

▼ Mercedes-Benz Greener Manufacturing

▼ Importing Library

```
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn import preprocessing
import xgboost as xgb
color = sns.color_palette()

%matplotlib inline

pd.options.mode.chained_assignment = None # default='warn'
pd.options.display.max_columns = 999

from subprocess import check_output

↳ /usr/local/lib/python3.6/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning:
      import pandas.util.testing as tm

print(check_output(["ls"]).decode("utf8"))

↳ sample_data
    test.csv
    train.csv
```

▼ Importing Dataset

```
train_df = pd.read_csv("train.csv")
test_df = pd.read_csv("test.csv")
print("Train shape : ", train_df.shape)
print("Test shape : ", test_df.shape)

↳ Train shape : (4209, 378)
    Test shape : (4209, 377)
```

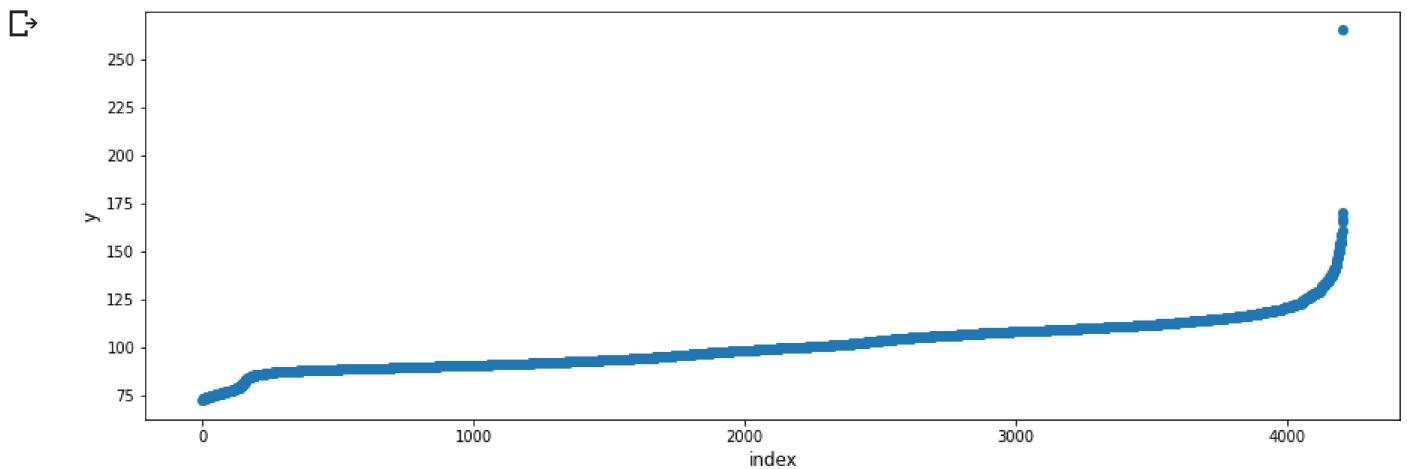
```
train_df.head()
```

```
train_df.groupby("X0").aggregate('count').reset_index()
```



| | ID | y | x0 | x1 | x2 | x3 | x4 | x5 | x6 | x8 | x10 | x11 | x12 | x13 | x14 | x15 | x16 | x17 | x |
|---|----|--------|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| 0 | 0 | 130.81 | k | v | at | a | d | u | j | o | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | |
| 1 | 6 | 88.53 | k | t | av | e | d | y | l | o | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 2 | 7 | 76.26 | az | w | n | c | d | x | j | x | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | |
| 3 | 9 | 80.62 | az | t | n | f | d | x | l | e | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 4 | 13 | 78.02 | az | v | n | f | d | h | d | n | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |

```
plt.figure(figsize=(15,5))
plt.scatter(range(train_df.shape[0]), np.sort(train_df.y.values))
plt.xlabel('index', fontsize=12)
plt.ylabel('y', fontsize=12)
plt.show()
```



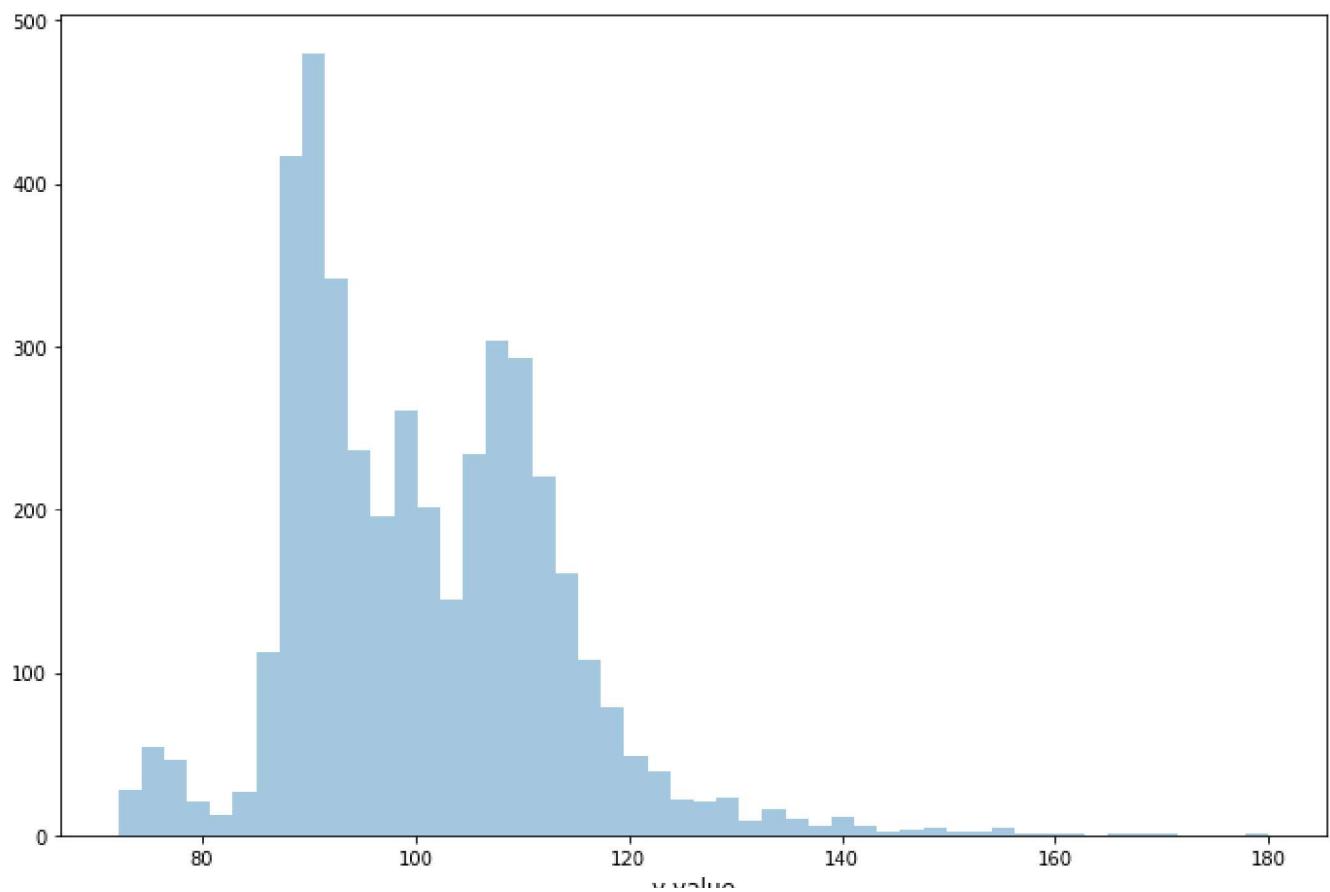
```
train_df['y'].max()
```

265.32

```
ulimit = 180
train_df['y'].loc[train_df['y']>ulimit] = ulimit

plt.figure(figsize=(12,8))
sns.distplot(a=train_df.y.values, bins=50, kde=False)
plt.xlabel('y value', fontsize=12)
plt.show()
```

→



▼ Data_Type For each Column

```
dtype_df = train_df.dtypes.reset_index()
dtype_df.columns = ["Count", "Column Type"]
dtype_df.groupby("Column Type").aggregate('count').reset_index()
```

| Column Type | Count |
|-------------|-------|
| int64 | 369 |
| float64 | 1 |
| object | 8 |

```
dtype_df.loc[:10,:]
```

```
→
```

| | Count | Column | Type |
|---|-------|--------|---------|
| 0 | ID | | int64 |
| 1 | y | | float64 |
| 2 | X0 | | object |
| 3 | X1 | | object |
| 4 | X2 | | object |
| 5 | X3 | | object |

▼ Missing values

```
missing_df = train_df.isnull().sum(axis=0).reset_index()
missing_df.columns = ['column_name', 'missing_count']
missing_df = missing_df.loc[missing_df['missing_count']>0]
missing_df = missing_df.sort_values(by='missing_count')
missing_df
```

⇨ column_name missing_count

▼ Integer Columns

```
unique_values_dict = {}
for col in train_df.columns:
    if col not in ["ID", "y", "X0", "X1", "X2", "X3", "X4", "X5", "X6", "X8"]:
        unique_value = str(np.sort(train_df[col].unique()).tolist())
        tlist = unique_values_dict.get(unique_value,[])
        tlist.append(col)
        unique_values_dict[unique_value] = tlist[:]
for unique_val, columns in unique_values_dict.items():
    print("Columns containing the unique values : ",unique_val)
    print(columns)
    print("-----")
```

⇨ Columns containing the unique values : [0, 1]

 Columns containing the unique values : [0]

