-score	support
0	0.96	0.95	0.95	33204
1	0.90	0.91	0.91	16602
accuracy			0.94	49806
macro avg	0.93	0.93	0.93	49806
weighted avg	0.94	0.94	0.94	49806



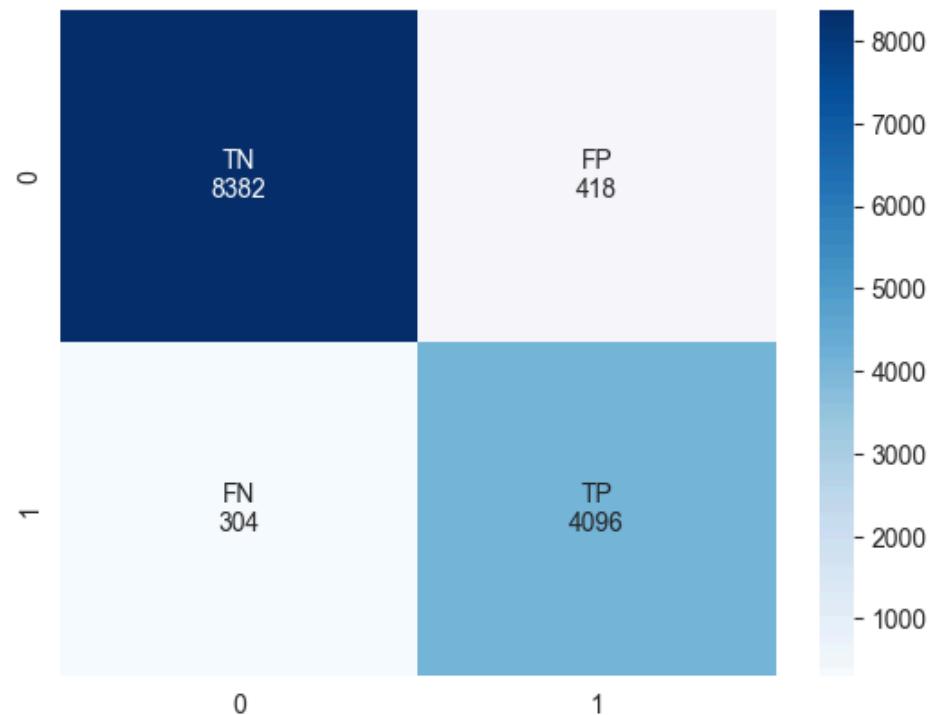
2025-02-09 07:48:09,045 - INFO - Testing Accuracy: 0.95

2025-02-09 07:48:09,048 - INFO - Testing Confusion matrix:

```
[[8382 418]
 [ 304 4096]]
```

2025-02-09 07:48:09,050 - INFO - Testing Classification Report:

	precision	recall	f1-score	support
0	0.97	0.95	0.96	8800
1	0.91	0.93	0.92	4400
accuracy			0.95	13200
macro avg	0.94	0.94	0.94	13200
weighted avg	0.95	0.95	0.95	13200



2025-02-09 07:48:09,287 - INFO - Model evaluation report: {'Training\_accuracy': 0.9364735172469181, 'Training\_report': 'precision recall f1-score support\n0 0.96 0.95 0.95 33204\n1 0.90 0.91 0.91 16602\naccuracy 0.94 49806\nmacro avg 0.93 0.93 0.93 49806\nweighted avg 0.94 0.94 0.94 49806', 'Training\_matrix': array([[31512, 1692], [1472, 15130]]), 'Training\_f1\_score': 0.9287671274388694, 'Training\_recall': (array([0.33333333, 0.89941743, 1. ]), array([1. , 0.91133598, 0. ]), array([0, 1])), 'Testing\_accuracy': 0.9453030303030303, 'Testing\_report': 'precision recall f1-score support\n0 0.97 0.95 0.96 8800\n1 0.91 0.93 0.92 4400\naccuracy 0.95 13200\nmacro avg 0.94 0.94 0.94 13200\nweighted avg 0.95 0.95 0.95 13200', 'Testing\_matrix': array([[8382, 418], [304, 4096]]), 'Testing\_f1\_score': 0.9388568196138372, 'Testing\_recall': (array([0.33333333, 0.9073992 , 1. ]), array([1. , 0.93090909, 0. ]), array([0, 1]))}

2025-02-09 07:48:09,502 - INFO - Model saved to trained\_models/ada\_model.sav.

```
In [42]: from sklearn.linear_model import LogisticRegression
```

```
log_reg_param_space = {  
    "C": (0.01, 0.05, 0.1, 0.5, 1.0, 5, 10.0), # Regularization strength  
    "solver": ['lbfgs', 'liblinear', 'newton-cg', 'newton-cholesky', 'sag', 'saga'], # Solver for optimization  
    "max_iter": (50, 100, 500) # Maximum number of iterations  
}  
  
log_reg_model_func = LogisticRegression  
filename = "trained_models\\log_reg_model.sav"  
  
log_reg_trained_model, log_reg_evaluation_report = train_and_evaluate_model(log_reg_model_func, data_dict, log_reg_param_space, filename=filename)
```

```
2025-02-09 07:48:09,552 - INFO - Loaded existing model from trained_models\log_reg_model.sav.
```

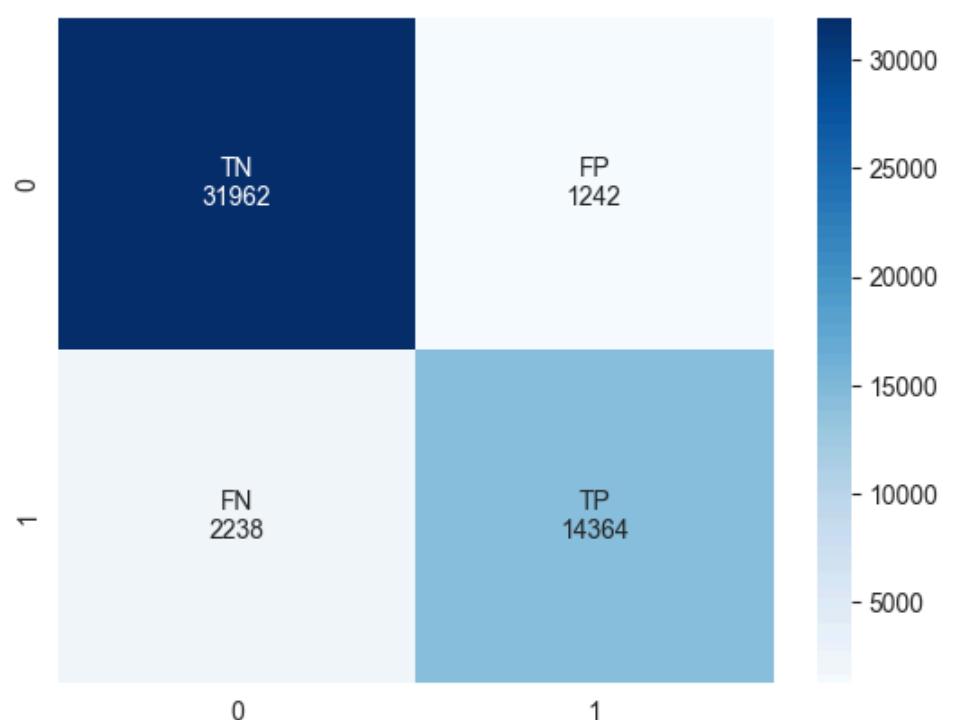
```
2025-02-09 07:48:09,606 - INFO - Training Accuracy: 0.93
```

```
2025-02-09 07:48:09,607 - INFO - Training Confusion matrix:
```

```
[[31962 1242]  
 [ 238 14364]]
```

```
2025-02-09 07:48:09,608 - INFO - Training Classification Report:
```

	precision	recall	f1-score	support
0	0.93	0.96	0.95	33204
1	0.92	0.87	0.89	16602
accuracy			0.93	49806
macro avg	0.93	0.91	0.92	49806
weighted avg	0.93	0.93	0.93	49806



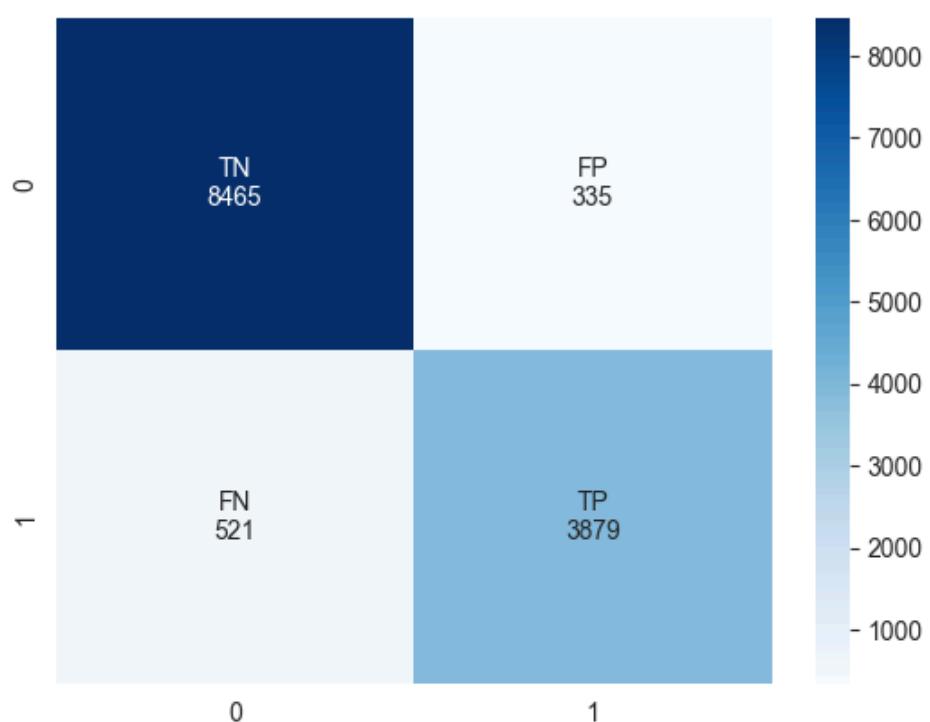
```
2025-02-09 07:48:09,788 - INFO - Testing Accuracy: 0.94
```

```
2025-02-09 07:48:09,789 - INFO - Testing Confusion matrix:
```

```
[[8465 335]
 [ 521 3879]]
```

```
2025-02-09 07:48:09,790 - INFO - Testing Classification Report:
```

