

**Check the total percentage of missing values of full dataset after dropping columns with more than 70% of missing values**

## Sensor Component Failure Prediction

In [ ]:

In [3]:

```
import pandas as pd
import seaborn as sns
import numpy as np
from statistics import mean
import matplotlib.pyplot as plt
import warnings
from sklearn.preprocessing import PowerTransformer
import numpy as np
from sklearn.preprocessing import LabelEncoder
from sklearn.pipeline import Pipeline
from sklearn.utils import resample

from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier, GradientBoostingClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, classification_report, ConfusionMatrixDisplay, \
    precision_score, recall_score, f1_score, roc_auc_score, roc_curve, confusion_matrix

from sklearn import metrics
from sklearn.model_selection import train_test_split, RepeatedStratifiedKFold, cross_val_score
from sklearn.preprocessing import OneHotEncoder, MinMaxScaler
from sklearn.compose import ColumnTransformer
from sklearn.impute import SimpleImputer, KNNImputer
from xgboost import XGBClassifier
from sklearn.preprocessing import StandardScaler, MinMaxScaler, RobustScaler
from sklearn.compose import ColumnTransformer
from catboost import CatBoostClassifier

warnings.filterwarnings("ignore")
%matplotlib inline
```

In [2]:

```
pip install catboost
```

Collecting catboost

Downloading catboost-1.1.1-cp39-none-win\_amd64.whl (74.0 MB)

Requirement already satisfied: numpy>=1.16.0 in c:\users\tando\anaconda3\lib\site-packages (from catboost) (1.22.4)

Requirement already satisfied: six in c:\users\tando\anaconda3\lib\site-packages (from catboost) (1.16.0)

Requirement already satisfied: scipy in c:\users\tando\anaconda3\lib\site-packages (from catboost) (1.7.3)

Requirement already satisfied: matplotlib in c:\users\tando\anaconda3\lib\site-packages (from catboost) (3.5.1)

Collecting graphviz

Downloading graphviz-0.20.1-py3-none-any.whl (47 kB)

Requirement already satisfied: plotly in c:\users\tando\anaconda3\lib\site-packages (from



```

'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']

```

We have 1 categorical features : ['class']

**As this is a Sensor data. Interpretation of the data is not required**

### Checking missing values

In [9]:

```

# Plotting Missing values count for each column
fig, ax = plt.subplots(figsize=(15,5))

missing = df.isna().sum().div(df.shape[0]).mul(100).to_frame().sort_values(by=0, ascending = False)

ax.bar(missing.index, missing.values.T[0])
plt.xticks([])
plt.ylabel("Percentage missing")
plt.show()

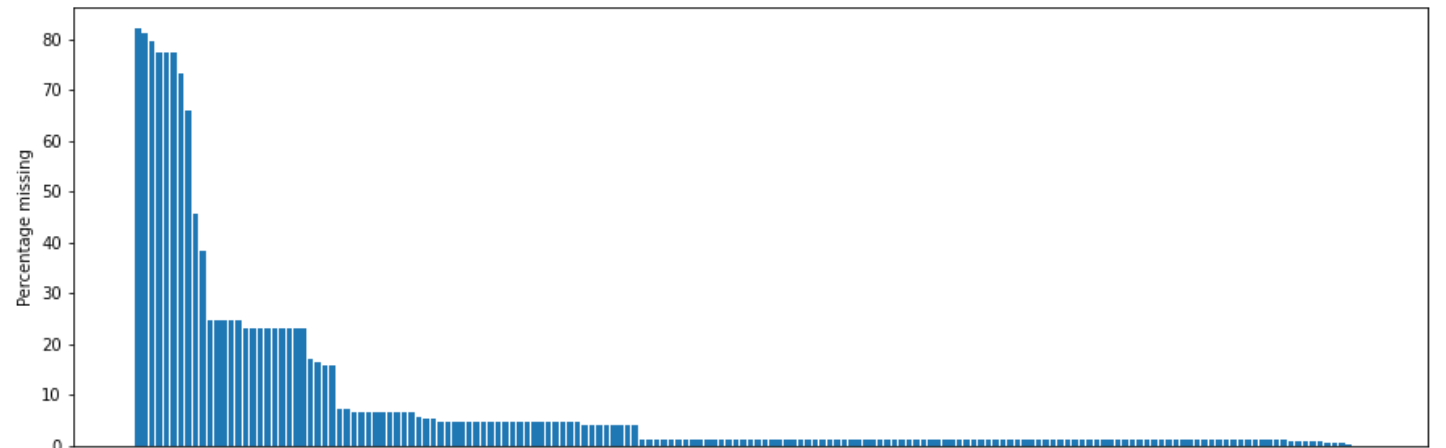
```

C:\Users\tando\AppData\Local\Temp\ipykernel\_23620\2680106201.py:7: MatplotlibDeprecationWarning: Support for passing numbers through unit converters is deprecated since 3.5 and support will be removed two minor releases later; use Axis.convert\_units instead.

```

plt.xticks([])

```



**Dropping Columns which has more than 70% of missing values.**

In [10]:

```

## Dropping columns which has more than 70% of missing values
dropcols = missing[missing[0]>70]
dropcols

```

Out[10]:

	0
br_000	82.106667

<b>bq_000</b>	81.203333
<b>bp_000</b>	79.566667
<b>bo_000</b>	77.221667
<b>ab_000</b>	77.215000
<b>cr_000</b>	77.215000
<b>bn_000</b>	73.348333

In [11]:

```
df.drop(list(dropcols.index), axis=1, inplace=True)
```

In [12]:

```
# Check shape of the dataset after dropping columns
df.shape
```

Out[12]:

```
(60000, 164)
```

In [13]:

```
missing_values_count= df.isnull().sum()
total_cells = np.product(df.shape)
total_missing = missing_values_count.sum()

# percent of data that is missing
print(f"Percentage of total missing cells in the data {(total_missing/total_cells) * 100}%")
```