▼ Categorical Columns

The number of 0's and 1's in each of these variables.

```
zero_count_list = []
one_count_list = []
cols_list = unique_values_dict['[0, 1]']
for col in cols_list:
    zero_count_list.append((train_df[col]==0).sum())
    one_count_list.append((train_df[col]==1).sum())
N = len(cols_list)
ind = np.arange(N)

width = 0.35
plt.figure(figsize=(6,100))
p1 = plt.barh(ind, zero_count_list, width, color='red')
p2 = plt.barh(ind, one_count_list, width, left=zero_count_list, color="blue")
plt.yticks(ticks=ind,labels=cols_list)
plt.legend((p1[0], p2[0]), ('Zero count', 'One Count'))
plt.show()
```



```
lbl.fit(list(train_df[f].values))
train_df[f] = lbl.transform(list(train_df[f].values))

train_y = train_df['y'].values
train_X = train_df.drop(["ID", "y", "eval_set"], axis=1)

# Thanks to anokas for this #
def xgb_r2_score(preds, dtrain):
    labels = dtrain.get_label()
    return 'r2', r2_score(labels, preds)

xgb_params = {
    'eta': 0.05,
    'max_depth': 6,
    'subsample': 0.7,
    'colsample_bytree': 0.7,
    'objective': 'reg:linear',
    'silent': 1
}
dtrain = xgb.DMatrix(train_X, train_y, feature_names=train_X.columns.values)
model = xgb.train(dict(xgb_params, silent=0), dtrain, num_boost_round=100, feval=xgb_r2_score

# plot the important features #
fig, ax = plt.subplots(figsize=(12,18))
xgb.plot_importance(model, max_num_features=50, height=0.8, ax=ax)
plt.show()
```



| | X0 | ID | y | X1 | X2 | X3 | X4 | X5 | X6 | X8 | X10 | X11 | X12 | X13 | X14 | X15 | X16 |
|----|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0 | a | 21 | 21 | 21 | 21 | 21 | 21 | 21 | 21 | 21 | 21 | 21 | 21 | 21 | 21 | 21 | |
| 1 | aa | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | |
| 2 | ab | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
| 3 | ac | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
| 4 | ad | 14 | 14 | 14 | 14 | 14 | 14 | 14 | 14 | 14 | 14 | 14 | 14 | 14 | 14 | 14 | |
| 5 | af | 35 | 35 | 35 | 35 | 35 | 35 | 35 | 35 | 35 | 35 | 35 | 35 | 35 | 35 | 35 | |
| 6 | ai | 34 | 34 | 34 | 34 | 34 | 34 | 34 | 34 | 34 | 34 | 34 | 34 | 34 | 34 | 34 | |
| 7 | aj | 151 | 151 | 151 | 151 | 151 | 151 | 151 | 151 | 151 | 151 | 151 | 151 | 151 | 151 | 151 | |
| 8 | ak | 349 | 349 | 349 | 349 | 349 | 349 | 349 | 349 | 349 | 349 | 349 | 349 | 349 | 349 | 349 | |
| 9 | al | 67 | 67 | 67 | 67 | 67 | 67 | 67 | 67 | 67 | 67 | 67 | 67 | 67 | 67 | 67 | |
| 10 | am | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | |
| 11 | ao | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | |
| 12 | ap | 103 | 103 | 103 | 103 | 103 | 103 | 103 | 103 | 103 | 103 | 103 | 103 | 103 | 103 | 103 | |
| 13 | aq | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 | |
| 14 | as | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | |
| 15 | at | 25 | 25 | 25 | 25 | 25 | 25 | 25 | 25 | 25 | 25 | 25 | 25 | 25 | 25 | 25 | |
| 16 | au | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | 11 | |
| 17 | aw | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | 16 | |
| 18 | ax | 19 | 19 | 19 | 19 | 19 | 19 | 19 | 19 | 19 | 19 | 19 | 19 | 19 | 19 | 19 | |

column -> X0

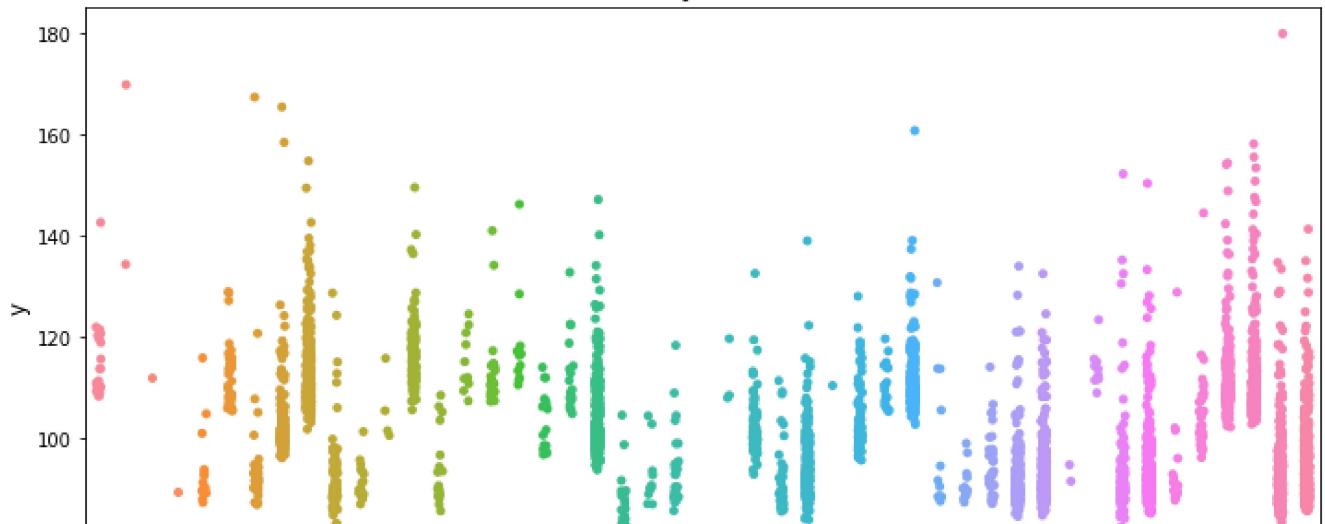
```

20  27  175  175  175  175  175  175  175  175  175  175  175  175  175  175  175
var_name = "X0"
col_order = np.sort(train_df[var_name].unique()).tolist()
plt.figure(figsize=(12,6))
sns.stripplot(x=var_name, y='y', data=train_df, order=col_order)
plt.xlabel(var_name, fontsize=12)
plt.ylabel('y', fontsize=12)
plt.title("Distribution of y variable with "+var_name, fontsize=15)
plt.show()

```



Distribution of y variable with X0



column -> X1

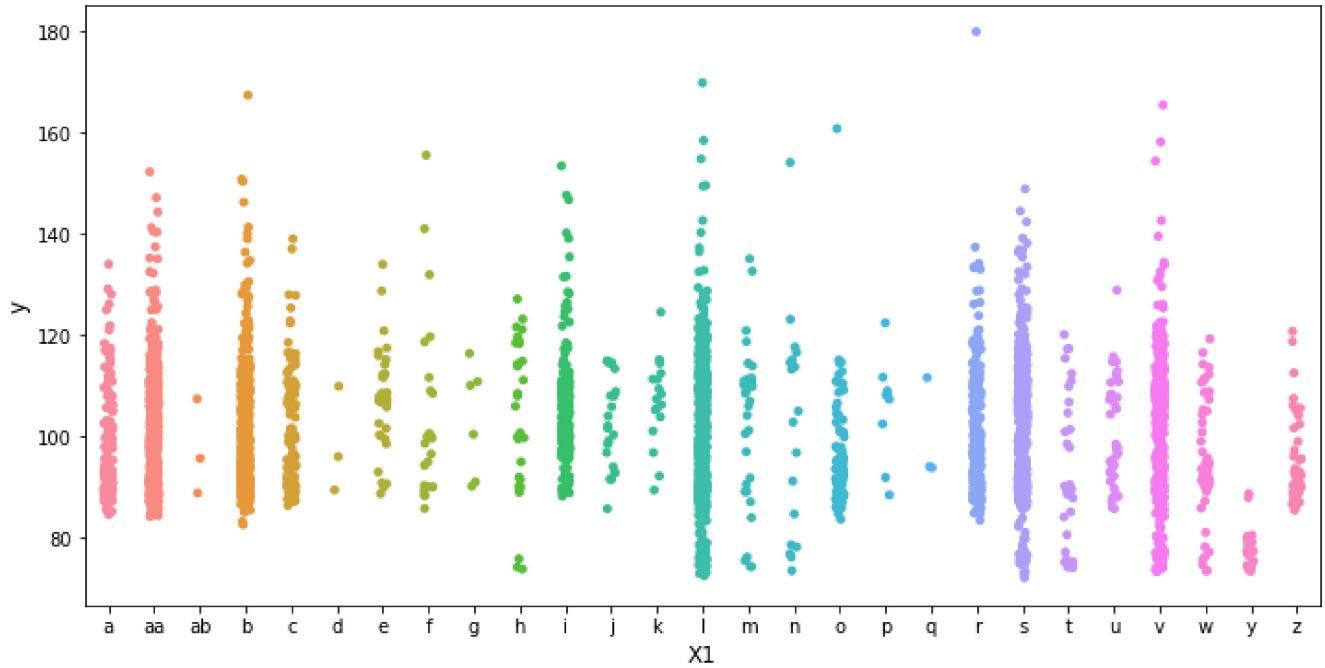
```

a aa ab ac ad af ai ai ak al am ao ap as at au aw ax ay az b ba bc c d e f a h i i k l m n o a r s t u v w x y z
var_name = "X1"
col_order = np.sort(train_df[var_name].unique()).tolist()
plt.figure(figsize=(12,6))
sns.stripplot(x=var_name, y='y', data=train_df, order=col_order)
plt.xlabel(var_name, fontsize=12)
plt.ylabel('y', fontsize=12)
plt.title("Distribution of y variable with "+var_name, fontsize=15)
plt.show()