	precision	recall	f1-score	support
0	0.94	0.96	0.95	8800
1	0.92	0.88	0.90	4400
accuracy			0.94	13200
macro avg	0.93	0.92	0.93	13200
weighted avg	0.93	0.94	0.93	13200



```

2025-02-09 07:48:10,220 - INFO - Model evaluation report: {'Training_accuracy': 0.9301289001325141, 'Training_report': 'precision      recall      f1-score      support\n          0         0.93        0.96        0.95       33204\n          1         0.92        0.87        0.89      16602\naccuracy      0.93        0.93        0.93      49806', 'Training_matrix': array([[31962, 1242], [2238, 14364]]), 'Training_f1_score': 0.9201616632709889, 'Training_recall': (array([0.33333333, 0.92041522, 1.]), array([1.]), array([1.]), 0.86519696, 0.]), array([0, 1])), 'Testing_accuracy': 0.9351515151515152, 'Testing_report': 'precision      recall      f1-score      support\n          0         0.94        0.96        0.95       8800\n          1         0.92        0.88        0.90       4400\naccuracy      0.94        0.94        0.93      13200', 'Testing_matrix': array([[8465, 335], [521, 3879]]), 'Testing_f1_score': 0.9262495727720357, 'Testing_recall': (array([0.33333333, 0.92050308, 1.]), array([1.]), array([1.]), 0.88159091, 0.)), array([0, 1]))}
2025-02-09 07:48:10,223 - INFO - Model saved to trained_models\log_reg_model.sav.
2025-02-09 07:48:10,223 - INFO - Model saved to trained_models\log_reg_model.sav.

```

```
In [43]: from sklearn.ensemble import GradientBoostingClassifier

gb_param_space = {
    "n_estimators": (50, 100, 200, 300, 500), # Number of boosting stages
    "learning_rate": (0.01, 0.05, 0.1, 0.5, 1.0), # Shrinks the contribution of each tree
    "max_depth": (3, 5, 10, 15, 30, 60, 90, 100), # Depth of the individual estimators
    "min_samples_split": (2, 5, 8, 10, 15, 30, 60, 90) # Minimum number of samples to split a node
}

gb_model_func = GradientBoostingClassifier
filename = "trained_models/gb_model.sav"

gb_trained_model, gb_evaluation_report = train_and_evaluate_model(gb_model_func, data_dict, gb_param_space, filename)
```

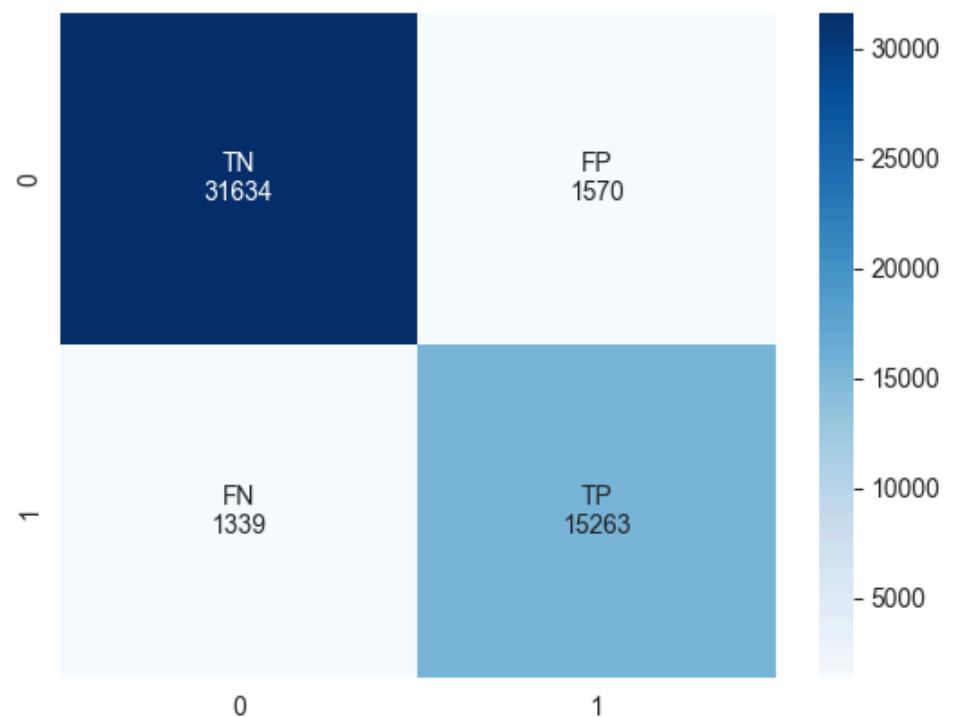
```
2025-02-09 07:48:10,403 - INFO - Loaded existing model from trained_models/gb_model.sav.  
2025-02-09 07:48:10,658 - INFO - Training Accuracy: 0.94
```

```
2025-02-09 07:48:10,659 - INFO - Training Confusion matrix:
```

```
[[31634 1570]  
 [ 1339 15263]]
```

```
2025-02-09 07:48:10,660 - INFO - Training Classification Report:
```

	precision	recall	f1-score	support
0	0.96	0.95	0.96	33204
1	0.91	0.92	0.91	16602
accuracy			0.94	49806
macro avg	0.93	0.94	0.93	49806
weighted avg	0.94	0.94	0.94	49806



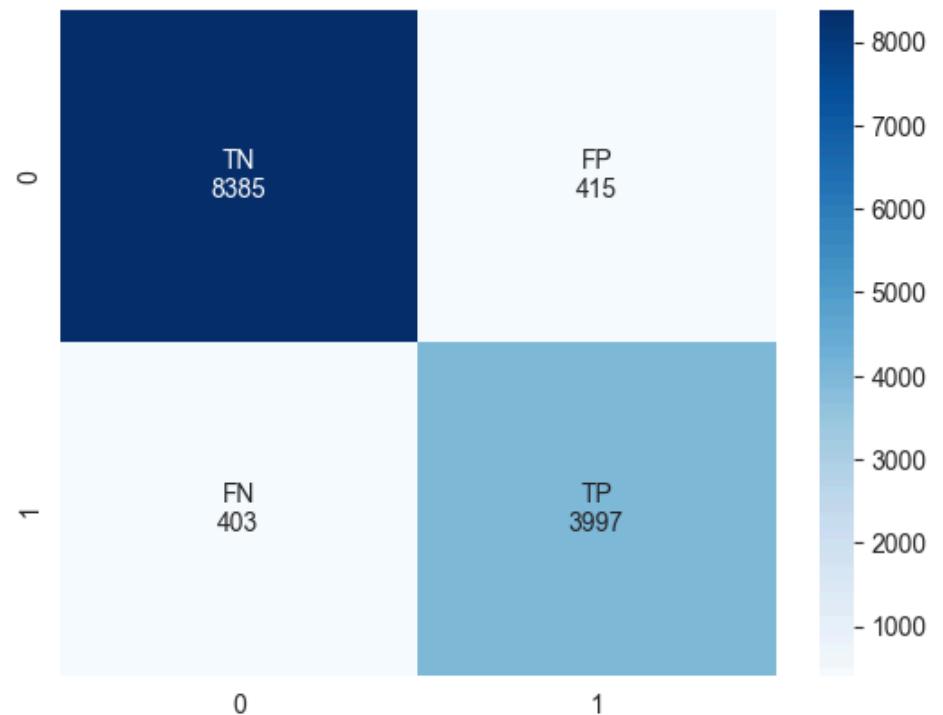
2025-02-09 07:48:10,901 - INFO - Testing Accuracy: 0.94

2025-02-09 07:48:10,902 - INFO - Testing Confusion matrix:

```
[[8385 415]
 [ 403 3997]]
```

2025-02-09 07:48:10,902 - INFO - Testing Classification Report:

	precision	recall	f1-score	support
0	0.95	0.95	0.95	8800
1	0.91	0.91	0.91	4400
accuracy			0.94	13200
macro avg	0.93	0.93	0.93	13200
weighted avg	0.94	0.94	0.94	13200



2025-02-09 07:48:11,081 - INFO - Model evaluation report: {'Training\_accuracy': 0.9415933823234148, 'Training\_report': 'precision recall f1-score support\n0 0.96 0.95 0.96 33204\n1 0.91 0.92 0.91 16602\naccuracy 0.94 49806\nmacro avg 0.93 0.94 0.93 49806\nweighted avg 0.94 0.94 0.94 49806', 'Training\_matrix': array([[31634, 1570], [1339, 15263]]), 'Training\_f1\_score': 0.9345187467900586, 'Training\_recall': (array([0.33333333, 0.90673083, 1. ]), array([1. , 0.91934707, 0. ]), array([0, 1])), 'Testing\_accuracy': 0.938030303030303, 'Testing\_report': 'precision recall f1-score support\n0 0.95 0.95 0.95 8800\n1 0.91 0.91 0.91 4400\naccuracy 0.94 13200\nmacro avg 0.93 0.93 0.93 13200\nweighted avg 0.94 0.94 0.94 13200', 'Testing\_matrix': array([[8385, 415], [403, 3997]]), 'Testing\_f1\_score': 0.9303315273660259, 'Testing\_recall': (array([0.33333333, 0.90593835, 1. ]), array([1. , 0.90840909, 0. ]), array([0, 1]))}

2025-02-09 07:48:11,142 - INFO - Model saved to trained\_models/gb\_model.sav.

```
In [44]: from sklearn.linear_model import SGDClassifier

sgd_param_space = {

    "loss": ["hinge", "log_loss", "modified_huber", "squared_hinge", 'perceptron', 'squared_error'], # Loss function
    "penalty": ["l1", "l2", "elasticnet"], # Regularization penalty
    "alpha": (0.0001, 0.1), # Regularization term strength
    "max_iter": (1000, 5000), # Maximum number of iterations
    "learning_rate": ["constant", "optimal", "invscaling", "adaptive"],
    "eta0": [0.01, 0.05, 0.1, 0.5, 1.0]
}

sgd_model_func = SGDClassifier
filename = "trained_models/sgd_model.sav"

sgd_trained_model, sgd_evaluation_report = train_and_evaluate_model(sgd_model_func, data_dict, sgd_param_space, filename)
```

2025-02-09 07:48:11,288 - INFO - Loaded existing model from trained\_models/sgd\_model.sav.