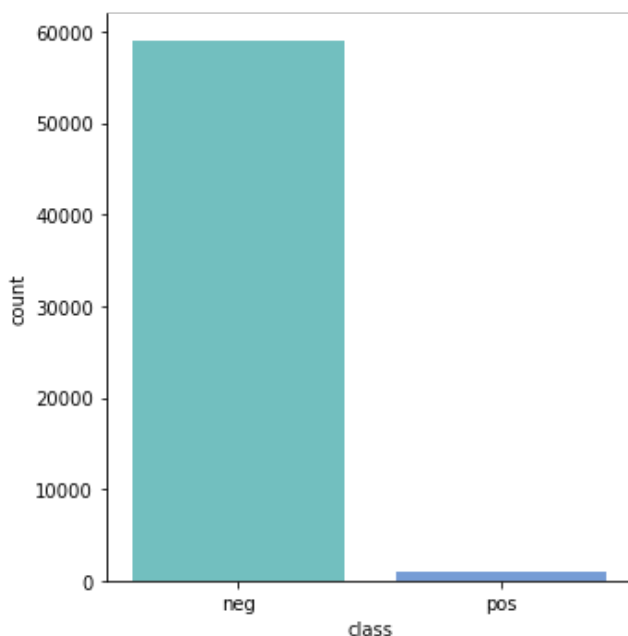
Percentage of total missing cells in the data 5.29765243902439%

## Visualization of unique values in Target variable

In [14]:

```
pos = df[df['class']=='pos'].shape[0]
neg = df[df['class']=='neg'].shape[0]
print("Positive: " + str(pos) + ", Negative: " + str(neg))
sns.catplot(data=df, x="class", kind="count", palette="winter_r", alpha=.6)
plt.show()
```

Positive: 1000, Negative: 59000



- The target classes are highly imbalanced
- Class imbalance is a scenario that arises when we have unequal distribution of class in a dataset i.e. the no. of data points in the negative class (majority class) very large compared to that of the positive class (minority class)
- If the imbalanced data is not treated beforehand, then this will degrade the performance of the classifier model.
- Hence we should handle imbalanced data with certain methods.

## Create Functions for model training and evaluation

In [15]:

```
def evaluate_clf(true, predicted):

    acc = accuracy_score(true, predicted) # Calculate Accuracy
    f1 = f1_score(true, predicted) # Calculate F1-score
    precision = precision_score(true, predicted) # Calculate Precision
    recall = recall_score(true, predicted) # Calculate Recall
    roc_auc = roc_auc_score(true, predicted) #Calculate Roc
    return acc, f1 , precision, recall, roc_auc
```

In [16]:

```
# Create cost of the model as per data description
def total_cost(y_true, y_pred):

    tn, fp, fn, tp = confusion_matrix(y_true, y_pred).ravel()
    cost = 10*fp + 500*fn
    return cost
```

In [17]:

```
# Create a function which can evaluate models and return a report
def evaluate_models(X, y, models):

    # separate dataset into train and test
    X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2,random_state=4
2)

    cost_list=[]
    models_list = []
    accuracy_list = []

    for i in range(len(list(models))):
        model = list(models.values())[i]
        model.fit(X_train, y_train) # Train model

        # Make predictions
        y_train_pred = model.predict(X_train)
        y_test_pred = model.predict(X_test)

        # Training set performance
        model_train_accuracy, model_train_f1,model_train_precision,\
        model_train_recall,model_train_rocauc_score=evaluate_clf(y_train ,y_train_pred)
        train_cost = total_cost(y_train, y_train_pred)

        # Test set performance
        model_test_accuracy,model_test_f1,model_test_precision,\
        model_test_recall,model_test_rocauc_score=evaluate_clf(y_test, y_test_pred)
        test_cost = total_cost(y_test, y_test_pred)

        print(list(models.keys())[i])
        models_list.append(list(models.keys())[i])

        print('Model performance for Training set')
```

```

print("- Accuracy: {:.4f}".format(model_train_accuracy))
print("- F1 score: {:.4f}".format(model_train_f1))
print("- Precision: {:.4f}".format(model_train_precision))
print("- Recall: {:.4f}".format(model_train_recall))
print("- Roc Auc Score: {:.4f}".format(model_train_rocauc_score))
print(f'- COST: {train_cost}.')

print('-----')

print('Model performance for Test set')
print("- Accuracy: {:.4f}".format(model_test_accuracy))
print("- F1 score: {:.4f}".format(model_test_f1))
print("- Precision: {:.4f}".format(model_test_precision))
print("- Recall: {:.4f}".format(model_test_recall))
print("- Roc Auc Score: {:.4f}".format(model_test_rocauc_score))
print(f'- COST: {test_cost}.')
cost_list.append(test_cost)
print('='*35)
print('\n')

report=pd.DataFrame(list(zip(models_list, cost_list)), columns=['Model Name', 'Cost'
]).sort_values(by=["Cost"])

return report

```

## Plot distribution of all Independent Numerical variables

In [18]:

```

numeric_features = [feature for feature in df.columns if df[feature].dtype != 'O']

plt.figure(figsize=(15, 100))
for i, col in enumerate(numeric_features):
    plt.subplot(60, 3, i+1)
    sns.distplot(x=df[col], color='indianred')
    plt.xlabel(col, weight='bold')
    plt.tight_layout()

```

C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)

C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)

C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)

C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)

C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)

C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)

C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)

[illegible]



[illegible]



[illegible]

[illegible]