```



Distribution of y variable with X1

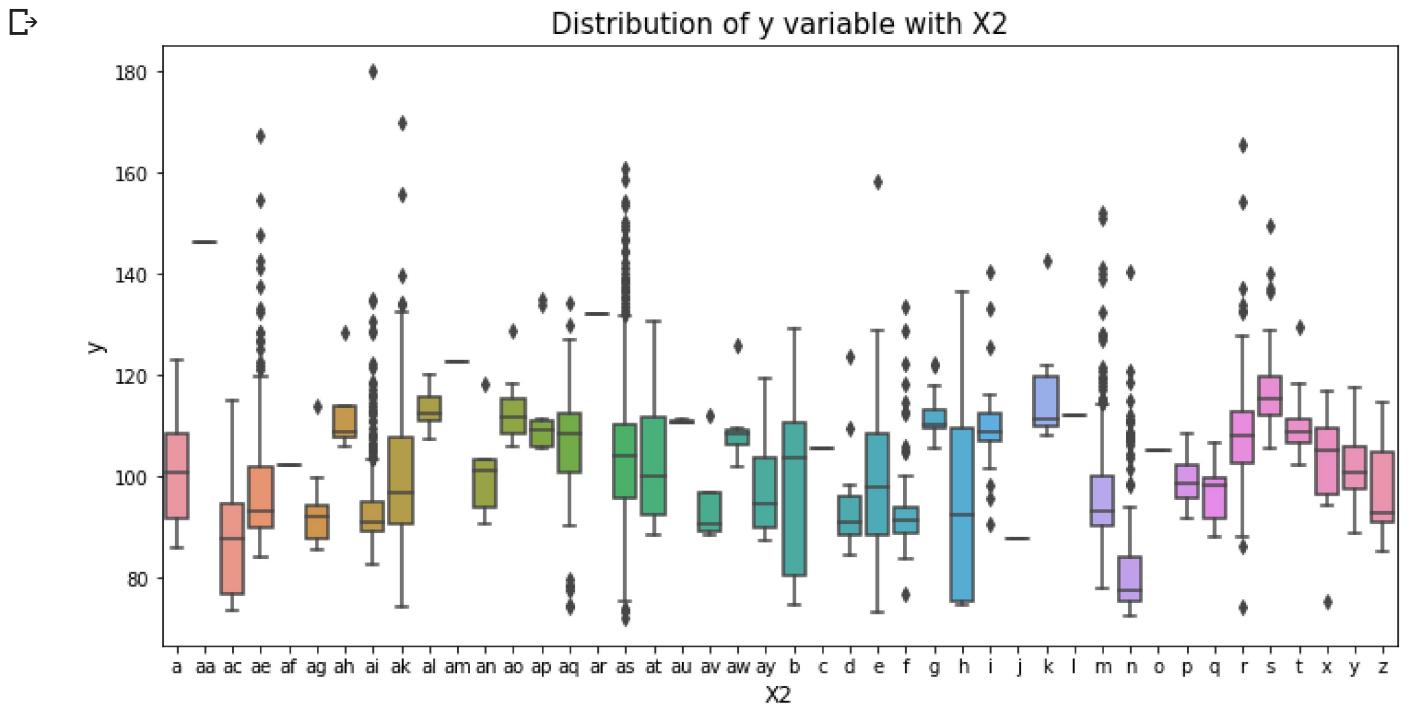


column -> X2

```

var_name = "X2"
col_order = np.sort(train_df[var_name].unique()).tolist()
plt.figure(figsize=(12,6))
sns.boxplot(x=var_name, y='y', data=train_df, order=col_order)
plt.xlabel(var_name, fontsize=12)
plt.ylabel('y', fontsize=12)
plt.title("Distribution of y variable with "+var_name, fontsize=15)
plt.show()

```



column -> X3

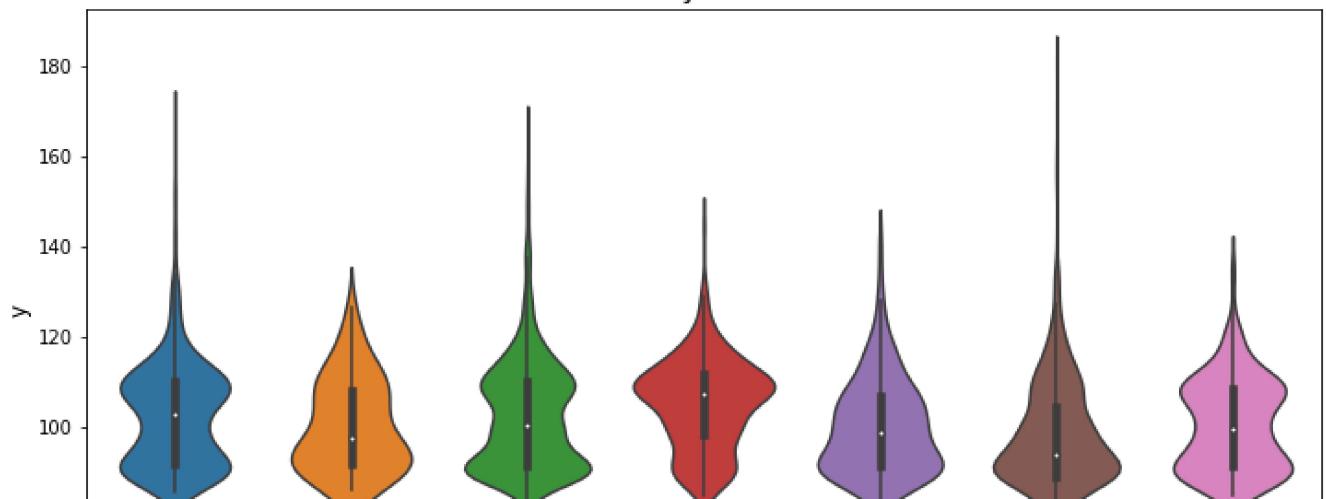
```

var_name = "X3"
col_order = np.sort(train_df[var_name].unique()).tolist()
plt.figure(figsize=(12,6))
sns.violinplot(x=var_name, y='y', data=train_df, order=col_order)
plt.xlabel(var_name, fontsize=12)
plt.ylabel('y', fontsize=12)
plt.title("Distribution of y variable with "+var_name, fontsize=15)
plt.show()

```

↪

Distribution of y variable with X3



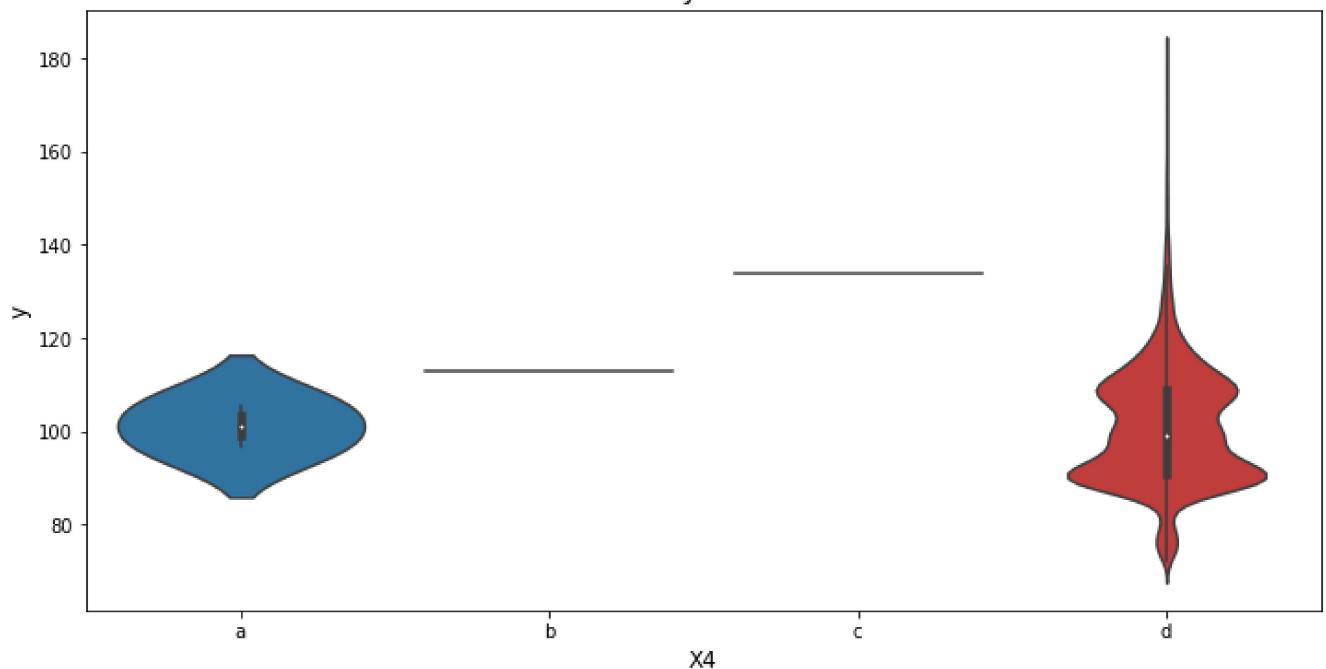
column -> X4

```
var_name = "X4"
```

```
col_order = np.sort(train_df[var_name].unique()).tolist()
plt.figure(figsize=(12,6))
sns.violinplot(x=var_name, y='y', data=train_df, order=col_order)
plt.xlabel(var_name, fontsize=12)
plt.ylabel('y', fontsize=12)
plt.title("Distribution of y variable with "+var_name, fontsize=15)
plt.show()
```

⇨

Distribution of y variable with X4

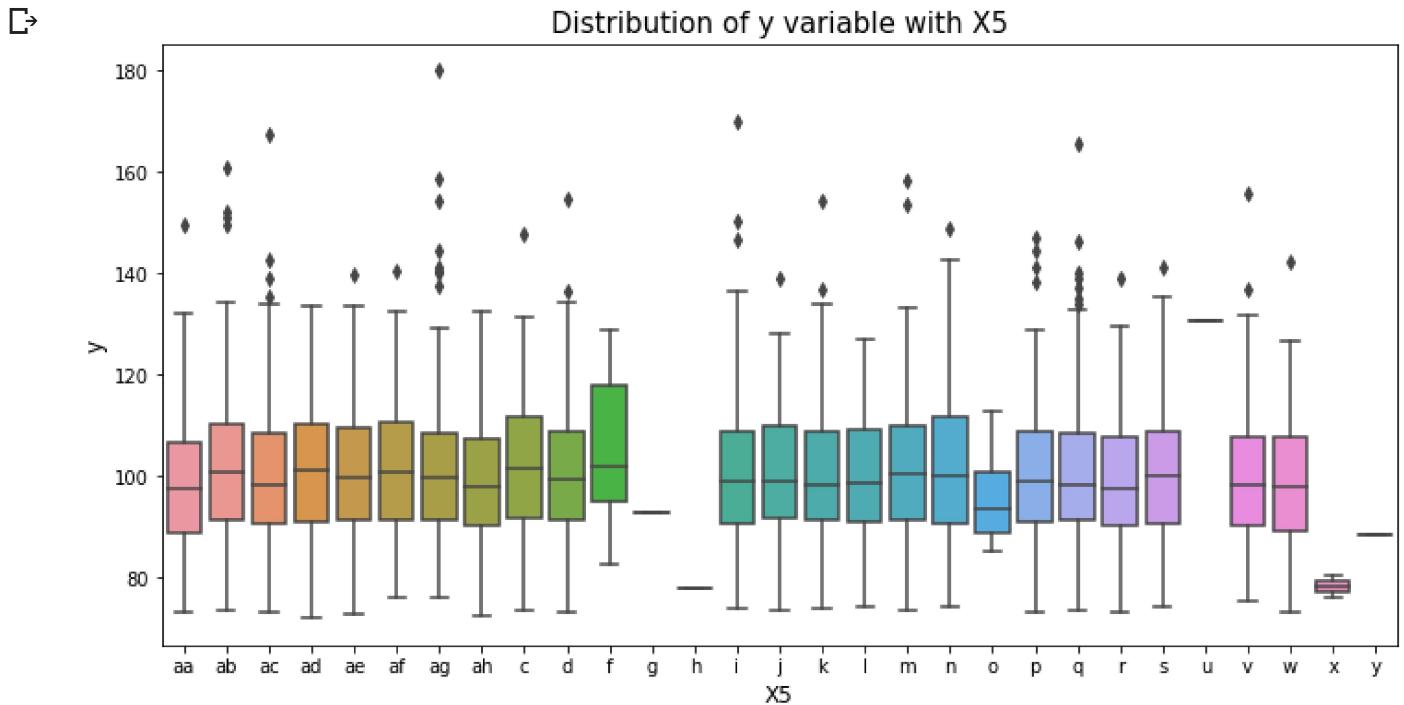


column -> X5

```

var_name = "X5"
col_order = np.sort(train_df[var_name].unique()).tolist()
plt.figure(figsize=(12,6))
sns.boxplot(x=var_name, y='y', data=train_df, order=col_order)
plt.xlabel(var_name, fontsize=12)
plt.ylabel('y', fontsize=12)
plt.title("Distribution of y variable with "+var_name, fontsize=15)
plt.show()