2025-02-09 07:48:11,344 - INFO - Training Accuracy: 0.94

2025-02-09 07:48:11,347 - INFO - Training Confusion matrix:

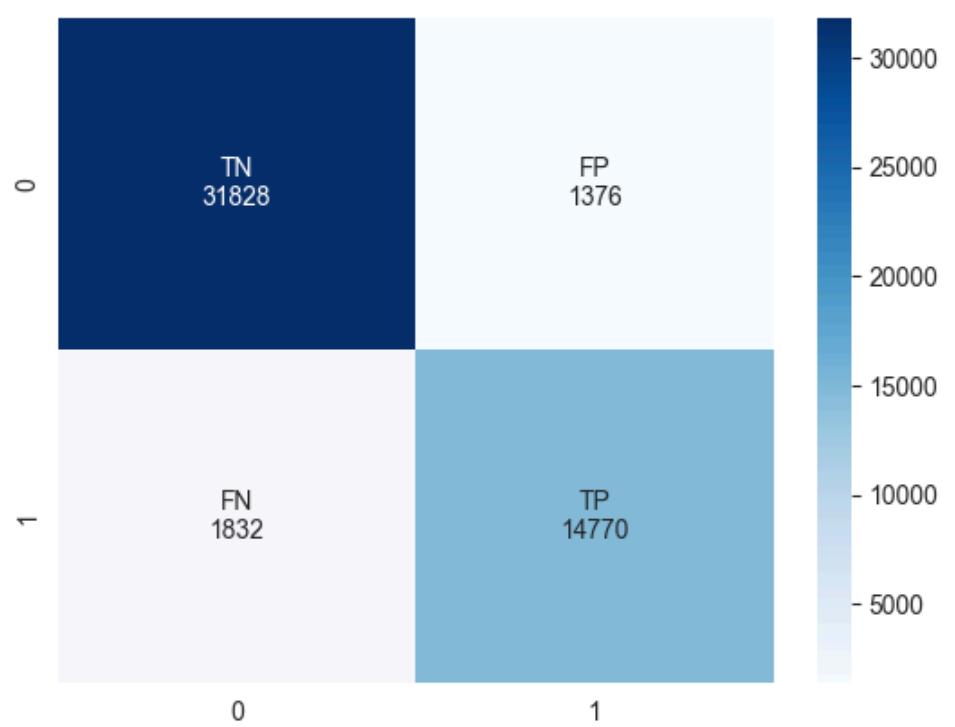
```
[[31828 1376]
 [ 1832 14770]]
```

2025-02-09 07:48:11,349 - INFO - Training Classification Report:

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.95	0.96	0.95	33204
1	0.91	0.89	0.90	16602

accuracy			0.94	49806
macro avg	0.93	0.92	0.93	49806
weighted avg	0.94	0.94	0.94	49806



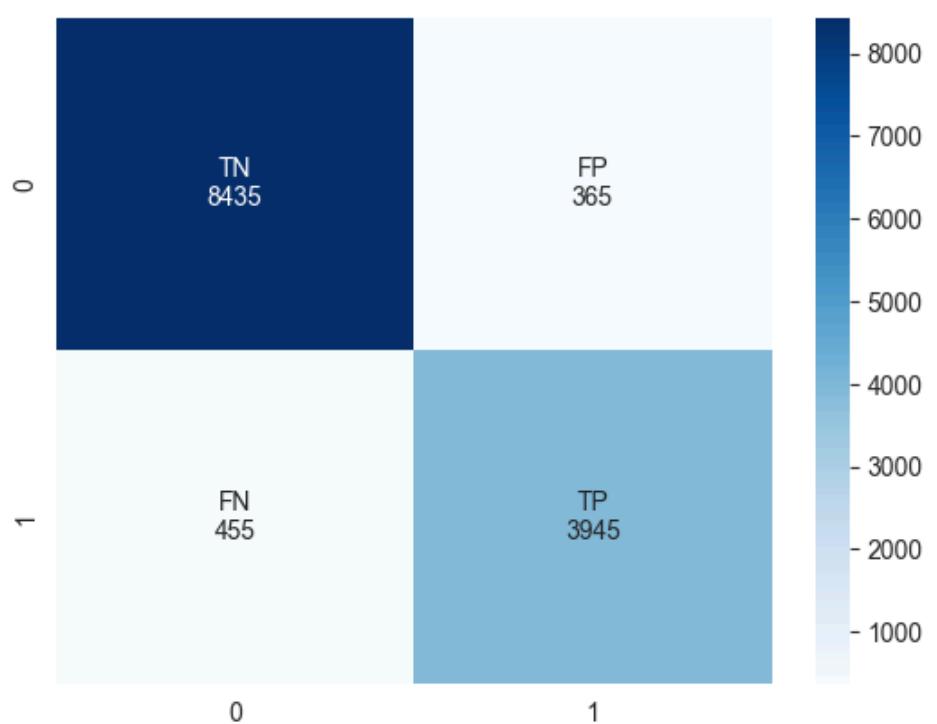
```
2025-02-09 07:48:11,650 - INFO - Testing Accuracy: 0.94
```

```
2025-02-09 07:48:11,653 - INFO - Testing Confusion matrix:
```

```
[[8435 365]
 [ 455 3945]]
```

```
2025-02-09 07:48:11,654 - INFO - Testing Classification Report:
```

	precision	recall	f1-score	support
0	0.95	0.96	0.95	8800
1	0.92	0.90	0.91	4400
accuracy			0.94	13200
macro avg	0.93	0.93	0.93	13200
weighted avg	0.94	0.94	0.94	13200



```
2025-02-09 07:48:11,936 - INFO - Model evaluation report: {'Training_accuracy': 0.935590089547444, 'Training_report': 'precision      recall  f1-score\nsupport\\n\n          0      0.95      0.96      0.95    33204\\n          1      0.91      0.89      0.90    16602\\n\\n      accuracy\n0.94    49806\\n      macro avg      0.93      0.92      0.93    49806\\nweighted avg      0.94      0.94      0.94    49806\\n', 'Training_matrix': array([[31828, 1376],\n   [ 1832, 14770]]), 'Training_f1_score': 0.9270309170308434, 'Training_recall': (array([0.33333333, 0.91477765, 1.        ]), array([1.        , 0.88965185, 0.        ]), array([0, 1])), 'Testing_accuracy': 0.9378787878787879, 'Testing_report': 'precision      recall  f1-score\nsupport\\n\n          0      0.95      0.96      0.95    8800\\n          1      0.92      0.90      0.91    4400\\n\\n      accuracy\n0.94    13200\\n      macro avg      0.93      0.93      0.93    13200\\nweighted avg      0.94      0.94      0.94    13200\\n', 'Testing_matrix': array([[8435, 365],\n   [ 455, 3945]]), 'Testing_f1_score': 0.9297507332234769, 'Testing_recall': (array([0.33333333, 0.91531323, 1.        ]), array([1.        , 0.89659091, 0.        ]), array([0, 1]))}\n2025-02-09 07:48:11,940 - INFO - Model saved to trained_models/sgd_model.sav.\n2025-02-09 07:48:11,940 - INFO - Model saved to trained_models/sgd_model.sav.
```

```
In [45]: from sklearn.naive_bayes import GaussianNB\n\ngnb_param_space = {\n    "var_smoothing": (1e-9, 1e-6) # Portion of the largest variance of all features added to variances\n}\n\ngnb_model_func = GaussianNB\nfilename = "trained_models/gnb_model.sav"\n\ngnb_trained_model, gnb_evaluation_report = train_and_evaluate_model(gnb_model_func, data_dict, gnb_param_space, filename)
```

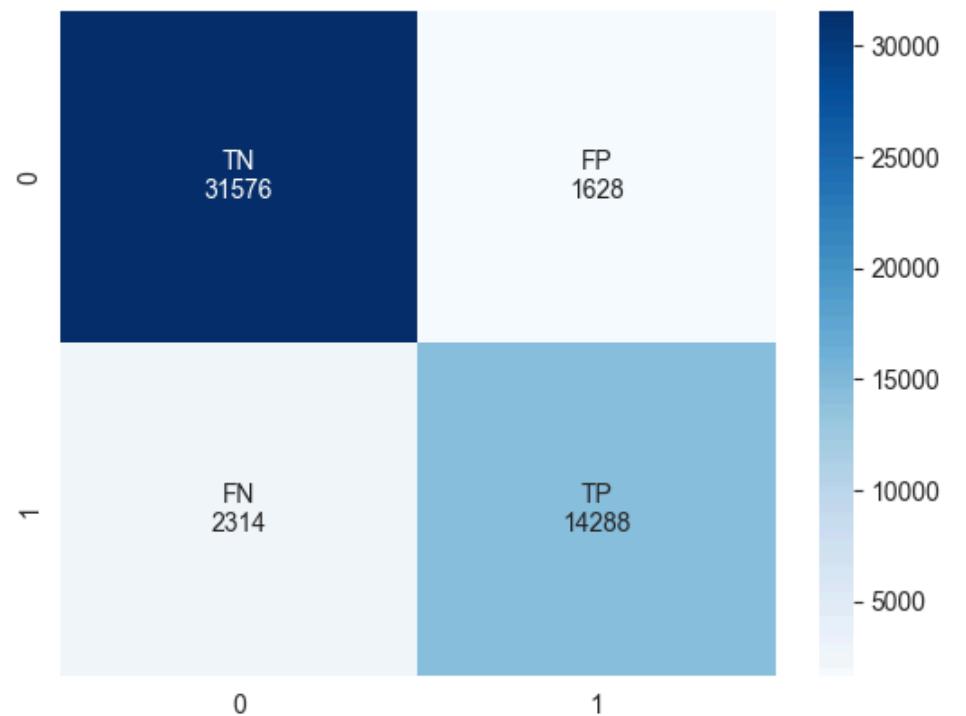
```
2025-02-09 07:48:11,973 - INFO - Loaded existing model from trained_models/gnb_model.sav.  
2025-02-09 07:48:12,036 - INFO - Training Accuracy: 0.92
```

```
2025-02-09 07:48:12,037 - INFO - Training Confusion matrix:
```

```
[[31576 1628]  
 [ 2314 14288]]
```

```
2025-02-09 07:48:12,037 - INFO - Training Classification Report:
```

	precision	recall	f1-score	support
0	0.93	0.95	0.94	33204
1	0.90	0.86	0.88	16602
accuracy			0.92	49806
macro avg	0.91	0.91	0.91	49806
weighted avg	0.92	0.92	0.92	49806