[illegible]

```
Your code has used distplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)  
C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:  
`distplot` is a deprecated function and will be removed in a future version. Please adapt  
your code to use either `displot` (a figure-level function with similar flexibility) or `  
histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)  
C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:  
`distplot` is a deprecated function and will be removed in a future version. Please adapt  
your code to use either `displot` (a figure-level function with similar flexibility) or `  
histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)  
C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:  
`distplot` is a deprecated function and will be removed in a future version. Please adapt  
your code to use either `displot` (a figure-level function with similar flexibility) or `  
histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)  
C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:  
`distplot` is a deprecated function and will be removed in a future version. Please adapt  
your code to use either `displot` (a figure-level function with similar flexibility) or `  
histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)  
C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:316: UserWarning: Dat  
aset has 0 variance; skipping density estimate. Pass `warn_singular=False` to disable thi  
s warning.  
warnings.warn(msg, UserWarning)  
C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:  
`distplot` is a deprecated function and will be removed in a future version. Please adapt  
your code to use either `displot` (a figure-level function with similar flexibility) or `  
histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)  
C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:  
`distplot` is a deprecated function and will be removed in a future version. Please adapt  
your code to use either `displot` (a figure-level function with similar flexibility) or `  
histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)  
C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:  
`distplot` is a deprecated function and will be removed in a future version. Please adapt  
your code to use either `displot` (a figure-level function with similar flexibility) or `  
histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)  
C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:  
`distplot` is a deprecated function and will be removed in a future version. Please adapt  
your code to use either `displot` (a figure-level function with similar flexibility) or `  
histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)  
C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:  
`distplot` is a deprecated function and will be removed in a future version. Please adapt  
your code to use either `displot` (a figure-level function with similar flexibility) or `  
histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)  
C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:  
`distplot` is a deprecated function and will be removed in a future version. Please adapt  
your code to use either `displot` (a figure-level function with similar flexibility) or `  
histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)  
C:\Users\tando\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning:
```



[illegible]

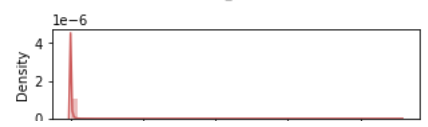
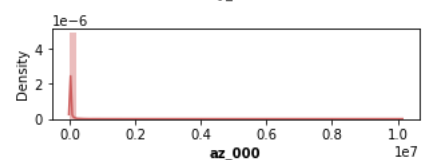
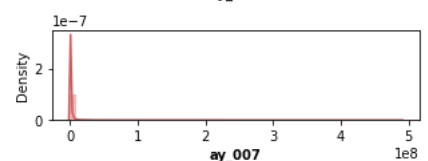
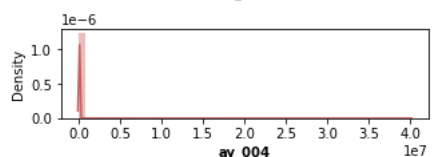
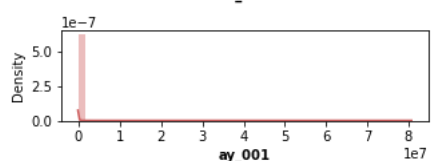
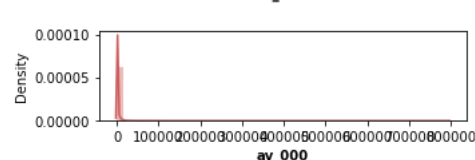
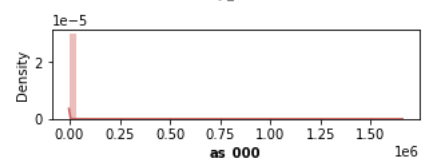
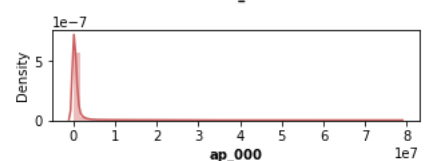
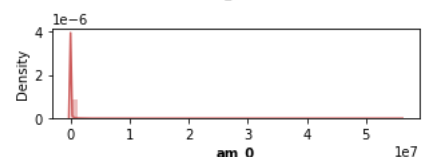
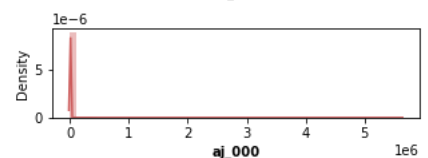
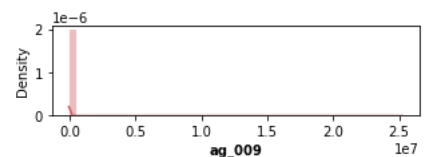
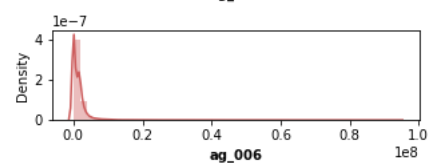
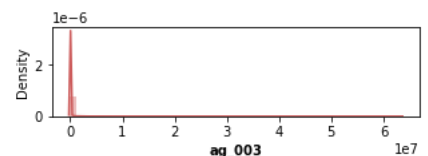
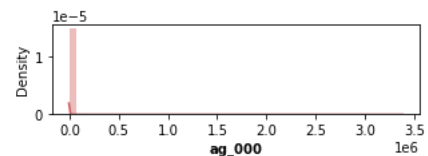
[illegible]