```



column -> X6

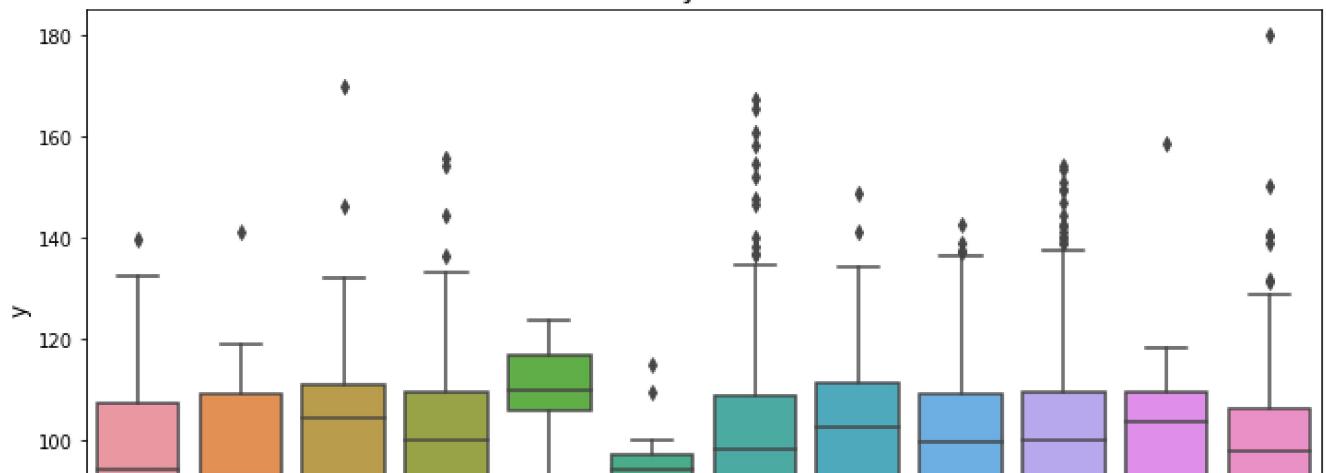
```

var_name = "X6"
col_order = np.sort(train_df[var_name].unique()).tolist()
plt.figure(figsize=(12,6))
sns.boxplot(x=var_name, y='y', data=train_df, order=col_order)
plt.xlabel(var_name, fontsize=12)
plt.ylabel('y', fontsize=12)
plt.title("Distribution of y variable with "+var_name, fontsize=15)
plt.show()

```

↪

Distribution of y variable with X6



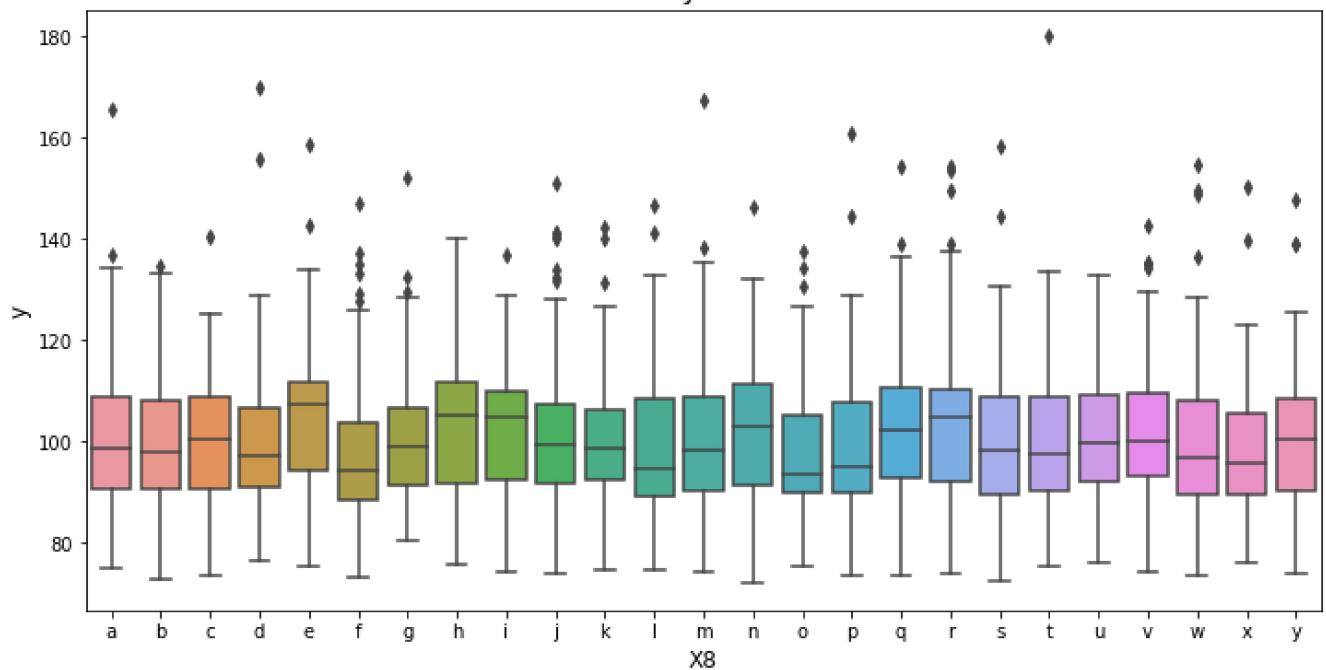
column -> X8

```

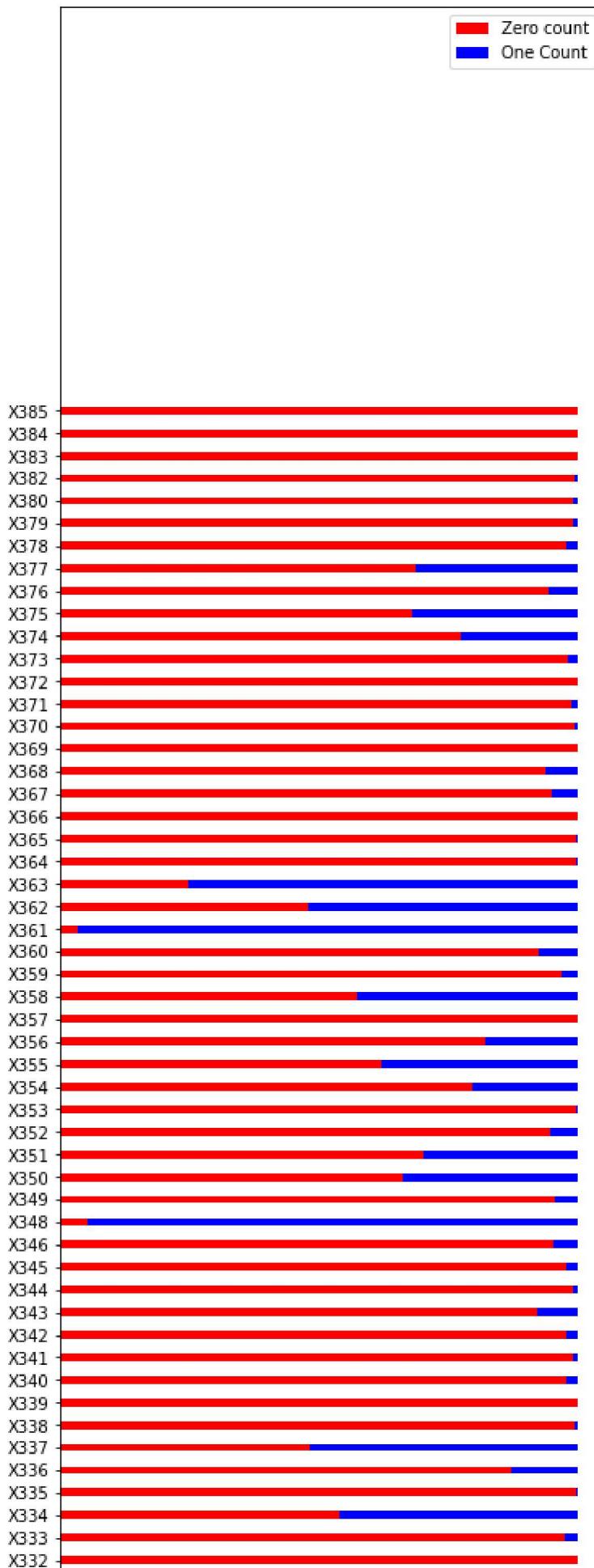
var_name = "X8"
col_order = np.sort(train_df[var_name].unique()).tolist()
plt.figure(figsize=(12,6))
sns.boxplot(x=var_name, y='y', data=train_df, order=col_order)
plt.xlabel(var_name, fontsize=12)
plt.ylabel('y', fontsize=12)
plt.title("Distribution of y variable with "+var_name, fontsize=15)
plt.show()

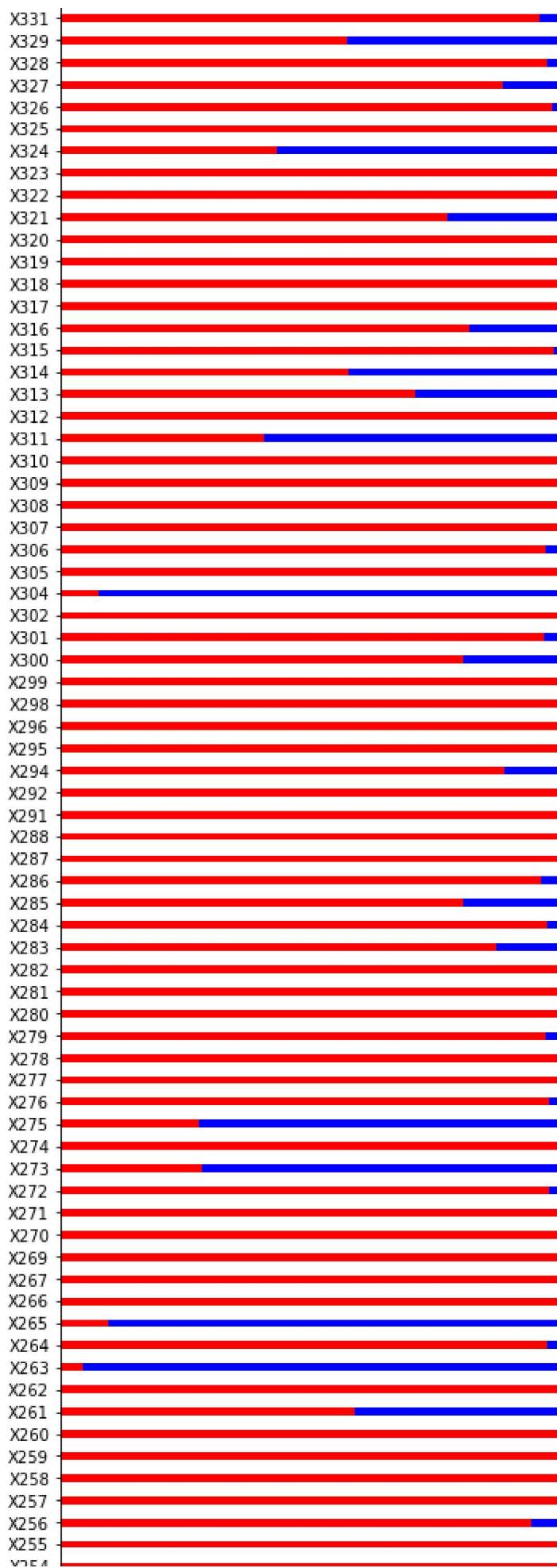
```