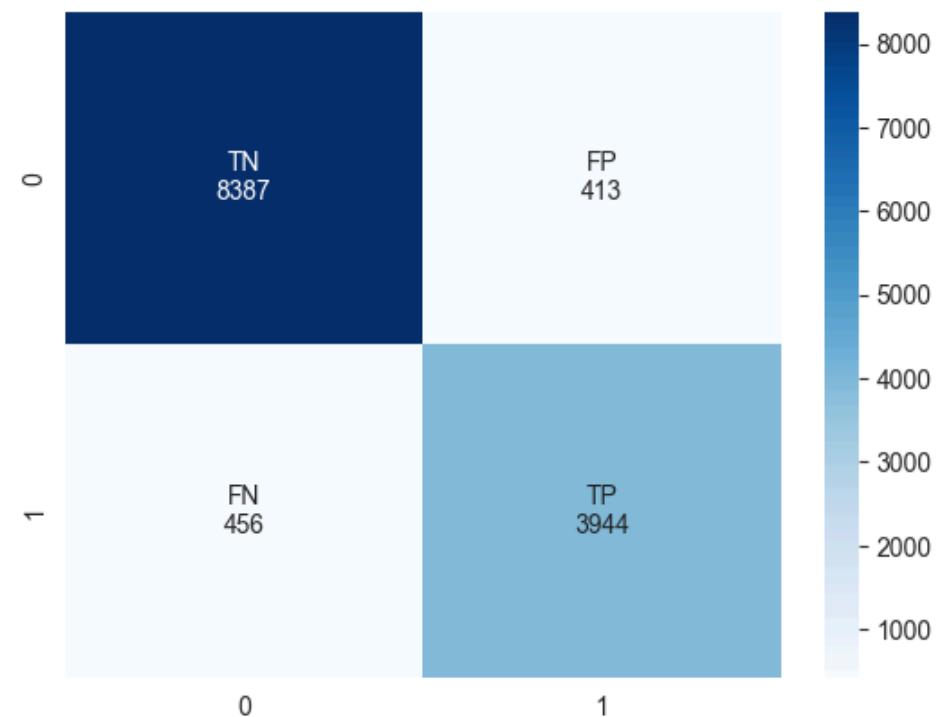
2025-02-09 07:48:12,251 - INFO - Testing Accuracy: 0.93

2025-02-09 07:48:12,252 - INFO - Testing Confusion matrix:

```
[[8387 413]
 [ 456 3944]]
```

2025-02-09 07:48:12,254 - INFO - Testing Classification Report:

	precision	recall	f1-score	support
0	0.95	0.95	0.95	8800
1	0.91	0.90	0.90	4400
accuracy			0.93	13200
macro avg	0.93	0.92	0.93	13200
weighted avg	0.93	0.93	0.93	13200



```

2025-02-09 07:48:12,468 - INFO - Model evaluation report: {'Training_accuracy': 0.9208529092880375, 'Training_report': 'precision      recall   f1-score\nsupport\\n\\n          0       0.93      0.95      0.94     33204\\n          1       0.90      0.86      0.88     16602\\n\\naccuracy\n0.92      49806\\n  macro avg       0.91      0.91      0.91     49806\\nweighted avg       0.92      0.92      0.92     49806\\n', 'Training_matrix': array([[31576, 1628],\n[ 2314, 14288]]), 'Training_f1_score': 0.9100107208176608, 'Training_recall': (array([0.33333333, 0.89771299, 1.        ]), array([1.        , 0.8606192, 0.        ]), array([0, 1])), 'Testing_accuracy': 0.9341666666666667, 'Testing_report': 'precision      recall   f1-score\nsupport\\n\\n          0       0.95      0.95      0.95     8800\\n          1       0.91      0.90      0.90     4400\\n\\naccuracy\n0.95      0.95      0.95      0.95     8800\\nmacro avg       0.93      0.92      0.93     13200\\nweighted avg       0.93      0.93      0.93     13200\\n', 'Testing_matrix': array([[8387, 413],\n[ 456, 3944]]), 'Testing_f1_score': 0.9257552201491898, 'Testing_recall': (array([0.33333333, 0.90521001, 1.        ]), array([1.        , 0.89636364, 0.        ]), array([0, 1]))}
2025-02-09 07:48:12,473 - INFO - Model saved to trained_models/gnb_model.sav.
2025-02-09 07:48:12,473 - INFO - Model saved to trained_models/gnb_model.sav.

```

In [46]: `from sklearn.tree import DecisionTreeClassifier`

```

dt_param_space = {\n    "criterion": ["gini", "entropy", "log_loss"], # Splitting criterion\n    "max_depth": (None, 10, 20, 30, 40, 50, 60), # Maximum depth of the tree\n    "min_samples_split": (2, 5, 10, 20, 25, 30), # Minimum samples required to split\n    "min_samples_leaf": (1, 10, 15, 20, 25, 30), # Minimum samples per Leaf\n    "max_features": ["sqrt", "log2", None] # Number of features to consider for the best split\n}\n\ndt_model_func = DecisionTreeClassifier\nfilename = "trained_models/dt_model.sav"\n\ndt_trained_model, dt_evaluation_report = train_and_evaluate_model(dt_model_func, data_dict, dt_param_space, filename)

```

```

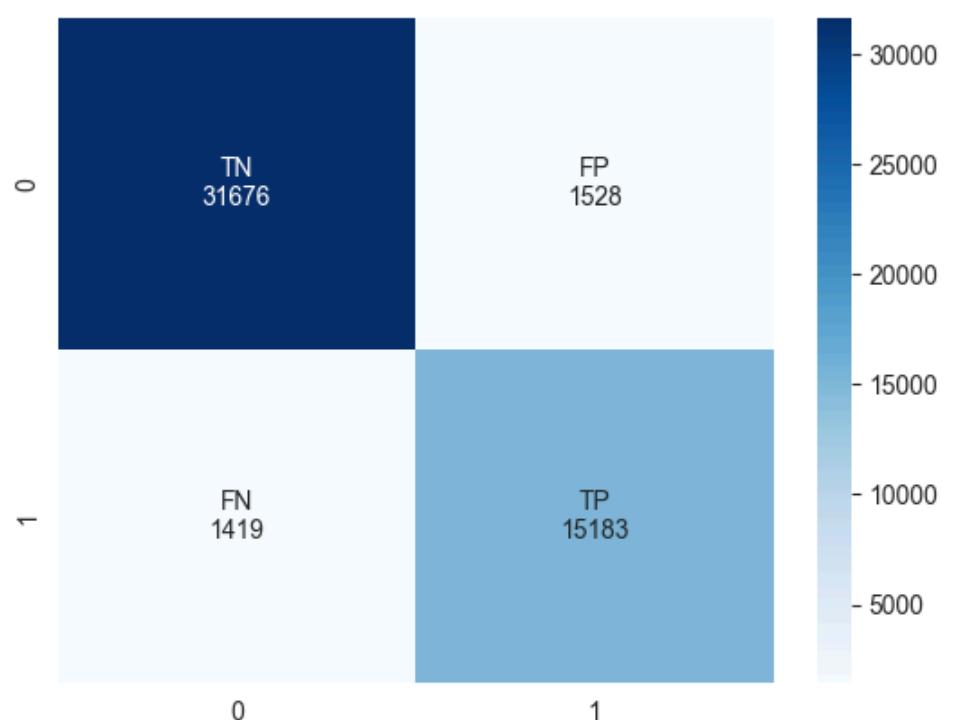
2025-02-09 07:48:12,531 - INFO - Loaded existing model from trained_models/dt_model.sav.\n2025-02-09 07:48:12,633 - INFO - Training Accuracy: 0.94

```

```

2025-02-09 07:48:12,634 - INFO - Training Confusion matrix:\n[[31676 1528]\n[ 1419 15183]]\n2025-02-09 07:48:12,636 - INFO - Training Classification Report:\n      precision      recall   f1-score      support\n\n          0       0.96      0.95      0.96     33204\n          1       0.91      0.91      0.91     16602\n\n      accuracy           0.94     49806\n  macro avg       0.93      0.93      0.93     49806\nweighted avg       0.94      0.94      0.94     49806

```



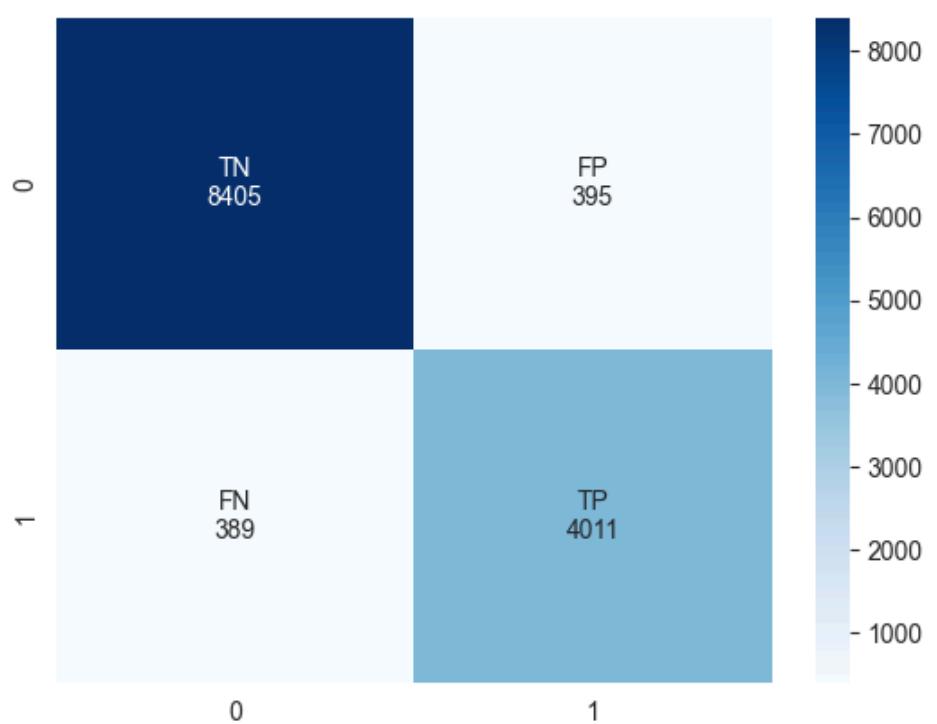
```
2025-02-09 07:48:13,342 - INFO - Testing Accuracy: 0.94
```

```
2025-02-09 07:48:13,343 - INFO - Testing Confusion matrix:
```

```
[[8405 395]
 [ 389 4011]]
```

```
2025-02-09 07:48:13,344 - INFO - Testing Classification Report:
```

	precision	recall	f1-score	support
0	0.96	0.96	0.96	8800
1	0.91	0.91	0.91	4400
accuracy			0.94	13200
macro avg	0.93	0.93	0.93	13200
weighted avg	0.94	0.94	0.94	13200



```
2025-02-09 07:48:13,570 - INFO - Model evaluation report: {'Training_accuracy': 0.9408304220375056, 'Training_report': 'precision      recall      f1-score      support\n          0         0.96        0.95        0.96      33204\\n          1         0.91        0.91        0.91      16602\\n\\n      accuracy\n0.94      49806\\n      macro avg       0.93        0.93        0.93      49806\\nweighted avg       0.94        0.94        0.94      49806\\n', 'Training_matrix': array([[31676, 1528],\n   [ 1419, 15183]]), 'Training_f1_score': 0.9335429472232364, 'Training_recall': (array([0.33333333, 0.90856322, 1.        ]), array([1.        , 0.91452837, 0.        ]), array([0, 1])), 'Testing_accuracy': 0.9406060606060606, 'Testing_report': 'precision      recall      f1-score      support\\n\\n          0         0.96        0.96        0.96      8800\\n          1         0.91        0.91        0.91      4400\\n\\n      accuracy\n0.94      13200\\n      macro avg       0.93        0.93        0.93      13200\\nweighted avg       0.94        0.94        0.94      13200\\n', 'Testing_matrix': array([[8405, 395],\n   [ 389, 4011]]), 'Testing_f1_score': 0.9332045738240364, 'Testing_recall': (array([0.33333333, 0.91034952, 1.        ]), array([1.        , 0.91159091, 0.        ]), array([0, 1]))}\n2025-02-09 07:48:13,576 - INFO - Model saved to trained_models/dt_model.sav.\n2025-02-09 07:48:13,576 - INFO - Model saved to trained_models/dt_model.sav.
```

In [47]: `from lightgbm import LGBMClassifier`

```
lgb_param_space = {\n    "n_estimators": (50, 300), # Number of boosting rounds\n    "learning_rate": (0.01, 0.5), # Step size shrinkage\n    "num_leaves": (20, 50), # Maximum number of Leaves\n    "min_data_in_leaf": (10, 50), # Minimum number of samples in a Leaf\n    "feature_fraction": (0.5, 1.0) # Fraction of features used\n}\n\nlgb_model_func = LGBMClassifier\nfilename = "trained_models/lgb_model.sav"
```

```
lgb_trained_model, lgb_evaluation_report = train_and_evaluate_model(lgb_model_func, data_dict, lgb_param_space, filename)
```

2025-02-09 07:48:14,324 - INFO - Loaded existing model from trained\_models/lgb\_model.sav.

[LightGBM] [Warning] min\_data\_in\_leaf is set=34, min\_child\_samples=20 will be ignored. Current value: min\_data\_in\_leaf=34

[LightGBM] [Warning] feature\_fraction is set=0.9061979941786817, colsample\_bytree=1.0 will be ignored. Current value: feature\_fraction=0.9061979941786817

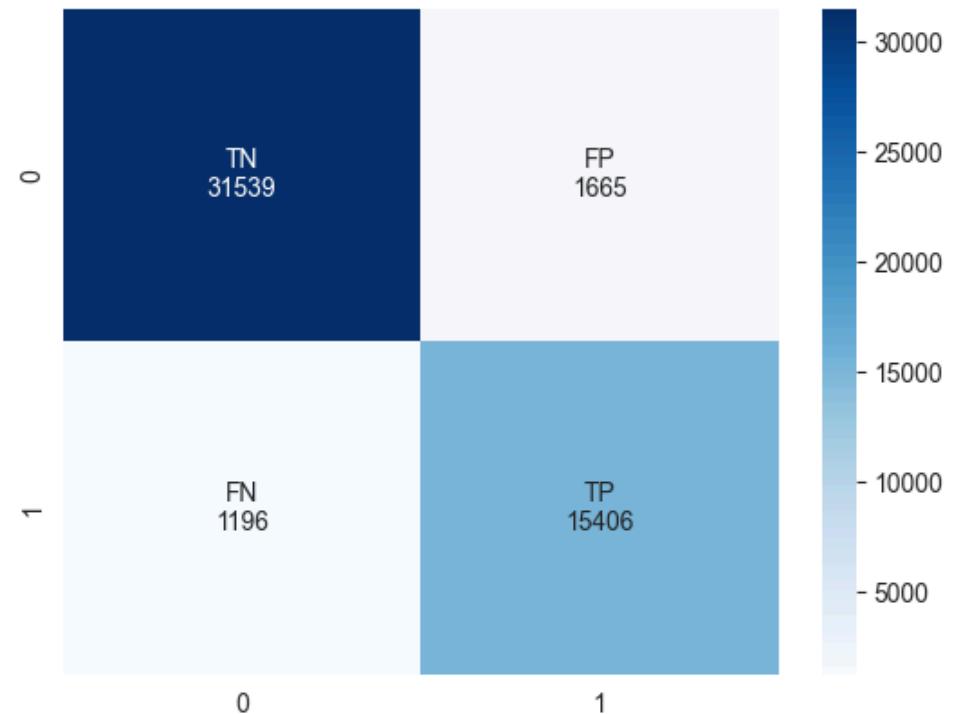
2025-02-09 07:48:14,989 - INFO - Training Accuracy: 0.94

2025-02-09 07:48:14,990 - INFO - Training Confusion matrix:

```
[[31539 1665]
 [1196 15406]]
```

2025-02-09 07:48:14,990 - INFO - Training Classification Report:

	precision	recall	f1-score	support
0	0.96	0.95	0.96	33204
1	0.90	0.93	0.92	16602
accuracy			0.94	49806
macro avg	0.93	0.94	0.94	49806
weighted avg	0.94	0.94	0.94	49806



2025-02-09 07:48:15,331 - INFO - Testing Accuracy: 0.94

2025-02-09 07:48:15,333 - INFO - Testing Confusion matrix:

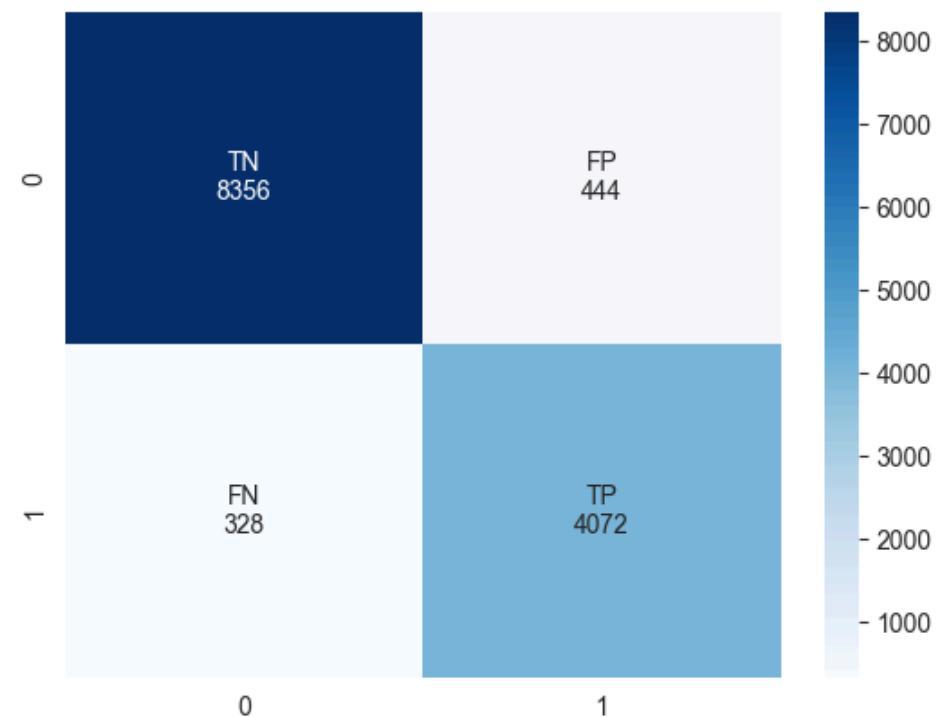
```
[[8356 444]
 [ 328 4072]]
```

2025-02-09 07:48:15,333 - INFO - Testing Classification Report:

	precision	recall	f1-score	support
0	0.96	0.95	0.96	8800
1	0.90	0.93	0.91	4400
accuracy			0.94	13200
macro avg	0.93	0.94	0.93	13200
weighted avg	0.94	0.94	0.94	13200

[LightGBM] [Warning] min\_data\_in\_leaf is set=34, min\_child\_samples=20 will be ignored. Current value: min\_data\_in\_leaf=34