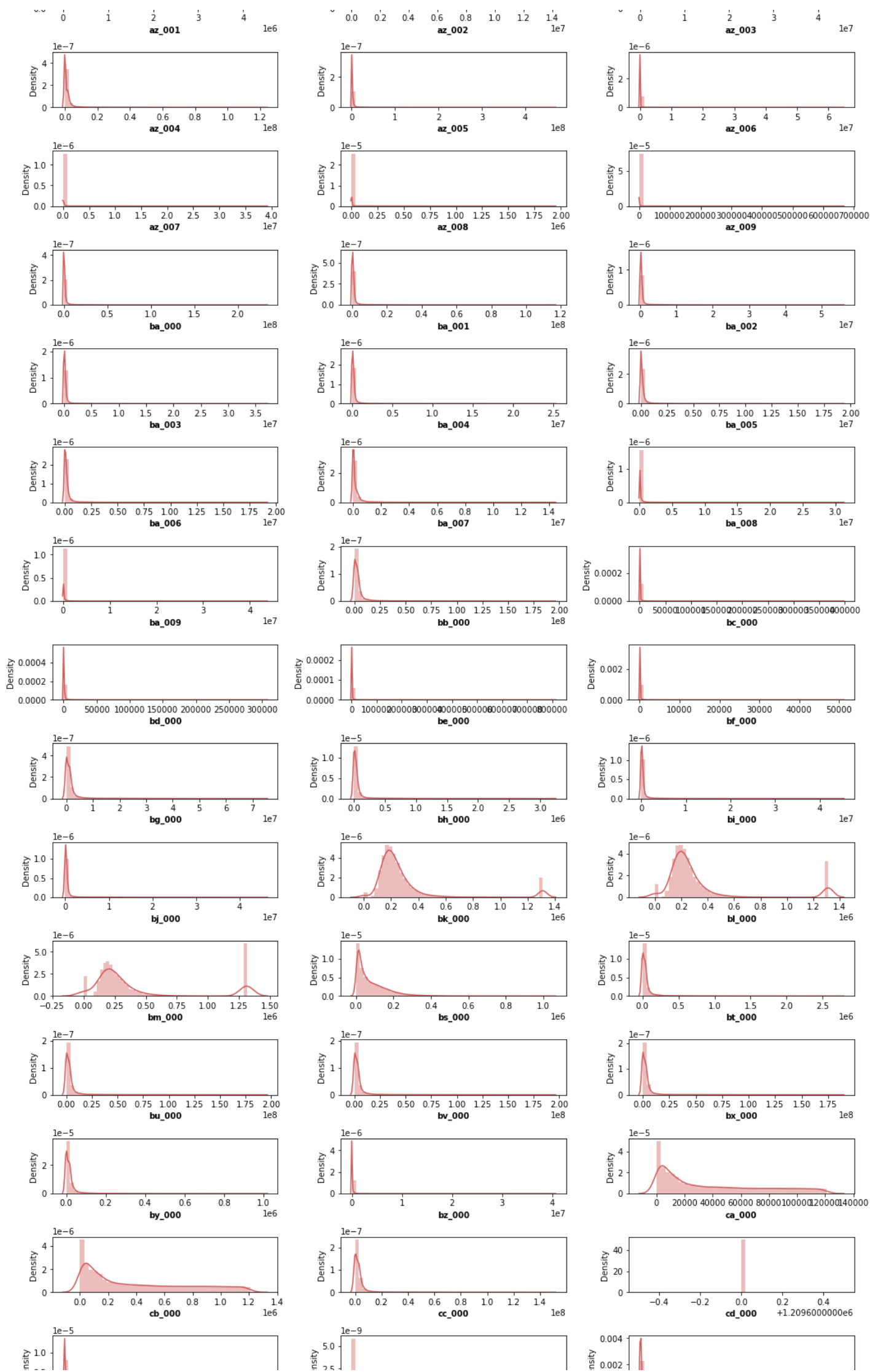
[illegible]