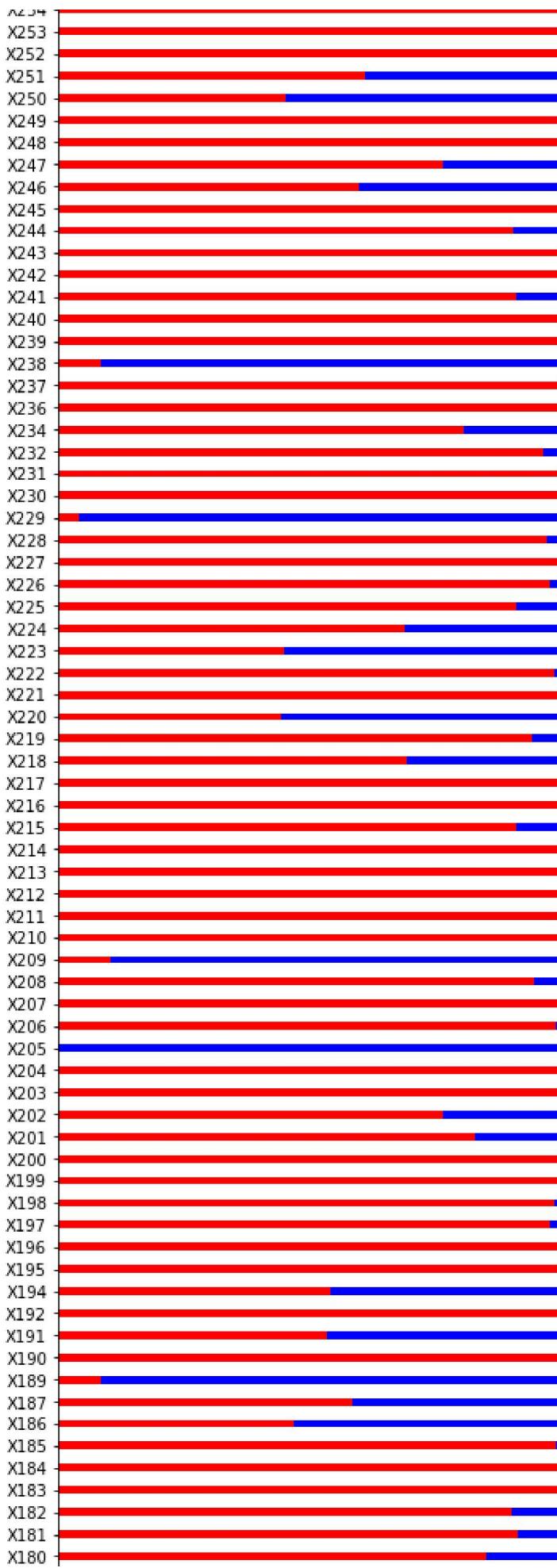
Distribution of y variable with X8

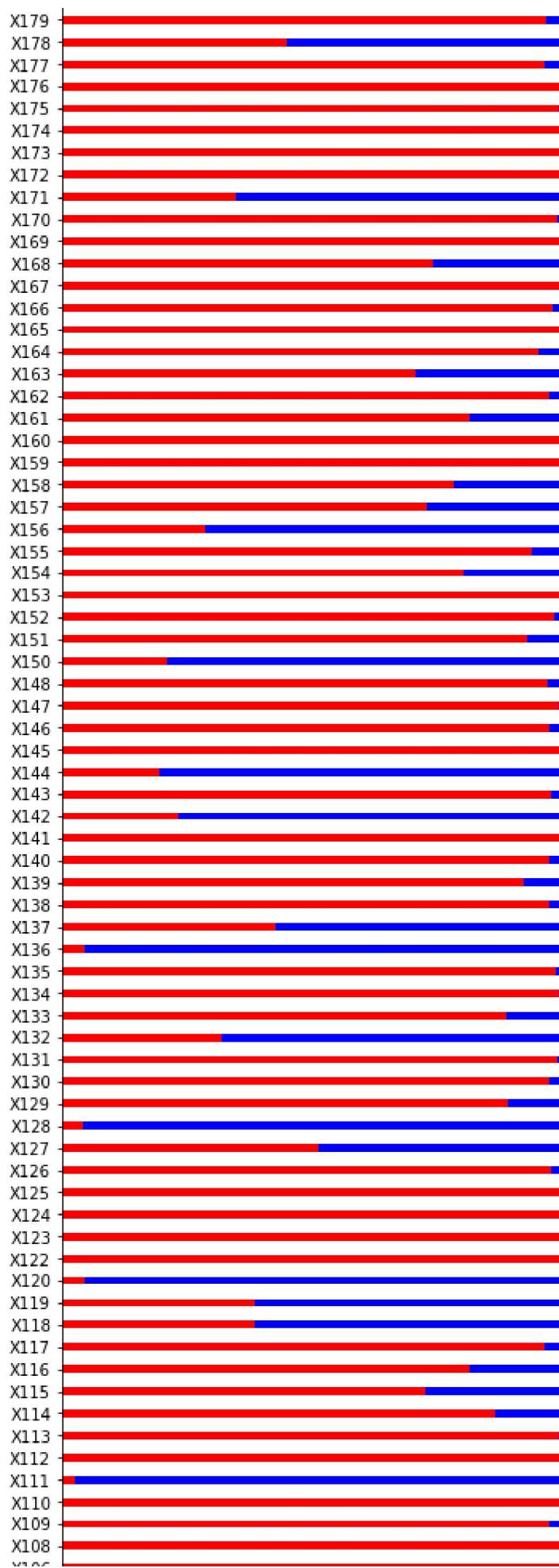


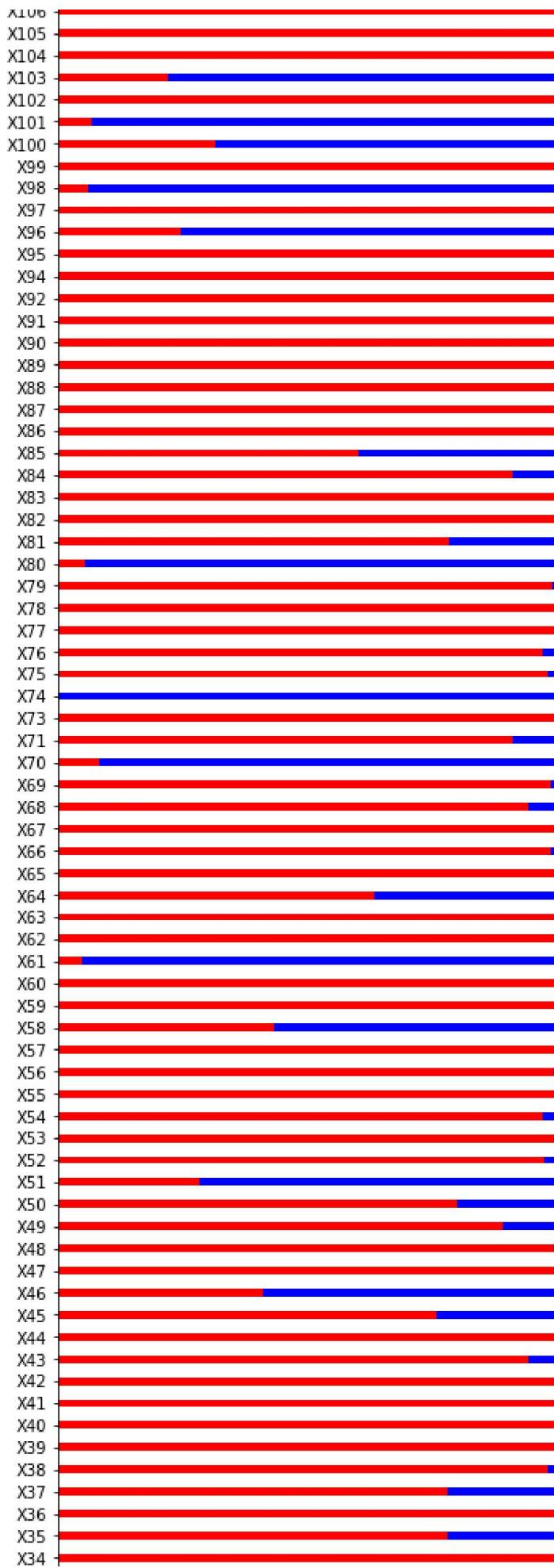
▼ Binary Variables











X33 |

The mean y value in each of the binary variable

X29 |

```
zero_mean_list = []
one_mean_list = []
cols_list = unique_values_dict['[0, 1]']
for col in cols_list:
    zero_mean_list.append(train_df.loc[train_df[col]==0].y.mean())
    one_mean_list.append(train_df.loc[train_df[col]==1].y.mean())
```