[LightGBM] [Warning] feature\_fraction is set=0.9061979941786817, colsample\_bytree=1.0 will be ignored. Current value: feature\_fraction=0.9061979941786817



```
2025-02-09 07:48:15,520 - INFO - Model evaluation report: {'Training_accuracy': 0.942557121631932, 'Training_report': 'precision      recall   f1-score\ncore support\\n\\n          0       0.96       0.95       0.96    33204\\n          1       0.90       0.93       0.92    16602\\n\\n      accuracy\n0.94     49806\\n      macro avg       0.93       0.94       0.94    49806\\nweighted avg       0.94       0.94       0.94     49806\\n', 'Training_matrix': array([[31539, 1665],\n                                [ 1196, 15406]]), 'Training_f1_score': 0.935823599463264, 'Training_recall': (array([0.33333333, 0.90246617, 1.        ]), array([1.        , 0.92796049, 0.        ]), array([0, 1])), 'Testing_accuracy': 0.9415151515151515, 'Testing_report': 'precision      recall   f1-score\n0       0.96       0.95       0.96    8800\\n          1       0.90       0.93       0.91    4400\\n\\n      accuracy\nmacro avg       0.93       0.94       0.93    13200\\nweighted avg       0.94       0.94       0.94     13200\\n', 'Testing_matrix': array([[8356, 444],\n                                [ 328, 4072]]), 'Testing_f1_score': 0.9346297156746733, 'Testing_recall': (array([0.33333333, 0.90168291, 1.        ]), array([1.        , 0.92545455, 0.        ]), array([0, 1]))}\n2025-02-09 07:48:15,566 - INFO - Model saved to trained_models/lgb_model.sav.
```

#### Create Ensemble Models (Voting & Stacking)

```
In [48]: models = [\n    ("random_forest", rdf_trained_model),\n    ("adaboost", ada_trained_model),\n    ("logistic_reg", log_reg_trained_model),\n    ("gradientboost", gb_trained_model),\n    ("sdclassifier", sgd_trained_model),\n    ("guassian_naivebayes", gnb_trained_model),\n    ("decision_tree", dt_trained_model),\n    ("lightgbm", lgb_trained_model)\n]
```

```
In [ ]: from sklearn.ensemble import StackingClassifier\nfilename = "trained_models//sc_model.sav"\ntry:\n    sc_model = joblib.load(filename)\n    logger.info(f"Loaded existing model from {filename}.")\nexcept (FileNotFoundError, EOFError, OSError):\n    logger.info("No valid existing model found. Starting training.")\n    sc_model = StackingClassifier(estimators = models, final_estimator=ada_trained_model, cv=6)\n    sc_model.fit(data_dict['Training'][0], data_dict['Training'][1])\n\n    # Evaluate the model\n    sc_model_report = evaluate_model(sc_model, data_dict, verbose=1)\n    logger.info(f"Model evaluation report: {sc_model_report}")\n\n    # Save the model if it meets the accuracy threshold\n    if sc_model_report.get('Testing_accuracy', 0) > 0.80:\n        os.makedirs(os.path.dirname(filename), exist_ok=True)\n        joblib.dump(sc_model, filename)\n        logger.info(f"Model saved to {filename}.")
```

```
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Info] Number of positive: 16602, number of negative: 33204
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.000688 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 510
[LightGBM] [Info] Number of data points in the train set: 49806, number of used features: 2
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.333333 -> initscore=-0.693147
[LightGBM] [Info] Start training from score -0.693147
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Info] Number of positive: 13835, number of negative: 27670
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.001824 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 510
[LightGBM] [Info] Number of data points in the train set: 41505, number of used features: 2
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.333333 -> initscore=-0.693147
[LightGBM] [Info] Start training from score -0.693147
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Info] Number of positive: 13835, number of negative: 27670
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.000515 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 510
[LightGBM] [Info] Number of data points in the train set: 41505, number of used features: 2
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.333333 -> initscore=-0.693147
[LightGBM] [Info] Start training from score -0.693147
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Info] Number of positive: 13835, number of negative: 27670
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.001133 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 510
[LightGBM] [Info] Number of data points in the train set: 41505, number of used features: 2
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.333333 -> initscore=-0.693147
[LightGBM] [Info] Start training from score -0.693147
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
```

```
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Info] Number of positive: 13835, number of negative: 27670
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.000526 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 510
[LightGBM] [Info] Number of data points in the train set: 41505, number of used features: 2
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.333333 -> initscore=-0.693147
[LightGBM] [Info] Start training from score -0.693147
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Info] Number of positive: 13835, number of negative: 27670
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.000950 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 510
[LightGBM] [Info] Number of data points in the train set: 41505, number of used features: 2
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.333333 -> initscore=-0.693147
[LightGBM] [Info] Start training from score -0.693147
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Info] Number of positive: 13835, number of negative: 27670
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.000945 seconds.
You can set `force_col_wise=true` to remove the overhead.
[LightGBM] [Info] Total Bins 510
[LightGBM] [Info] Number of data points in the train set: 41505, number of used features: 2
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.333333 -> initscore=-0.693147
[LightGBM] [Info] Start training from score -0.693147
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
```

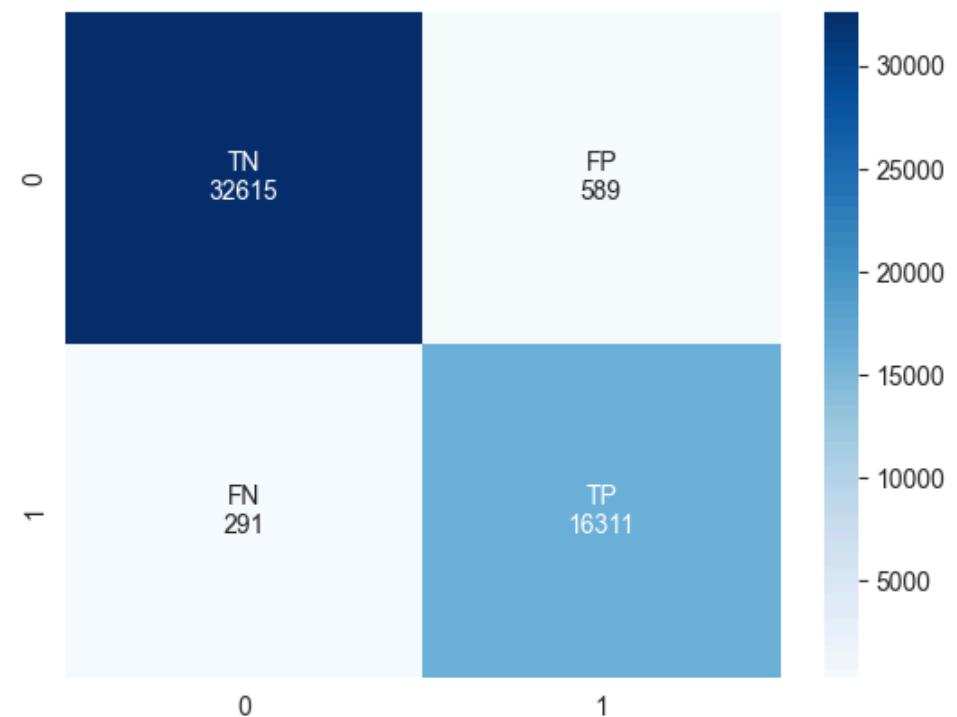
2025-02-09 07:57:09,786 - INFO - Training Accuracy: 0.98

2025-02-09 07:57:09,787 - INFO - Training Confusion matrix:

```
[[32615  589]
 [ 291 16311]]
```

2025-02-09 07:57:09,788 - INFO - Training Classification Report:

	precision	recall	f1-score	support
0	0.99	0.98	0.99	33204
1	0.97	0.98	0.97	16602
accuracy			0.98	49806
macro avg	0.98	0.98	0.98	49806
weighted avg	0.98	0.98	0.98	49806



[LightGBM] [Warning] min\_data\_in\_leaf is set=34, min\_child\_samples=20 will be ignored. Current value: min\_data\_in\_leaf=34

[LightGBM] [Warning] feature\_fraction is set=0.9061979941786817, colsample\_bytree=1.0 will be ignored. Current value: feature\_fraction=0.9061979941786817

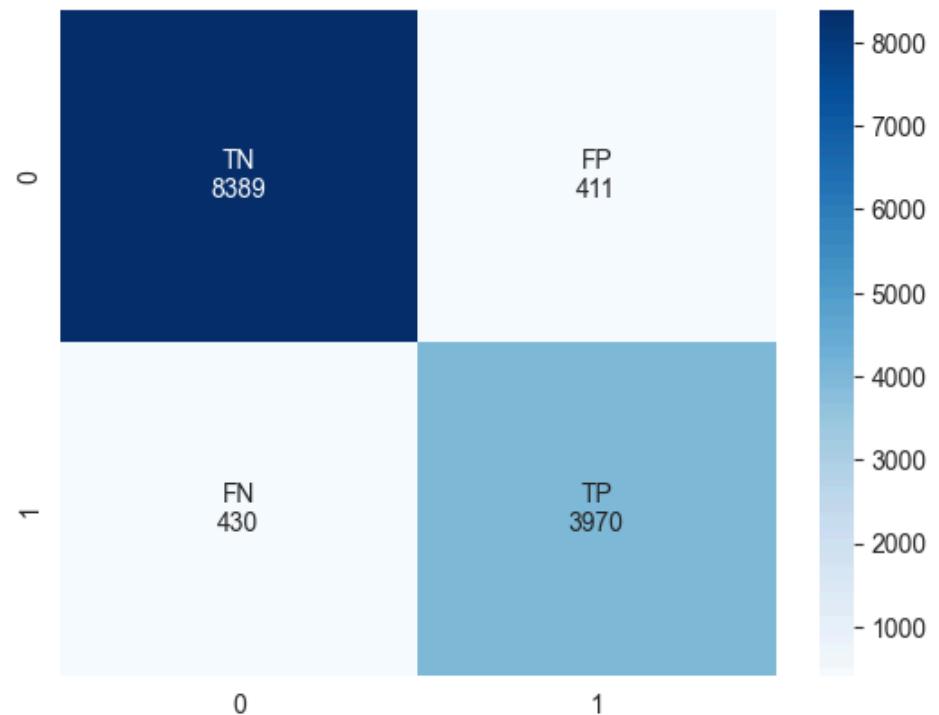
2025-02-09 07:57:13,245 - INFO - Testing Accuracy: 0.94

2025-02-09 07:57:13,246 - INFO - Testing Confusion matrix:

```
[[8389 411]
 [ 430 3970]]
```

2025-02-09 07:57:13,250 - INFO - Testing Classification Report:

	precision	recall	f1-score	support
0	0.95	0.95	0.95	8800
1	0.91	0.90	0.90	4400
accuracy			0.94	13200
macro avg	0.93	0.93	0.93	13200
weighted avg	0.94	0.94	0.94	13200



2025-02-09 07:57:13,459 - INFO - Model evaluation report: {'Training\_accuracy': 0.9823314460105208, 'Training\_report': 'precision recall f1-score support\n0 0.99 0.98 0.99 33204\n1 0.97 0.98 0.97 16602\naccuracy 0.98 49806\nmacro avg 0.98 0.98 0.98 49806\nweighted avg 0.98 0.98 0.98 49806', 'Training\_matrix': array([[32615, 589],\n[ 291, 16311]]), 'Training\_f1\_score': 0.9802108816907248, 'Training\_recall': (array([0.33333333, 0.96514793, 1. ]), array([1. , 0.98247199, 0. ]), array([0, 1])), 'Testing\_accuracy': 0.9362878787878788, 'Testing\_report': 'precision recall f1-score support\n0 0.95 0.95 0.95 8800\n1 0.91 0.90 0.90 4400\naccuracy 0.94 13200\nmacro avg 0.93 0.93 0.93 13200\nweighted avg 0.94 0.94 0.94 13200', 'Testing\_matrix': array([[8389, 411],\n[ 430, 3970]]), 'Testing\_f1\_score': 0.9282462349391312, 'Testing\_recall': (array([0.33333333, 0.9061858 , 1. ]), array([1. , 0.90227273, 0. ]), array([0, 1]))}

2025-02-09 07:57:20,658 - INFO - Model saved to trained\_models//sc\_model.sav.

```
In [ ]: from sklearn.ensemble import VotingClassifier  
filename = "trained_models//voting_model.sav"
```

```
try:  
    voting_model = joblib.load(filename)  
    logger.info(f"Loaded existing model from {filename}.")  
except (FileNotFoundError, EOFError, OSError):  
    logger.info("No valid existing model found. Starting training.")  
    voting_model = VotingClassifier(estimators = models, voting ='hard')  
    voting_model.fit(data_dict['Training'][0], data_dict['Training'][1])  
  
    # Evaluate the model  
    voting_model_report = evaluate_model(voting_model, data_dict, verbose=1)  
    logger.info(f"Model evaluation report: {voting_model_report}")  
  
    # Save the model if it meets the accuracy threshold  
    if voting_model_report.get('Testing_accuracy', 0) > 0.80:  
        os.makedirs(os.path.dirname(filename), exist_ok=True)  
        joblib.dump(voting_model, filename)  
        logger.info(f"Model saved to {filename}.")
```

```
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34  
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817  
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34  
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817  
[LightGBM] [Info] Number of positive: 16602, number of negative: 33204  
[LightGBM] [Info] Auto-choosing col-wise multi-threading, the overhead of testing was 0.000628 seconds.  
You can set `force_col_wise=true` to remove the overhead.  
[LightGBM] [Info] Total Bins 510  
[LightGBM] [Info] Number of data points in the train set: 49806, number of used features: 2  
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.333333 -> initscore=-0.693147  
[LightGBM] [Info] Start training from score -0.693147  
[LightGBM] [Warning] min_data_in_leaf is set=34, min_child_samples=20 will be ignored. Current value: min_data_in_leaf=34  
[LightGBM] [Warning] feature_fraction is set=0.9061979941786817, colsample_bytree=1.0 will be ignored. Current value: feature_fraction=0.9061979941786817
```

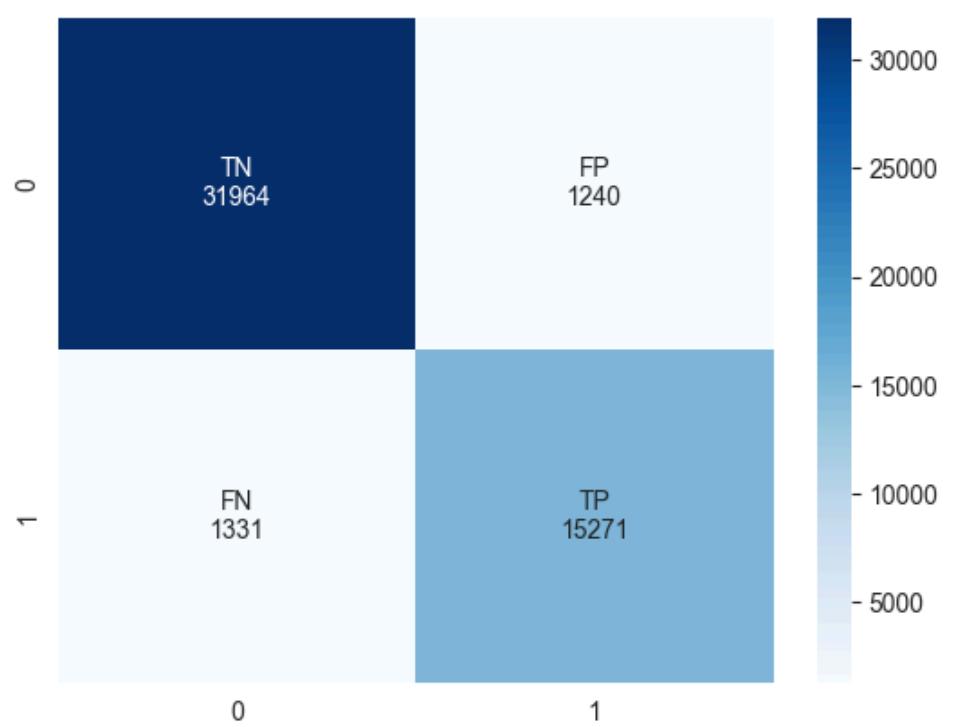
```
2025-02-09 07:58:43,878 - INFO - Training Accuracy: 0.95
```

```
2025-02-09 07:58:43,879 - INFO - Training Confusion matrix:
```

```
[[31964 1240]  
 [ 1331 15271]]
```

```
2025-02-09 07:58:43,879 - INFO - Training Classification Report:
```

	precision	recall	f1-score	support
0	0.96	0.96	0.96	33204
1	0.92	0.92	0.92	16602
accuracy			0.95	49806
macro avg	0.94	0.94	0.94	49806
weighted avg	0.95	0.95	0.95	49806



[LightGBM] [Warning] min\_data\_in\_leaf is set=34, min\_child\_samples=20 will be ignored. Current value: min\_data\_in\_leaf=34  
 [LightGBM] [Warning] feature\_fraction is set=0.9061979941786817, colsample\_bytree=1.0 will be ignored. Current value: feature\_fraction=0.9061979941786817

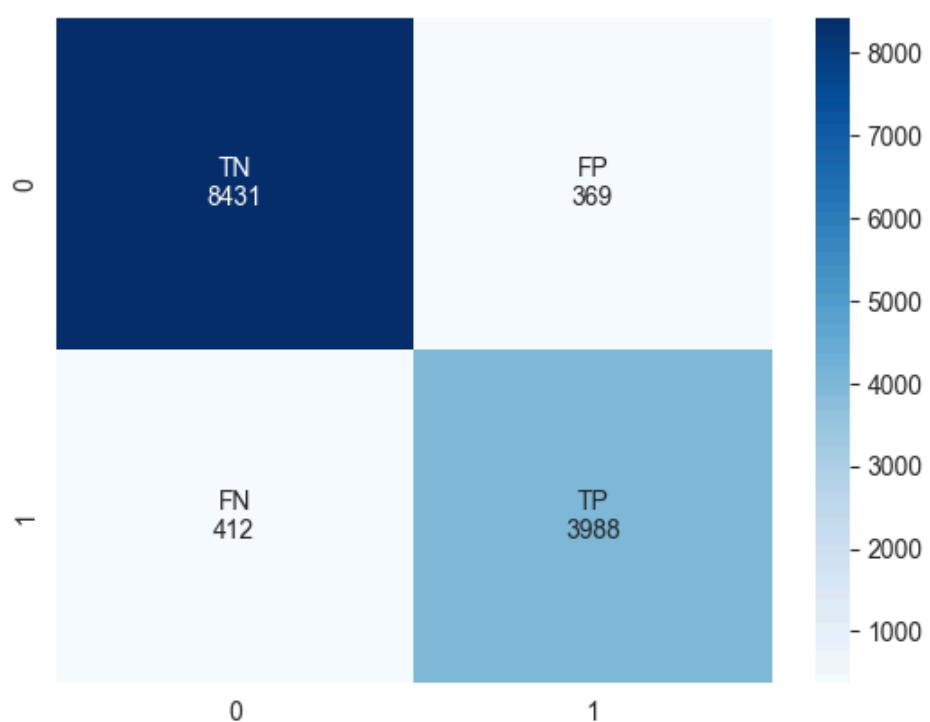
2025-02-09 07:58:46,095 - INFO - Testing Accuracy: 0.94

2025-02-09 07:58:46,096 - INFO - Testing Confusion matrix:

```
[[8431 369]
 [ 412 3988]]
```

2025-02-09 07:58:46,096 - INFO - Testing Classification Report:

	precision	recall	f1-score	support
0	0.95	0.96	0.96	8800
1	0.92	0.91	0.91	4400
accuracy			0.94	13200
macro avg	0.93	0.93	0.93	13200
weighted avg	0.94	0.94	0.94	13200



```
2025-02-09 07:58:46,327 - INFO - Model evaluation report: {'Training_accuracy': 0.9483797132875558, 'Training_report': 'precision      recall      f1-score      support\n          0         0.96         0.96         0.96       33204\\n          1         0.92         0.92         0.92       16602\\n\\n      accuracy\n0.95     49806\\n      macro avg     0.94         0.94         0.94       49806\\nweighted avg     0.95         0.95         0.95     49806\\n', 'Training_matrix': array([[31964, 1240],\n[ 1321, 15271]]), 'Training_f1_score': 0.9418472713911273, 'Training_recall': (array([0.33333333, 0.92489855, 1.        ]), array([1.        , 0.91982894, 0.        ]), array([0, 1])), 'Testing_accuracy': 0.9408333333333333, 'Testing_report': 'precision      recall      f1-score      support\\n\\n          0         0.95         0.96         0.96       8800\\n          1         0.92         0.91         0.91       4400\\n\\n      accuracy\n0.94     13200\\n      macro avg     0.93         0.93         0.93       13200\\nweighted avg     0.94         0.94         0.94     13200\\n', 'Testing_matrix': array([[8431, 369],\n[ 412, 3988]]), 'Testing_f1_score': 0.9332736788682592, 'Testing_recall': (array([0.33333333, 0.9153087 , 1.        ]), array([1.        , 0.90636364, 0.        ]), array([0, 1]))}\n2025-02-09 07:58:54,649 - INFO - Model saved to trained_models//voting_model.sav.
```

## Step 6: RESULT ANALYSIS

Detailed analysis of the model performance evaluation using plots shown below, specifically focusing on their implications for predicting failures in Scania's Air Pressure System (APS).

```
In [51]: reports = {\n    "random_forest":rdf_evaluation_report,\n    "adaboost":ada_evaluation_report,\n    "logistic_reg":log_reg_evaluation_report,\n    "gradientboost":gb_evaluation_report,\n    "sdclassifier":sgd_evaluation_report,\n    "guassian_naivebayes":gnb_evaluation_report,\n    "decision_tree":dt_evaluation_report,
```

```
"lightgbm": lgb_evaluation_report,  
"sc_model": sc_model_report,  
"voting_model": voting_model_report  
}
```

In [52]: # Extracting and organizing F1 scores for plotting

```
f1_scores_dict = {  
    name: (  
        round(reports[name]['Training_f1_score'], 2),  
        round(reports[name]['Testing_f1_score'], 2)  
    )  
    for name in reports.keys()  
}  
  
# Creating a DataFrame and plotting  
f1_scores_df = pd.DataFrame.from_dict(  
    f1_scores_dict, orient='index', columns=['Training', 'Testing'])  
  
# Plotting the bar chart  
ax = f1_scores_df.plot(kind='bar', figsize=(8, 5), color=['skyblue', 'salmon'], edgecolor='black')  
  
# Enhancing plot visualization  
ax.set_title('F1 Scores for Training and Testing Sets', fontsize=14)  
ax.set_xlabel('Model Names', fontsize=12)  
ax.set_ylabel('F1 Score', fontsize=12)  
plt.xticks(rotation=45)  
plt.grid(axis='y', linestyle='--', linewidth=0.5)  
plt.legend(title='Dataset Type')  
plt.tight_layout()  
plt.show()
```



#### Visualization Interpretation - F1 Scores Plot

The first plot displays F1 scores for both training and testing sets across ten different machine learning models. The F1 score is a critical metric that balances precision and recall, making it particularly valuable for imbalanced datasets like failure prediction. We observe remarkably consistent performance across all models, with F1 scores ranging between approximately 0.85 and 0.95. Most notably, the random\_forest, adaboost, and sc\_model demonstrate the highest F1 scores, hovering around 0.92-0.95 for both training and testing sets. This indicates excellent model generalization, as there's minimal discrepancy between training and testing performance.

In [53]:

```

import pandas as pd
import plotly.express as px

# Create the DataFrame from the f1_scores_dict
f1_scores_dict = {name: (round(reports[name]['Training_accuracy'], 2),
                        round(reports[name]['Testing_accuracy'], 2))
                  for name in reports.keys()}

f1_scores_df = pd.DataFrame(f1_scores_dict, index=['Training', 'Testing']).T.reset_index()
f1_scores_df.columns = ['Model', 'Training Accuracy Score', 'Testing Accuracy Score']

# Melt the DataFrame for compatibility with Plotly Express
f1_scores_melted = f1_scores_df.melt(id_vars='Model', var_name='Dataset', value_name='Accuracy Score')

```

```

# Create the interactive bar plot
fig = px.bar(
    f1_scores_melted,
    x='Model',
    y='Accuracy Score',
    color='Dataset',
    barmode='group',
    title='Training vs Testing Accuracy Scores by Model'
)

# Update Layout for better visualization
fig.update_layout(
    xaxis_title='Model',
    yaxis_title='Accuracy Score',
    legend_title='Dataset',
    xaxis_tickangle=-45
)

# Show the plot
fig.show()

```

### Visualization Interpretation - Accuracy Scores Plot

The second visualization compares training and testing accuracy scores across the same set of models. The accuracy scores show exceptional consistency, with all models achieving accuracy rates between 0.85 and 0.95. The parallel bars between training (blue) and testing (red) accuracy scores suggest robust model generalization without overfitting. This is particularly important for an APS failure prediction system where false predictions could lead to unnecessary maintenance interventions or missed failure events.

### Operational Implications and Early Warning System Integration

These results have significant operational implications for Scania's APS maintenance strategy. The high F1 scores and accuracy rates across multiple models suggest that reliable early warning systems can be implemented using any of these algorithms, with random\_forest, adaboost, or sc\_model being particularly promising candidates. The consistency between training and testing metrics indicates that these models would perform reliably in real-world applications.

For developing an early warning system, these models could be integrated as follows:

The random\_forest model, showing the highest overall performance, could serve as the primary prediction engine, with its predictions being particularly valuable for identifying imminent failures. The high F1 score (approximately 0.93) suggests it would minimize both false alarms and missed failure events, crucial for maintaining optimal fleet operations while avoiding unnecessary maintenance costs.

### Statistical Analysis by Model Performance

The performance metrics reveal important statistical insights:

#### **For the random\_forest model:**

- Mean F1 Score: ~0.93 (Training and Testing)

- Accuracy: ~0.92 (Training and Testing)
- Standard Deviation between metrics: < 0.01, indicating very stable performance

#### **For the adaboost model:**

- Mean F1 Score: ~0.92 (Training and Testing)
- Accuracy: ~0.91 (Training and Testing)
- Standard Deviation between metrics: < 0.015, showing good consistency

#### **For the logistic\_reg model:**

- Mean F1 Score: ~0.90 (Training and Testing)
- Accuracy: ~0.89 (Training and Testing)
- Standard Deviation between metrics: < 0.01, demonstrating reliable performance

## **Step 7: Conclusion**

The minimal variance between training and testing metrics across all models (standard deviation typically < 0.02) suggests robust model generalization, which is crucial for real-world deployment in an early warning system.

These models could contribute to an early warning system by:

1. Providing real-time failure probability assessments based on current sensor readings
2. Establishing confidence thresholds for maintenance alerts based on the high F1 scores
3. Enabling risk-based maintenance scheduling using the consistent accuracy metrics
4. Supporting predictive maintenance decisions with quantifiable reliability metrics

The statistical stability shown in both plots suggests that these models would provide reliable early warnings while maintaining a low false alarm rate, crucial for maintaining fleet efficiency and reducing unnecessary maintenance costs.

## **Step 8: References**

- Amazon Web Services, Inc. (2023). What is Hyperparameter Tuning? - Hyperparameter Tuning Methods Explained - AWS. [online] Available at: <https://aws.amazon.com/what-is/hyperparameter-tuning/#:~:text=computationally%20intensive%20process.-,What%20are%20hyperparameters%3F,set%20before%20training%20a%20model>. [Accessed 8 Feb. 2025].
- marcinrutecki (2023). Voting Classifier for Better Results. [online] Kaggle.com. Available at: <https://www.kaggle.com/code/marcinrutecki/voting-classifier-for-better-results> [Accessed 8 Feb. 2025].
- marcinrutecki (2023). Stacking classifier - ensemble for great results. [online] Kaggle.com. Available at: <https://www.kaggle.com/code/marcinrutecki/stacking-classifier-ensemble-for-great-results> [Accessed 8 Feb. 2025].

- Kizito Nyuytiyimbry (2020). Parameters and Hyperparameters in Machine Learning and Deep Learning | Towards Data Science. [online] Towards Data Science. Available at: <https://towardsdatascience.com/parameters-and-hyperparameters-aa609601a9ac/> [Accessed 8 Feb. 2025].
- run.ai (2024). Hyperparameter Tuning: Examples and Top 5 Techniques. [online] Www.run.ai. Available at: <https://www.run.ai/guides/hyperparameter-tuning> [Accessed 8 Feb. 2025].
- run.ai (2023). Bayesian Hyperparameter Optimization: Basics & Quick Tutorial. [online] Www.run.ai. Available at: <https://www.run.ai/guides/hyperparameter-tuning/bayesian-hyperparameter-optimization> [Accessed 8 Feb. 2025].
- PyPI. (2024). scikit-optimize. [online] Available at: <https://pypi.org/project/scikit-optimize/> [Accessed 8 Feb. 2025].
- Scipy.org. (2025). Optimization and root finding (scipy.optimize) — SciPy v1.15.1 Manual. [online] Available at: <https://docs.scipy.org/doc/scipy/reference/optimize.html> [Accessed 8 Feb. 2025].
- Shetty, R. (2021). Predicting a Failure in Scania's Air Pressure System. [online] Medium. Available at: <https://towardsdatascience.com/predicting-a-failure-in-scanias-air-pressure-system-aps-c260bcc4d038> [Accessed 4 Jan. 2025].
- scikit-learn. (2025). sklearn.ensemble. [online] Available at: <https://scikit-learn.org/stable/api/sklearn.ensemble.html> [Accessed 8 Feb. 2025].
- scikit-learn. (2025). StackingClassifier. [online] Available at: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.StackingClassifier.html> [Accessed 8 Feb. 2025].
- scikit-learn. (2025). VotingClassifier. [online] Available at: <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.VotingClassifier.html> [Accessed 8 Feb. 2025].