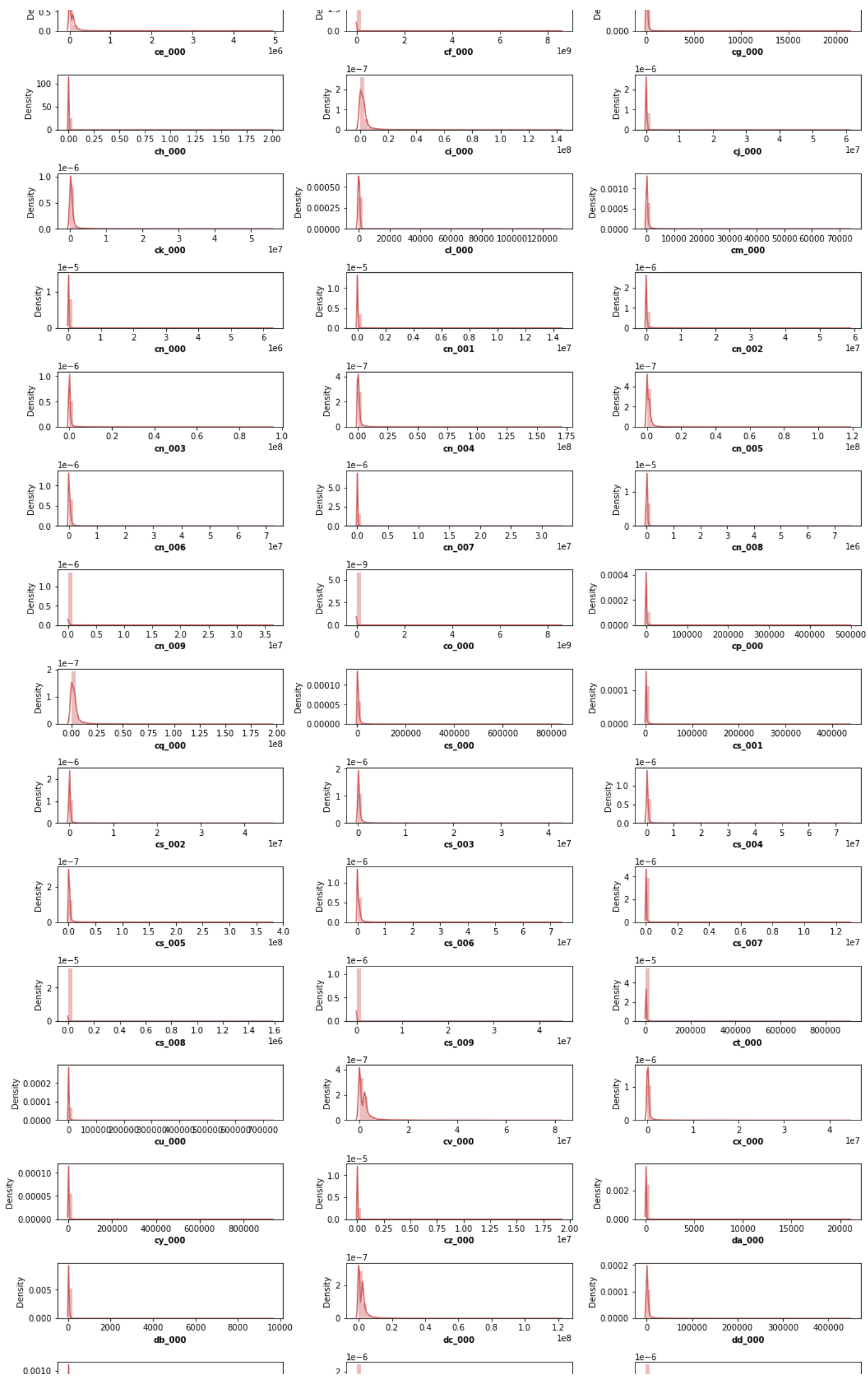


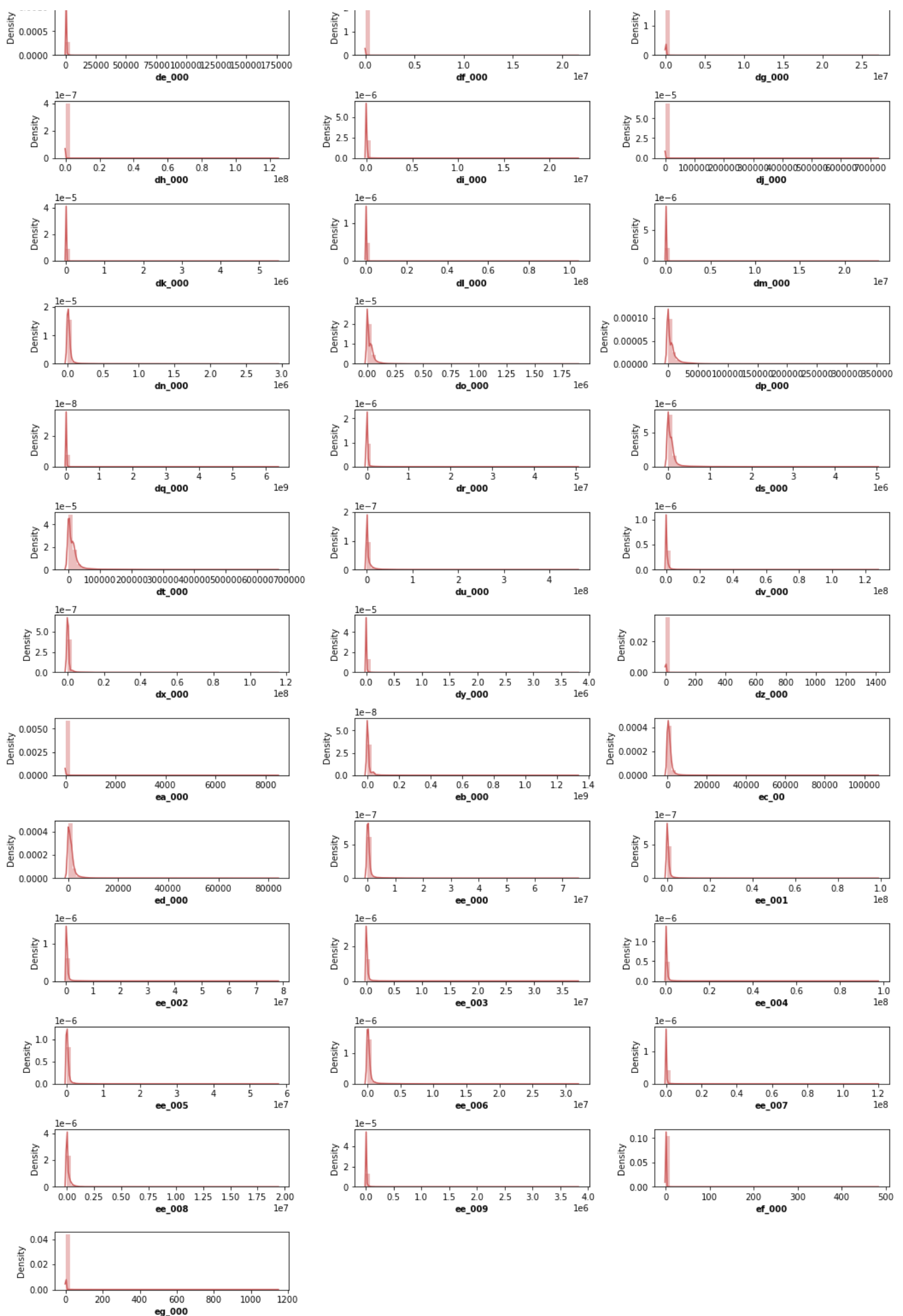
[illegible]











## Report

- As per the above plot most of the features are not normally distributed.
- Transformation of data is not of prime importance since it is a classification problem

- Transformation of data is not of prime importance since it is a classification problem.
- Interpreting each and every column is not necessary as this is sensor data.