X16 |

```
new_df = pd.DataFrame({"column_name":cols_list+cols_list, "value": [0]*len(cols_list) + [1]*len(cols_list),
                      "y_mean":zero_mean_list+one_mean_list})
new_df = new_df.pivot(index='column_name', columns='value', values='y_mean')
new_df
```

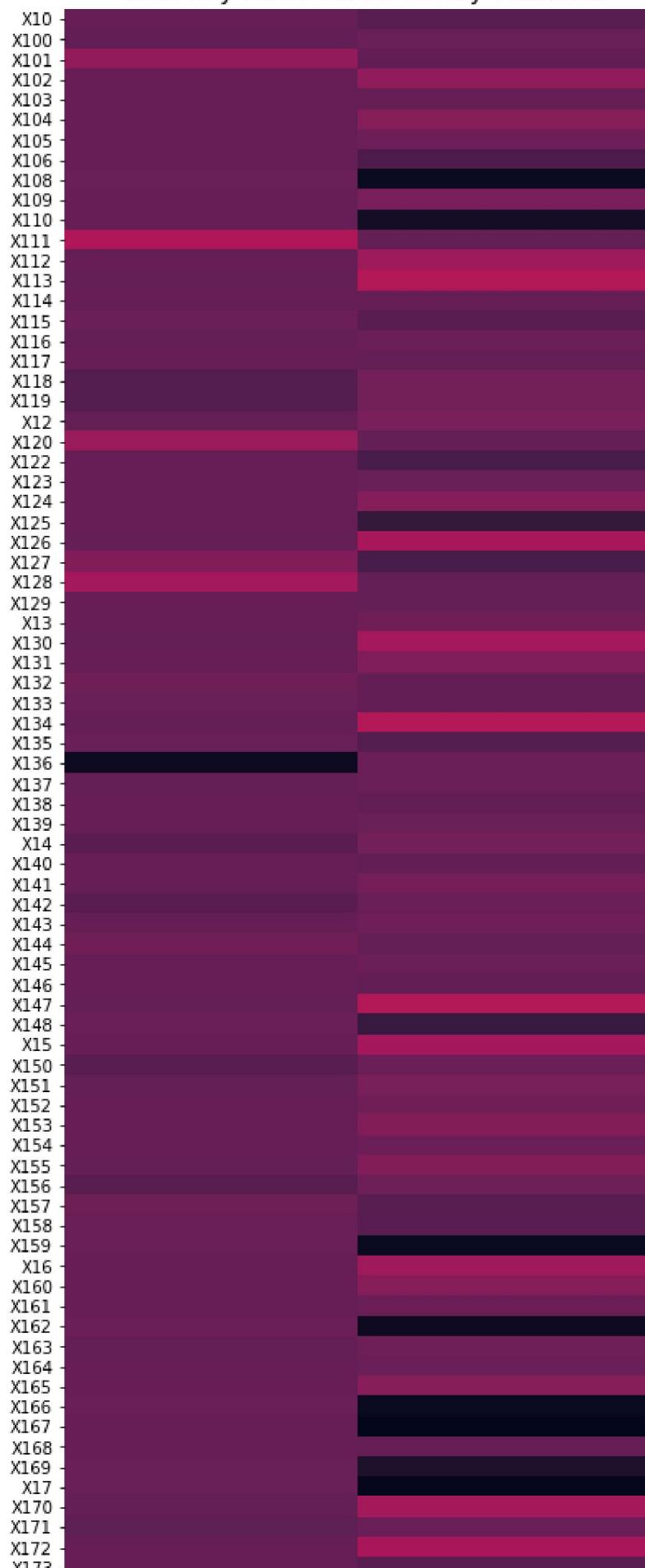
| | value | 0 | 1 |
|-------------|------------|------------|-----|
| column_name | | | |
| X10 | 100.688500 | 97.723214 | |
| X100 | 99.763083 | 101.046740 | |
| X101 | 110.198376 | 99.991894 | |
| X102 | 100.586823 | 109.617931 | |
| X103 | 100.472172 | 100.697632 | |
| ... | ... | ... | ... |
| X95 | 100.644223 | 120.950000 | |
| X96 | 97.968792 | 101.504108 | |
| X97 | 100.676829 | 94.180556 | |
| X98 | 110.267934 | 100.062264 | |
| X99 | 100.709446 | 93.647778 | |

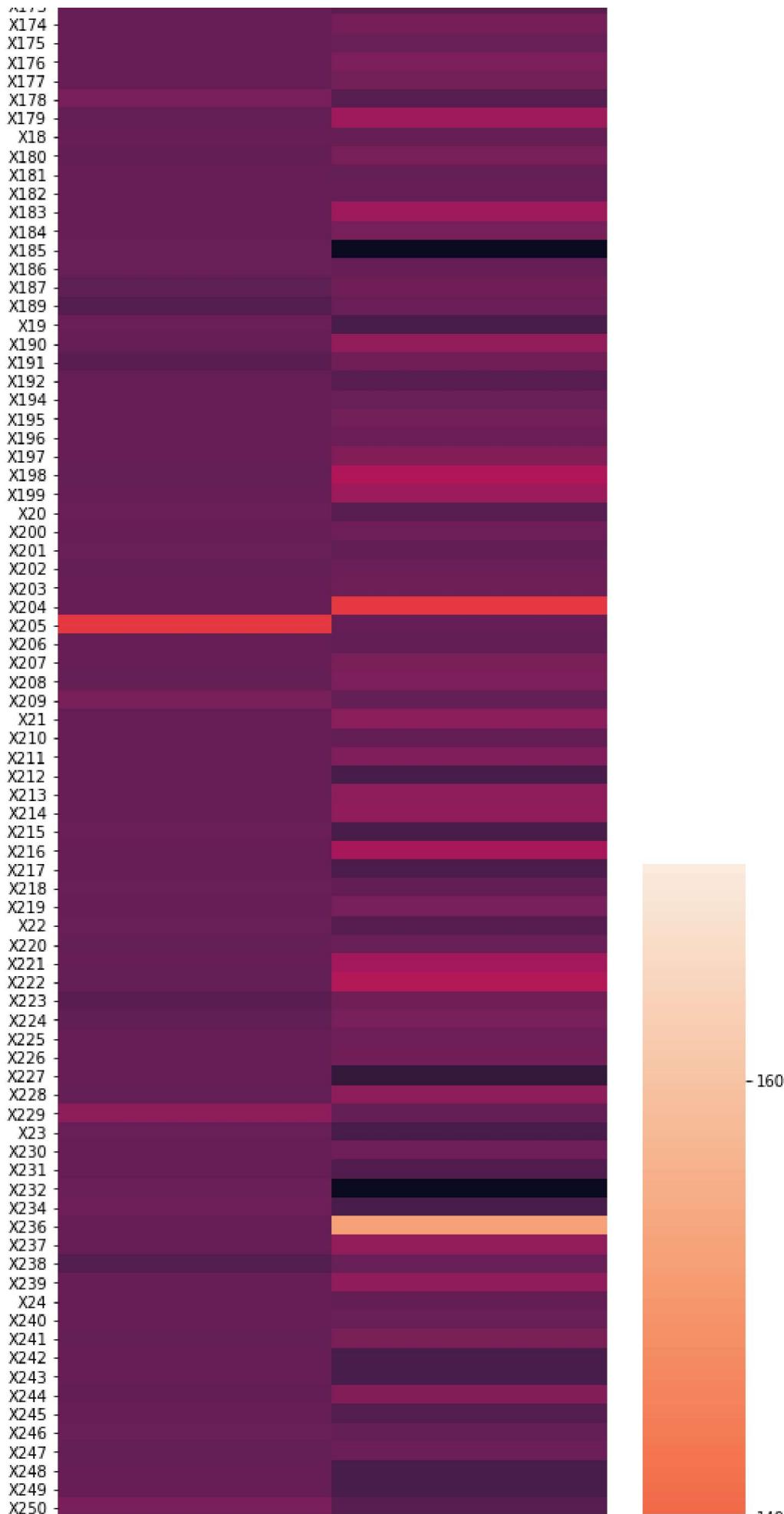
356 rows × 2 columns

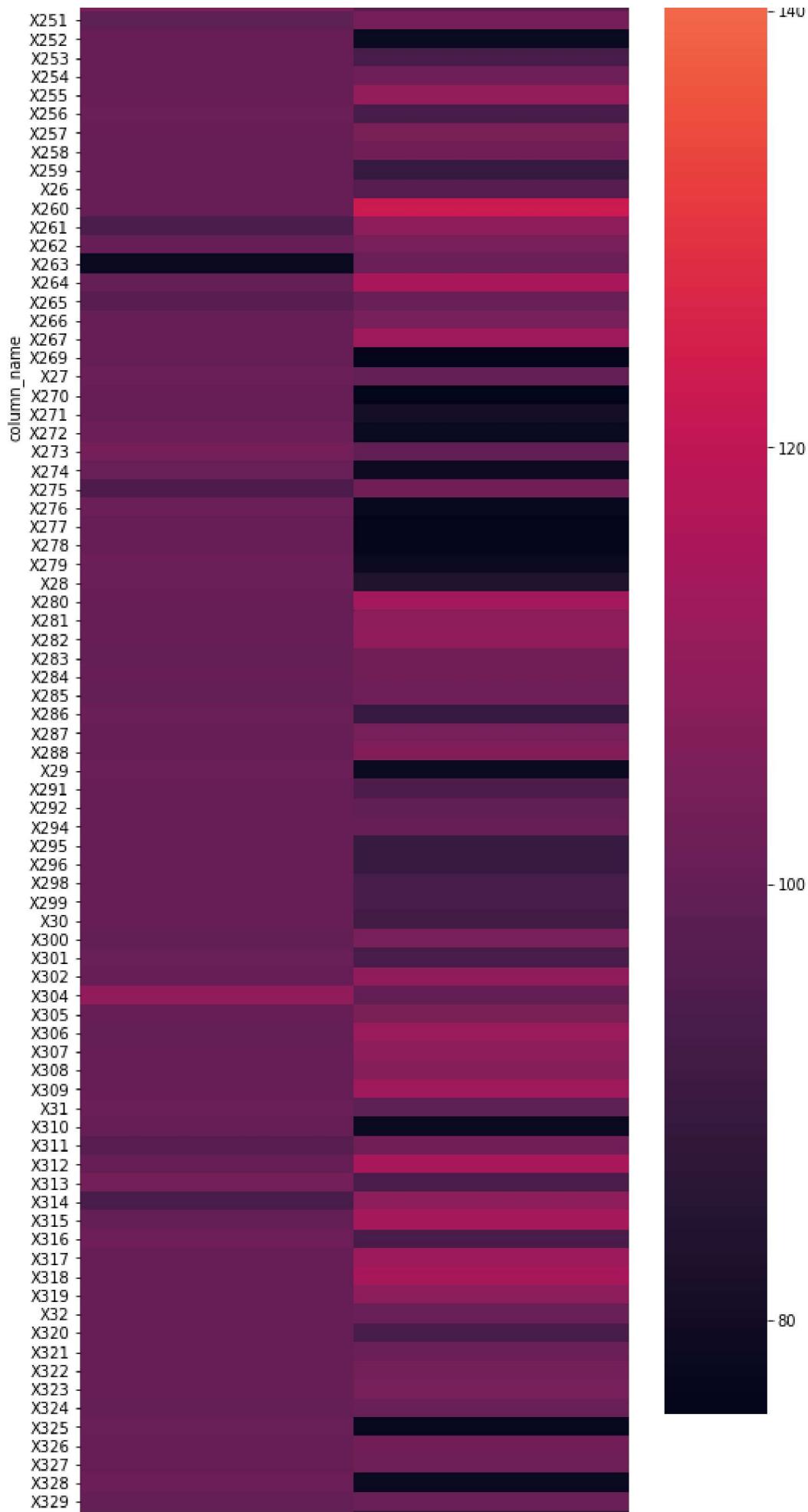
```
plt.figure(figsize=(8,80))
sns.heatmap(new_df)
plt.title("Mean of y value across binary variables", fontsize=15)
plt.show()
```

C→

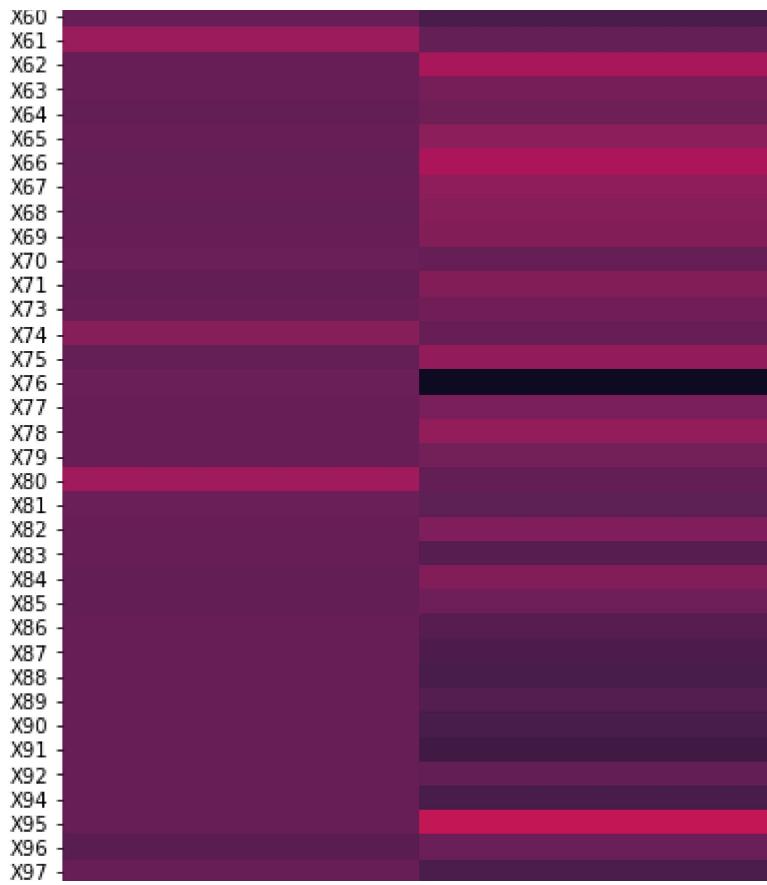
Mean of y value across binary variables











▼ ID variable

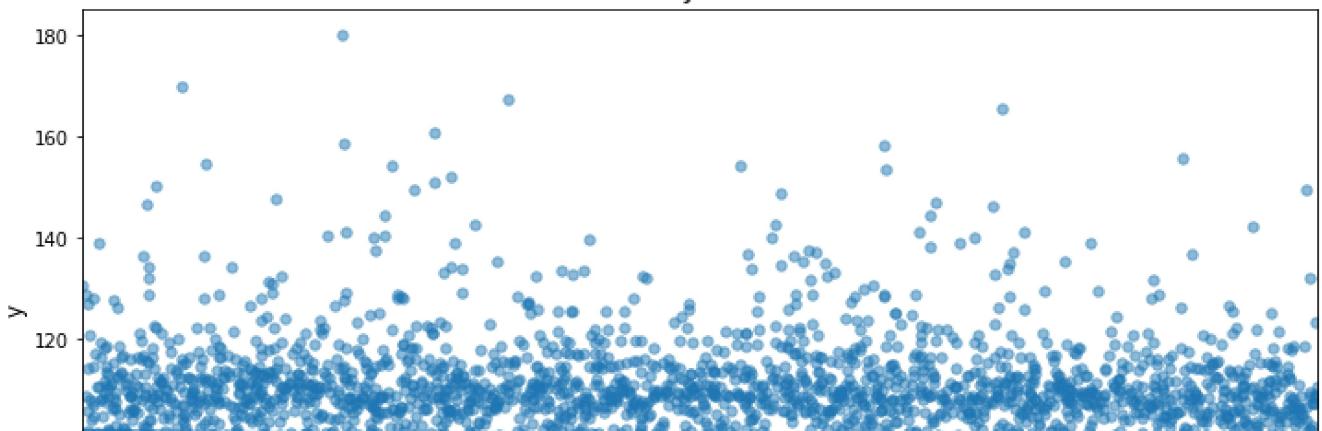
One more important thing we need to look at it is ID variable. This will give an idea of how the splits are done across train and test (random or id based) and also to help see if ID has some potential prediction capability (probably not so useful for business)

Let us first see how the 'y' variable changes with ID variable.

```
var_name = "ID"
plt.figure(figsize=(12,6))
sns.regplot(x=var_name, y='y', data=train_df, scatter_kws={'alpha':0.5, 's':30})
plt.xlabel(var_name, fontsize=12)
plt.ylabel('y', fontsize=12)
plt.title("Distribution of y variable with "+var_name, fontsize=15)
plt.show()
```



Distribution of y variable with ID



There seems to be a slight decreasing trend with respect to ID variable. Now let us see how the IDs are distributed across train and test.

```
|   .   .   .   .   .   .   .   .   .   .   .   .   .   .   .   .   .   .   .   .   . |
```

```
plt.figure(figsize=(6,10))
train_df['eval_set'] = "train"
test_df['eval_set'] = "test"
full_df = pd.concat([train_df[["ID","eval_set"]], test_df[["ID","eval_set"]]], axis=0)
full_df
```

| | ID | eval_set |
|------|------|----------|
| 0 | 0 | train |
| 1 | 6 | train |
| 2 | 7 | train |
| 3 | 9 | train |
| 4 | 13 | train |
| ... | ... | ... |
| 4204 | 8410 | test |
| 4205 | 8411 | test |
| 4206 | 8413 | test |
| 4207 | 8414 | test |
| 4208 | 8416 | test |

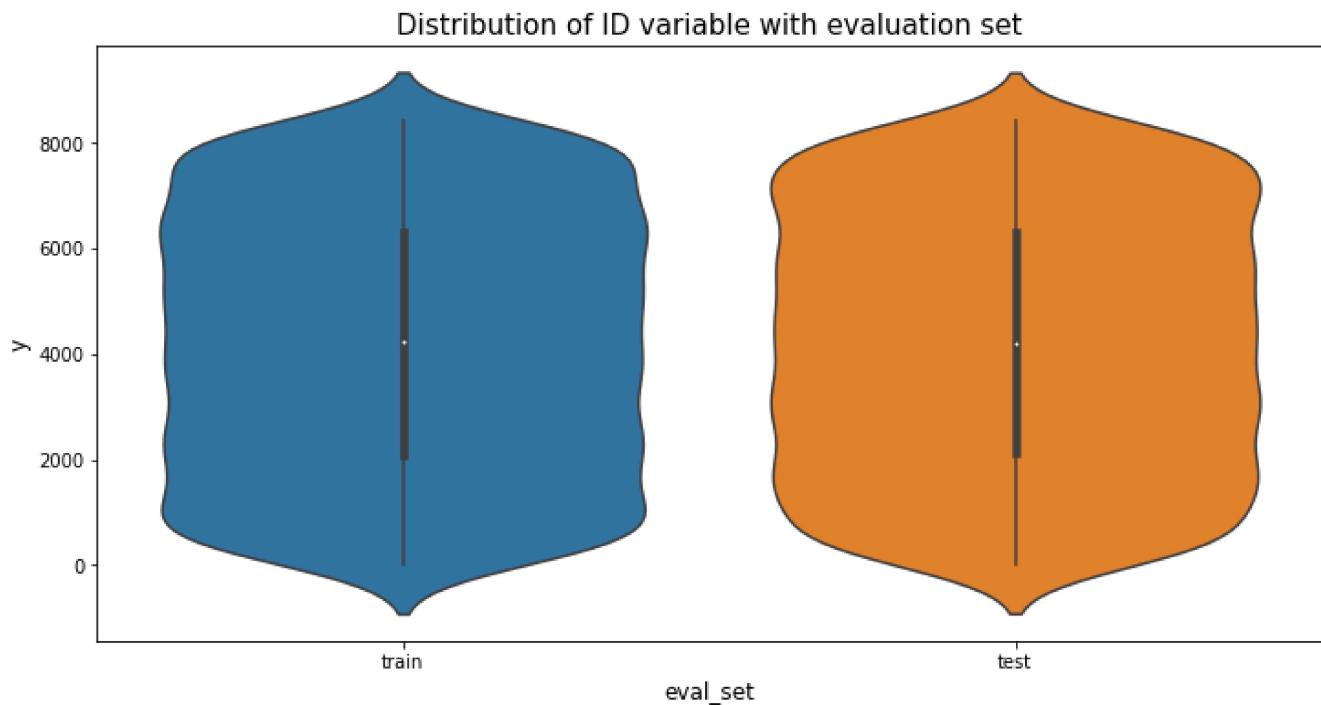
8418 rows × 2 columns

<Figure size 432x720 with 0 Axes>

```
plt.figure(figsize=(12,6))
sns.violinplot(x="eval_set", y='ID', data=full_df)
plt.xlabel("eval_set", fontsize=12)
plt.ylabel('y', fontsize=12)
```

```
plt.title("Distribution of ID variable with evaluation set", fontsize=15)
plt.show()
```

↳



Seems like a random split of ID variable between train and test samples.