## Evaluate Model on Different experiments

In [19]:

```
# Splitting X and y for all Experiments
X= df.drop('class', axis=1)
y = df['class']
```

- **Manually Encoding Target Variable**

In [20]:

```
y= y.replace({'pos': 1, 'neg': 0})
```

In [21]:

```
# Fit with robust scaler for KNN best K-selection experimet
robustscaler = RobustScaler()
X1 = robustscaler.fit_transform(X)
```

### Finding the optimal n\_neighbour value for KNN imputer

In [22]:

```
results=[]
# define imputer
imputer = KNNImputer(n_neighbors=5, weights='uniform', metric='nan_euclidean')
strategies = [str(i) for i in [1,3,5,7,9]]
for s in strategies:
    pipeline = Pipeline(steps=[('i', KNNImputer(n_neighbors=int(s))), ('m', LogisticRegression())])
    scores = cross_val_score(pipeline, X1, y, scoring='accuracy', cv=2, n_jobs=-1)
    results.append(scores)
    print('n_neighbors= %s || accuracy (%.4f)' % (s , mean(scores)))

n_neighbors= 1 || accuracy (0.7702)
n_neighbors= 3 || accuracy (0.7653)
n_neighbors= 5 || accuracy (0.7515)
n_neighbors= 7 || accuracy (0.7727)
n_neighbors= 9 || accuracy (0.7457)
```

**We can observe n\_neighbors=3 able to produce highest accuracy**

### Pipeline for KNN imputer

In [23]:

```
num_features = X.select_dtypes(exclude="object").columns

# Fit the KNN imputer with selected K-value
knn_pipeline = Pipeline(steps=[
    ('imputer', KNNImputer(n_neighbors=3)),
    ('RobustScaler', RobustScaler())
])
```

In [24]:

```
X_knn =knn_pipeline.fit_transform(X)
```

-----  
KeyboardInterrupt

Traceback (most recent call last)

Input In [24], in <cell line: 1>()



```
----> 1 x_knn = knn_pipeline.fit_transform(X)
```

```
File ~\anaconda3\lib\site-packages\sklearn\pipeline.py:414, in Pipeline.fit_transform(self, X, y, **fit_params)
```

```
    387 """Fit the model and transform with the final estimator.
    388
    389 Fits all the transformers one after the other and transform the
    (...)
    411     Transformed samples.
    412 """
    413 fit_params_steps = self._check_fit_params(**fit_params)
--> 414 Xt = self._fit(X, y, **fit_params_steps)
    416 last_step = self._final_estimator
    417 with _print_elapsed_time("Pipeline", self._log_message(len(self.steps) - 1)):
```

```
File ~\anaconda3\lib\site-packages\sklearn\pipeline.py:336, in Pipeline._fit(self, X, y, **fit_params_steps)
```

```
    334 cloned_transformer = clone(transformer)
    335 # Fit or load from cache the current transformer
--> 336 X, fitted_transformer = fit_transform_one_cached(
    337     cloned_transformer,
    338     X,
    339     y,
    340     None,
    341     message_clsname="Pipeline",
    342     message=self._log_message(step_idx),
    343     **fit_params_steps[name],
    344 )
    345 # Replace the transformer of the step with the fitted
    346 # transformer. This is necessary when loading the transformer
    347 # from the cache.
    348 self.steps[step_idx] = (name, fitted_transformer)
```

```
File ~\anaconda3\lib\site-packages\joblib\memory.py:349, in NotMemorizedFunc.__call__(self, *args, **kwargs)
```

```
    348 def __call__(self, *args, **kwargs):
--> 349     return self.func(*args, **kwargs)
```

```
File ~\anaconda3\lib\site-packages\sklearn\pipeline.py:870, in _fit_transform_one(transformer, X, y, weight, message_clsname, message, **fit_params)
```

```
    868 with _print_elapsed_time(message_clsname, message):
    869     if hasattr(transformer, "fit_transform"):
--> 870         res = transformer.fit_transform(X, y, **fit_params)
    871     else:
    872         res = transformer.fit(X, y, **fit_params).transform(X)
```

```
File ~\anaconda3\lib\site-packages\sklearn\base.py:867, in TransformerMixin.fit_transform(self, X, y, **fit_params)
```

```
    863 # non-optimized default implementation; override when a better
    864 # method is possible for a given clustering algorithm
    865 if y is None:
    866     # fit method of arity 1 (unsupervised transformation)
--> 867     return self.fit(X, **fit_params).transform(X)
    868 else:
    869     # fit method of arity 2 (supervised transformation)
    870     return self.fit(X, y, **fit_params).transform(X)
```

```
File ~\anaconda3\lib\site-packages\sklearn\impute\_knn.py:338, in KNNImputer.transform(self, X)
```

```
    329 # process in fixed-memory chunks
    330 gen = pairwise_distances_chunked(
    331     X[row_missing_idx, :],
    332     self._fit_X,
    (...)
    336     reduce_func=process_chunk,
    337 )
--> 338 for chunk in gen:
    339     # process_chunk modifies X in place. No return value.
    340     pass
    342 return super()._concatenate_indicator(X[:, valid_mask], X_indicator)
```

```
File ~\anaconda3\lib\site-packages\sklearn\metrics\pairwise.py:1859, in pairwise_distance
```

```

s_chunked(x, y, reduce_func, metric, n_jobs, working_memory, **kwds)
1857 if reduce_func is not None:
1858     chunk_size = D_chunk.shape[0]
-> 1859     D_chunk = reduce_func(D_chunk, sl.start)
1860     _check_chunk_size(D_chunk, chunk_size)
1861 yield D_chunk

```

File ~\anaconda3\lib\site-packages\sklearn\impute\\_knn.py:321, in KNNImputer.transform.<locals>.process\_chunk(dist\_chunk, start)

```

316     dist_subset = dist_chunk[dist_idx_map[receivers_idx] - start][
317         :, potential_donors_idx
318     ]
320 n_neighbors = min(self.n_neighbors, len(potential_donors_idx))
-> 321 value = self._calc_impute(
322     dist_subset,
323     n_neighbors,
324     self._fit_X[potential_donors_idx, col],
325     mask_fit_X[potential_donors_idx, col],
326 )
327 X[receivers_idx, col] = value

```

File ~\anaconda3\lib\site-packages\sklearn\impute\\_knn.py:159, in KNNImputer.\_calc\_impute(self, dist\_pot\_donors, n\_neighbors, fit\_X\_col, mask\_fit\_X\_col)

```

135 """Helper function to impute a single column.
136
137 Parameters
138 (...)
139     Imputed values for receiver.
140 """
141 # Get donors
-> 159 donors_idx = np.argpartition(dist_pot_donors, n_neighbors - 1, axis=1)[
160     :, :n_neighbors
161 ]
163 # Get weight matrix from distance matrix
164 donors_dist = dist_pot_donors[
165     np.arange(donors_idx.shape[0])[:, None], donors_idx
166 ]

```

File <\_\_array\_function\_\_ internals>:5, in argpartition(\*args, \*\*kwargs)

File ~\anaconda3\lib\site-packages\numpy\core\fromnumeric.py:839, in argpartition(a, kth, axis, kind, order)

```

763 @array_function_dispatch(_argpartition_dispatcher)
764 def argpartition(a, kth, axis=-1, kind='introspect', order=None):
765     """
766     Perform an indirect partition along the given axis using the
767     algorithm specified by the 'kind' keyword. It returns an array of
768     (...)
769     837
770     838     """
-> 839     return wrapfunc(a, 'argpartition', kth, axis=axis, kind=kind, order=order)

```

File ~\anaconda3\lib\site-packages\numpy\core\fromnumeric.py:57, in \_wrapfunc(obj, method, \*args, \*\*kwds)

```

54     return _wrapit(obj, method, *args, **kwds)
56 try:
-> 57     return bound(*args, **kwds)
58 except TypeError:
59     # A TypeError occurs if the object does have such a method in its
60     # class, but its signature is not identical to that of NumPy's. This
61     (...)
62     # Call _wrapit from within the except clause to ensure a potential
63     # exception has a traceback chain.
64     return _wrapit(obj, method, *args, **kwds)

```

KeyboardInterrupt:

## Handling Imbalanced data

In [ ]:

```
from imblearn.combine import SMOTETomek

# Resampling the minority class. The strategy can be changed as required.
smt = SMOTETomek(random_state=42,sampling_strategy='minority',n_jobs=-1)
# Fit the model to generate the data.
X_res, y_res = smt.fit_resample(X_knn, y)
```

## Initialize Default Models in a dictionary

In [ ]:

```
# Dictionary which contains models for experiment
models = {
    "Random Forest": RandomForestClassifier(),
    "Decision Tree": DecisionTreeClassifier(),
    "Gradient Boosting": GradientBoostingClassifier(),
    "Logistic Regression": LogisticRegression(),
    "K-Neighbors Classifier": KNeighborsClassifier(),
    "XGBClassifier": XGBClassifier(),
    "CatBoosting Classifier": CatBoostClassifier(verbose=False),
    "AdaBoost Classifier": AdaBoostClassifier()
}
```

## Fit KNN imputed data for models in dictionary

In [ ]:

```
report_knn = evaluate_models(X_res, y_res, models)
```

## Report for KNN Imputed data

In [ ]:

```
report_knn
```

## Insights

- For the Experiment 1: Knn imputer has XGBoost classifier as the best Model
- Proceeding with further experiments

## Experiment: 2 = Simple Imputer with Strategy Median

In [ ]:

```
num_features = X.select_dtypes(exclude="object").columns

# Fit the Simple imputer with strategy median
median_pipeline = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='median')),
    ('RobustScaler', RobustScaler())
])
```

In [ ]:

```
# Fit X with median pipeline
X_median = median_pipeline.fit_transform(X)
```

In [ ]:

```
# Resampling the minority class. The strategy can be changed as required.
smt = SMOTETomek(random_state=42,sampling_strategy='minority')
# Fit the model to generate the data.
X_res, y_res = smt.fit_resample(X_median, y)
```

In [ ]:

```
# Training the models
report_median = evaluate_models(X_res, y_res, models)
```

## Report for Simple Imputer with median strategy

In [ ]:

```
report_median
```

### Insights

- For the Experiment 2: Simple imputer with median strategy has Catboost classifier as the best Model
- Proceeding with further experiments

## Experiment: 3 = MICE for Imputing Null values

In [ ]:

```
import miceforest as mf

X_mice = X.copy()
kernel = mf.ImputationKernel(
    X_mice,
    save_all_iterations=True,
    random_state=1989
) # Run the MICE algorithm for 3 iterations kernel.mice(3)
```

In [ ]:

```
X_mice = kernel.complete_data()
```

In [ ]:

```
# fit robust scaler
mice_pipeline = Pipeline(steps=[
    ('RobustScaler', RobustScaler())
])
```

In [ ]:

```
# Fit X with Mice imputer
X_mice= mice_pipeline.fit_transform(X_mice)
```

In [ ]:

```
# Resampling the minority class. The strategy can be changed as required.
smt = SMOTETomek(random_state=42,sampling_strategy='minority', n_jobs=-1 )
# Fit the model to generate the data.
X_res, y_res = smt.fit_resample(X_mice, y)
```

In [ ]:

```
# Training the models
report_mice = evaluate_models(X_res, y_res, models)
```

## Report for MICE Imputer algorithm

In [ ]:

```
report_mice
```

... ..