▼ XGBoost

Important Variables:

Now let us run an xgboost model to get the important variables.

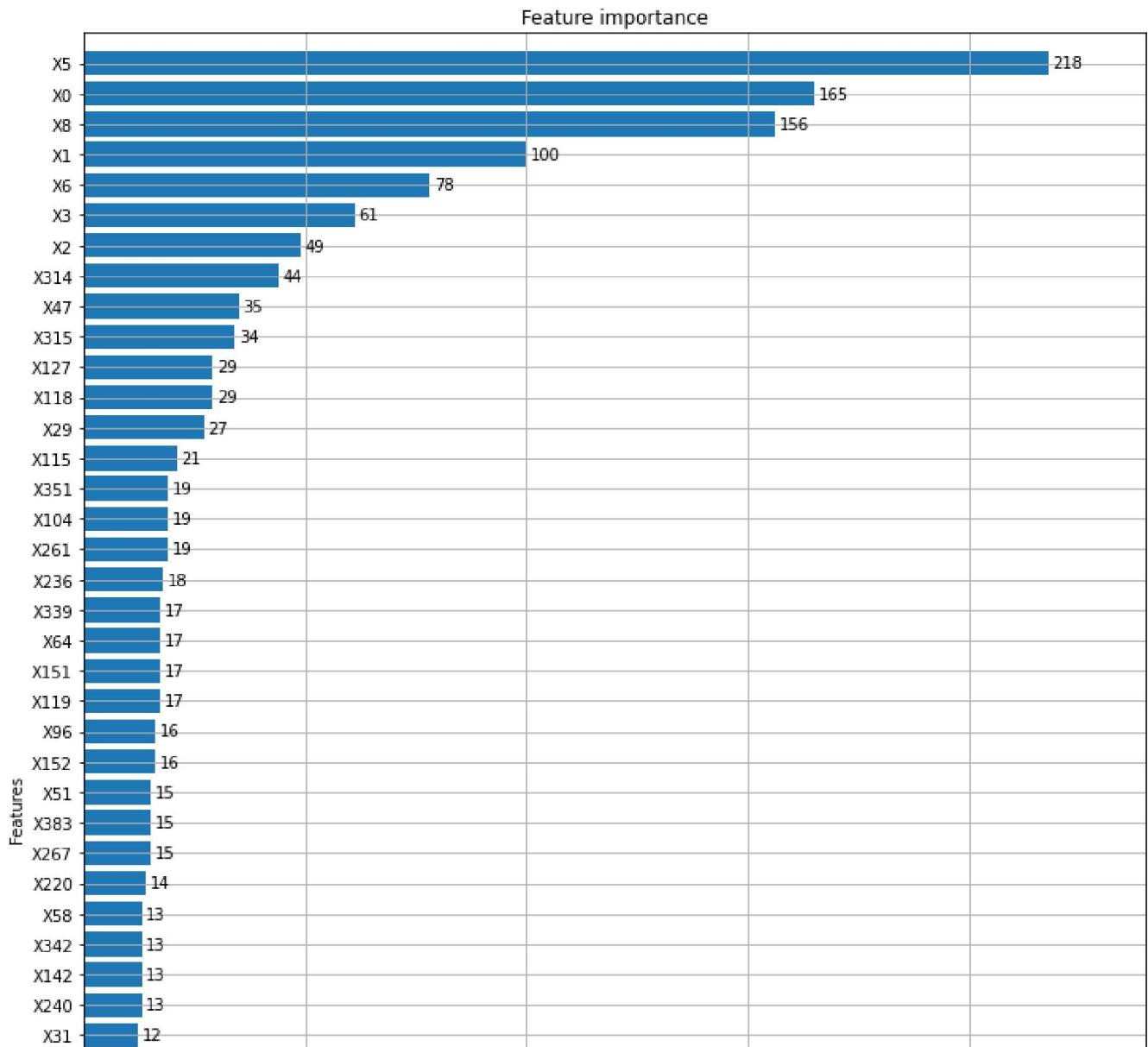
```
for f in ["X0", "X1", "X2", "X3", "X4", "X5", "X6", "X8"]:
    lbl = preprocessing.LabelEncoder()
    lbl.fit(list(train_df[f].values))
    train_df[f] = lbl.transform(list(train_df[f].values))

train_y = train_df['y'].values
train_X = train_df.drop(["ID", "y", "eval_set"], axis=1)

def xgb_r2_score(preds, dtrain):
    labels = dtrain.get_label()
    return 'r2', r2_score(labels, preds)

for f in ["X0", "X1", "X2", "X3", "X4", "X5", "X6", "X8"]:
    lbl = preprocessing.LabelEncoder()
```

[10:04:29] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now d



Categorical occupy the top spots followed by binary variables.

Let us also build a Random Forest model and check the important variables.

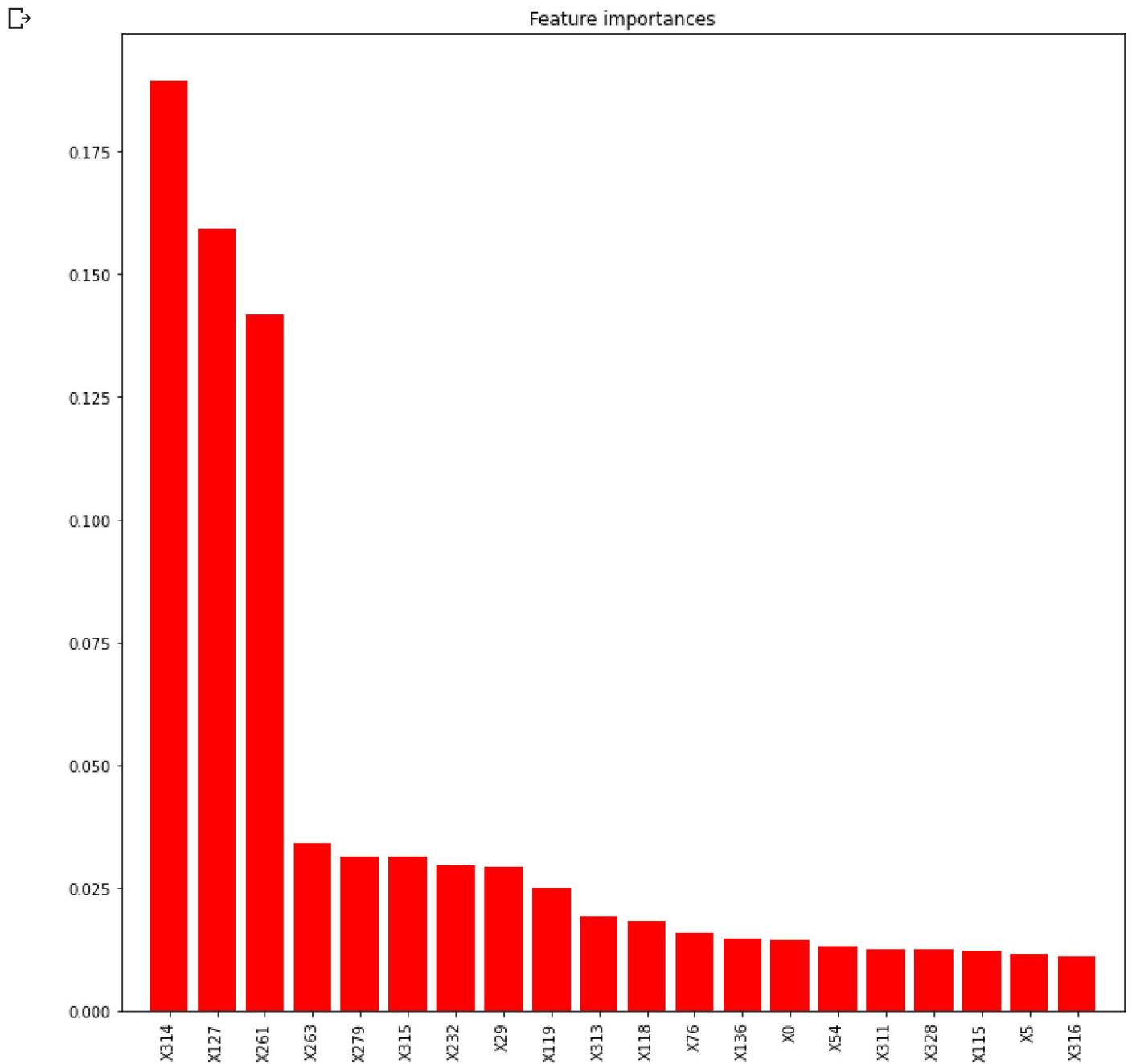


```
from sklearn.ensemble import RandomForestRegressor
model = RandomForestRegressor(n_estimators=200, max_depth=10, min_samples_leaf=4, max_features='sqrt')
model.fit(train_X, train_y)
feat_names = train_X.columns.values

## plot the importances ##
importances = model.feature_importances_
std = np.std([tree.feature_importances_ for tree in model.estimators_], axis=0)
indices = np.argsort(importances)[-1][:20]

plt.figure(figsize=(12,12))
plt.title("Feature importances")
plt.bar(range(len(indices)), importances[indices], color="r", align="center")
```

```
plt.xticks(range(len(indices)), feat_names[indices], rotation='vertical')
plt.xlim([-1, len(indices)])
plt.show()
```



Quite a few differences in the important variables between xgboost and random forest. Not sure why though.!