## Insights

- For the Experiment 3: Mice imputer has XGBoost classifier as the best Model
- Proceeding with further experiments

## Experiment: 4 = Simple Imputer with Strategy Constant

In [ ]:

```
# Create a pipeline with simple imputer with strategy constant and fill value 0
constant_pipeline = Pipeline(steps=[
    ('Imputer', SimpleImputer(strategy='constant', fill_value=0)),
    ('RobustScaler', RobustScaler())
])
```

In [ ]:

```
X_const = constant_pipeline.fit_transform(X)
```

In [ ]:

```
# Resampling the minority class. The strategy can be changed as required.
smt = SMOTETomek(random_state=42, sampling_strategy='minority', n_jobs=-1 )
# Fit the model to generate the data.
X_res, y_res = smt.fit_resample(X_const, y)
```

In [ ]:

```
# training the models
report_const = evaluate_models(X_res, y_res, models)
```

## Report for Simple Imputer with Constant strategy

In [ ]:

```
report_const
```

## Insights

- For the Experiment 4: Simple imputer with constant strategy has XGBoost classifier as the best Model
- Proceeding with further experiments

## Experiment: 5 = Simple Imputer with Strategy Mean

In [ ]:

```
# Create a pipeline with Simple imputer with strategy mean
mean_pipeline = Pipeline(steps=[
    ('Imputer', SimpleImputer(strategy='mean')),
    ('RobustScaler', RobustScaler())
])
```

In [ ]:

```
X_mean = mean_pipeline.fit_transform(X)
```

In [ ]:

```
# Resampling the minority class. The strategy can be changed as required.
smt = SMOTETomek(random_state=42, sampling_strategy='minority' , n_jobs=-1)
# Fit the model to generate the data.
X_res, y_res = smt.fit_resample(X_mean, y)
```

In [ ]:

```
In [ ]:
# Training all models
report_mean = evaluate_models(X_res, y_res, models)
```

## Report for Simple imputer with strategy mean

```
In [ ]:
report_mean
```

## Experiment: 5 = Principle component analysis with imputing median

```
In [ ]:
from sklearn.decomposition import PCA
```

```
In [ ]:
pca_pipeline = Pipeline(steps=[
    ('Imputer', SimpleImputer(strategy='constant', fill_value=0)),
    ('RobustScaler', RobustScaler())
])
```

```
In [ ]:
X_pca = pca_pipeline.fit_transform(X)
```

```
In [ ]:
#Applying PCA
from sklearn.decomposition import PCA
var_ratio={}
for n in range(2,150):
    pc=PCA(n_components=n)
    df_pca=pc.fit(X_pca)
    var_ratio[n]=sum(df_pca.explained_variance_ratio_)
```

## Variance Plot

```
In [ ]:
# plotting variance ratio
pd.Series(var_ratio).plot()
```

## Kneed algorithm to find the elbow point

```
In [ ]:
from kneed import KneeLocator

i = np.arange(len(var_ratio))
variance_ratio= list(var_ratio.values())
components= list(var_ratio.keys())
knee = KneeLocator(i, variance_ratio, S=1, curve='concave', interp_method='polynomial')

fig = plt.figure(figsize=(5, 5))
knee.plot_knee()
plt.xlabel("Points")
plt.ylabel("Distance")
plt.show()
k= components[knee.knee]
print('Knee Locator k =', k)
```

```
In [ ]:
# Reducing the dimensions of the data
```

```
pca_final=PCA(n_components=18,random_state=42).fit(X_res)

reduced=pca_final.fit_transform(X_pca)
```

In [ ]:

```
# Resampling the minority class. The strategy can be changed as required.
smt = SMOTETomek(random_state=42,sampling_strategy='minority', n_jobs=-1)
# Fit the model to generate the data.
X_res, y_res = smt.fit_resample(reduced, y)
```

In [ ]:

```
# Training all models
report_pca = evaluate_models(X_res,y_res, models)
```

## Report for PCA and Mean imputed data

In [ ]:

```
report_pca
```

## Final Model

In [ ]:

```
from prettytable import PrettyTable

pt=PrettyTable()
pt.field_names=["Model","Imputation_method","Total_cost"]
pt.add_row(["XGBClassifier","Simple Imputer-Constant","2950"])
pt.add_row(["XGBClassifier","Mice","3510"])
pt.add_row(["XGBClassifier","Knn-Imputer","4460"])
pt.add_row(["XGBClassifier","Simple Imputer-Mean","4950"])
pt.add_row(["CatBoostClassifier","Median","5760"])
pt.add_row(["Random Forest","PCA","34150"])
print(pt)
```

## Report

- From the final report we can see than XGBClassifier with Simple imputer with strategy constant has performed the best with cost of 2950

## Fitting the Final model and get reports

In [ ]:

```
final_model = XGBClassifier()

# Resampling the minority class. The strategy can be changed as required.
smt = SMOTETomek(random_state=42,sampling_strategy='minority', n_jobs=-1)
# Fit the model to generate the data.
X_res, y_res = smt.fit_resample(X_const, y)
```

In [ ]:

```
X_train, X_test, y_train, y_test = train_test_split(X_res,y_res,test_size=0.2,random_state=42)

final_model = final_model.fit(X_train, y_train)
y_pred = final_model.predict(X_test)
```

In [ ]:

```
print("Final XGBoost Classifier Accuracy Score (Train) :", final_model.score(X_train,y_train))
```



```
ain))  
print("Final XGBoost Classifier Accuracy Score (Test) :", accuracy_score(y_pred,y_test))
```

In [ ]:

```
print("Final XGBoost Classifier Cost Metric(Test) :",total_cost(y_test, y_pred))
```

In [ ]:

```
from sklearn.metrics import plot_confusion_matrix
```

```
#plots Confusion matrix
```

```
plot_confusion_matrix(final_model, X_test, y_test, cmap='Blues', values_format='d')
```

**The best Model is XGBoost Classifier with 99.6% accuracy and cost of 2950**

In [ ]:

In [ ]: