

```

In [7]: import multiprocessing
import os
import random
from datetime import datetime

import pandas as pd
import scipy
from sklearn import metrics
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler

# global parameters :)
train_path = 'train.csv'
test_path = 'eval.csv'
name = f"{datetime.now().strftime('%Y-%m-%d--%H:%M:%S')}"
xes = ['x1', 'x2', 'x3', 'x4', 'x5', 'x6', 'x7', 'x8', 'x9', 'x10']
grade_mapping = {'A': 7, 'B': 6, 'C': 5, 'D': 4, 'E': 3, 'F': 2, 'Fx': 1}
scale = False
scale_indices = xes[:4] + xes[6:]
y_mapping = {"Atsuto": 1, "Bob": 2, "Jörg": 3}
y_mapping_reversed = {v: k for k, v in y_mapping.items()}
true_false_mapping = {True: 1, False: 0, "True": 1, "False": 0}

def preprocess_train(df):
    df.dropna(inplace=True)
    df = df[~df[xes].isin(["?"]).any(axis=1)]
    df.loc[:, 'y'] = df.loc[:, 'y'].map(y_mapping)
    df.loc[:, ['x1', 'x2']] = df.loc[:, ['x1', 'x2']].astype(float)
    return df

def preprocess(df, scaler):
    df.loc[:, 'x6'] = df.loc[:, 'x6'].map(grade_mapping)
    df.loc[:, 'x5'] = df.loc[:, 'x5'].map(true_false_mapping)

    if scaler:
        scaler.transform(df[xes[:4] + xes[6:]])

    return df

if __name__ == '__main__':
    random.seed(72)
    pd.np.random.seed(72)
    scipy.random.seed(72)
    print(os.listdir())
    pd.set_option('display.max_rows', 500)
    pd.set_option('display.max_columns', 500)
    pd.set_option('display.width', 1000)

    train, test = pd.read_csv(train_path, comment="#"), pd.read_csv(test_path)

    # Preprocessing
    train = preprocess_train(train)
    scaler = StandardScaler().fit(train[scale_indices]) if scale else None
    train, test = preprocess(train, scaler), preprocess(test, scaler)

    X, y = train.loc[:, xes], train.loc[:, "y"]
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)

    # import matplotlib.pyplot as plt
    # import seaborn as sns
    # sns.set_style("whitegrid");

```

```
# sns.pairplot(train, hue="y", size=3);
# plt.show()
from autoviz.AutoViz_Class import AutoViz_Class
AV = AutoViz_Class()
dft = AV.AutoViz(
    "train_clean.csv",
    ",",
    "y",
    train,
    header=0,
    verbose=0,
    lowess=False,
    chart_format="svg",
    max_rows_analyzed=1500000,
    max_cols_analyzed=300,
)
```

```
['.git', '.ipynb_checkpoints', '.~lock.eval.csv#', '.~lock.train.csv#', 'demo.py', 'eval.csv', 'kth-ml-ch.txt', 'train.csv', 'train_clean.csv', 'Untitled.ipynb']
```

Shape of your Data Set: (996, 12)

Classifying variables in data set...

11 Predictors classified...

This does not include the Target column(s)

1 variables removed since they were ID or low-information variables

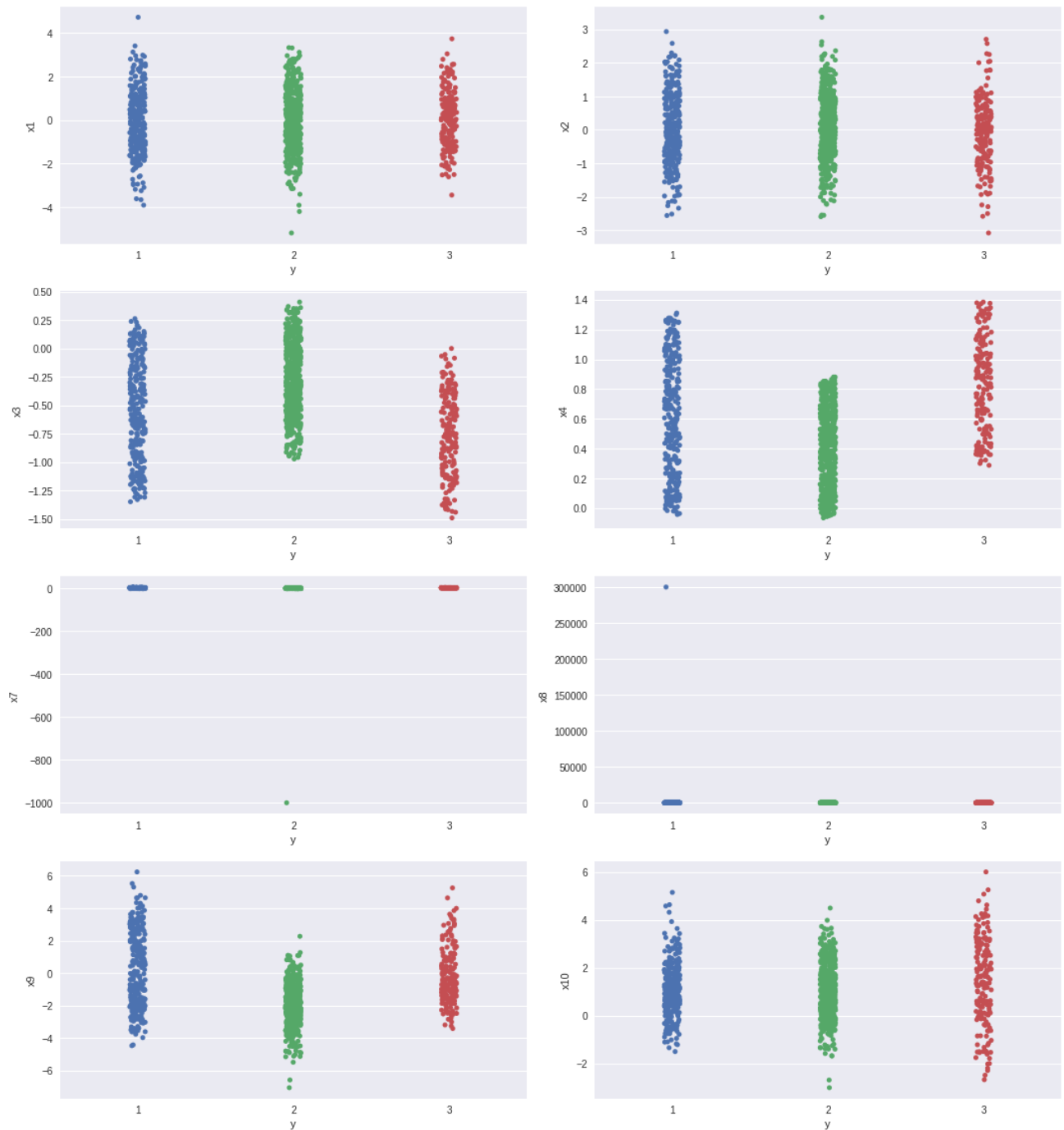
Total Number of Scatter Plots = 36

Could not draw Pivot Charts against Dependent Variable

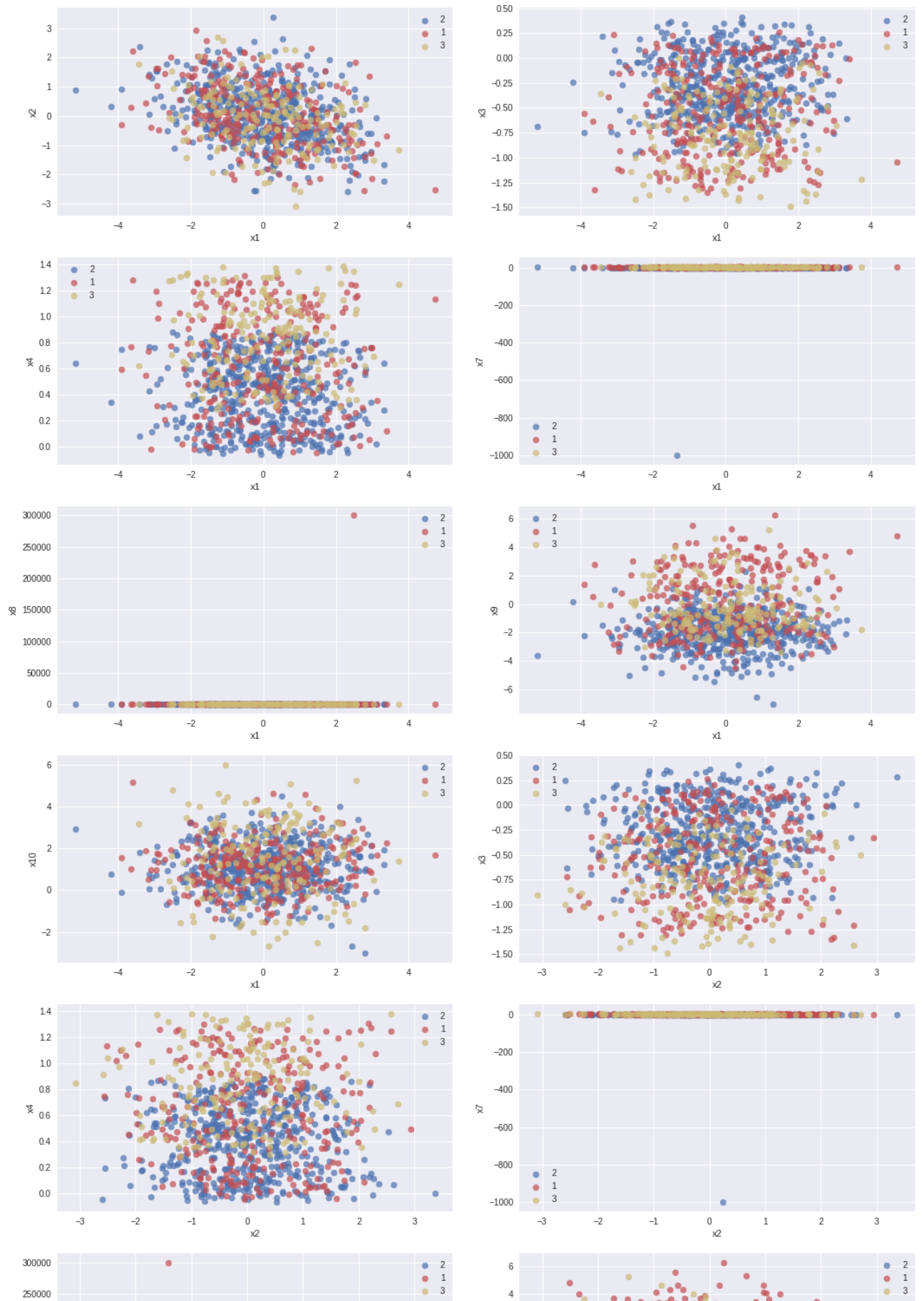
All plots done

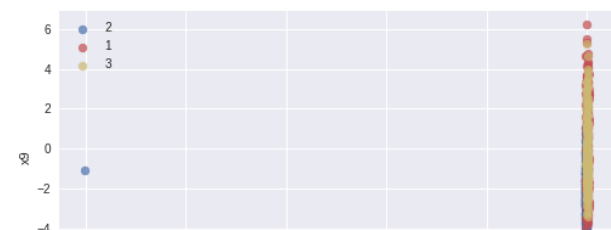
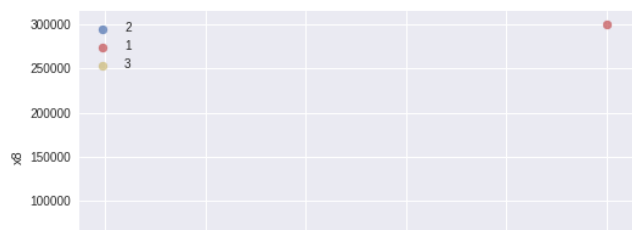
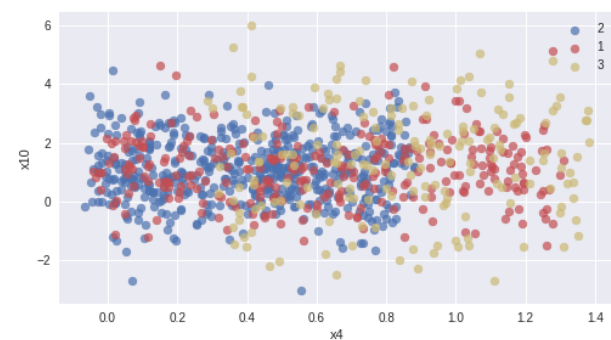
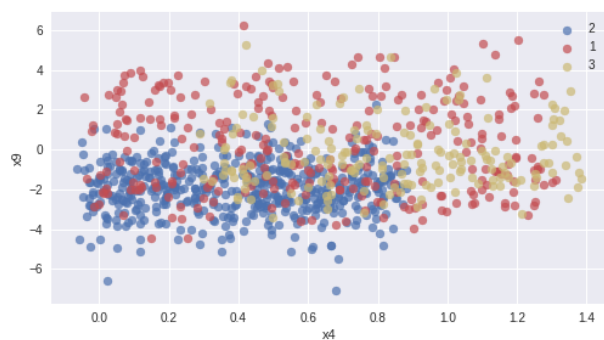
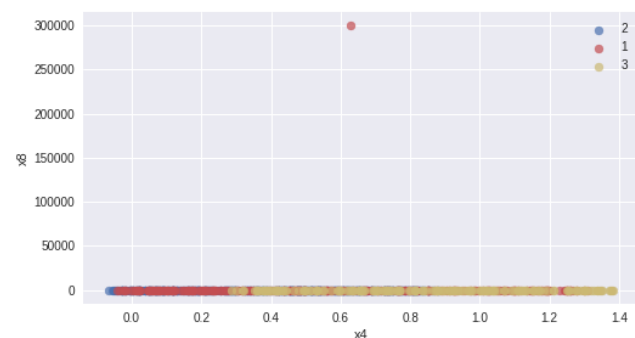
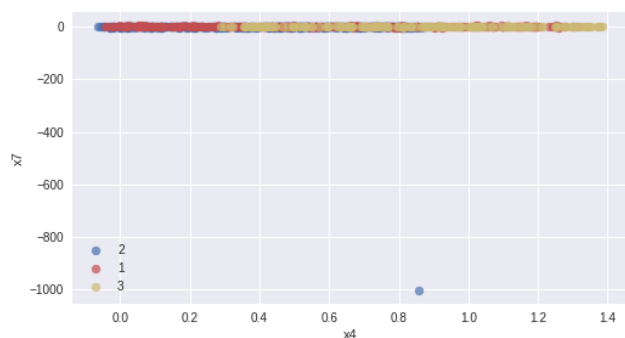
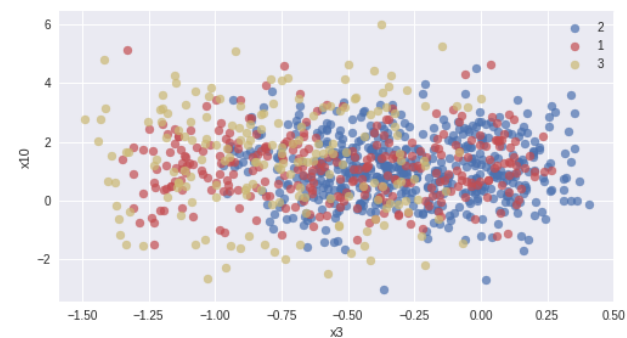
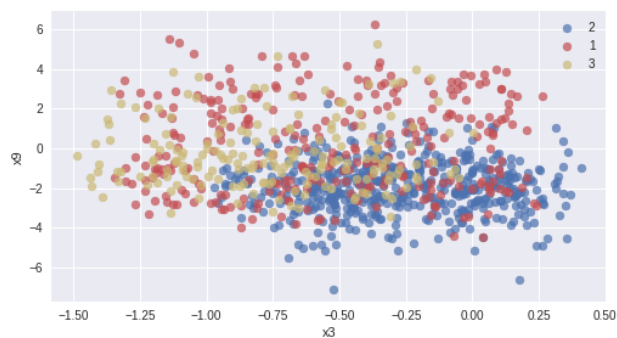
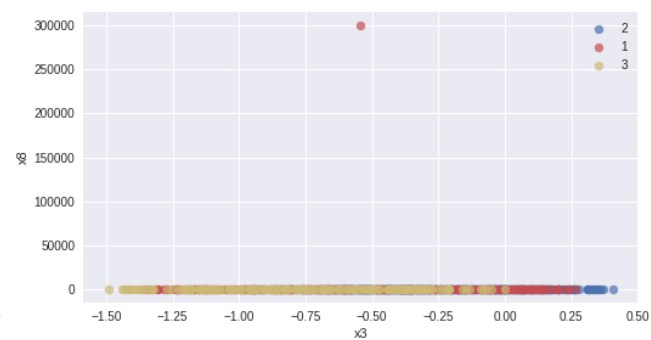
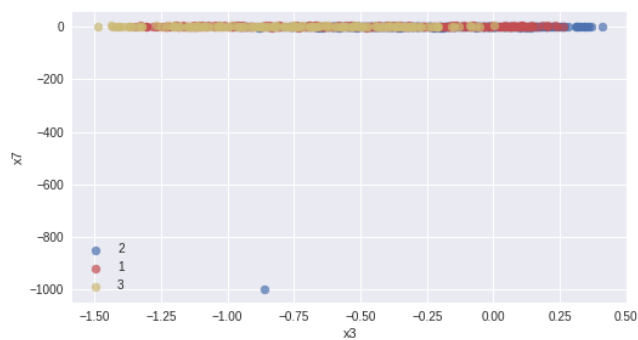
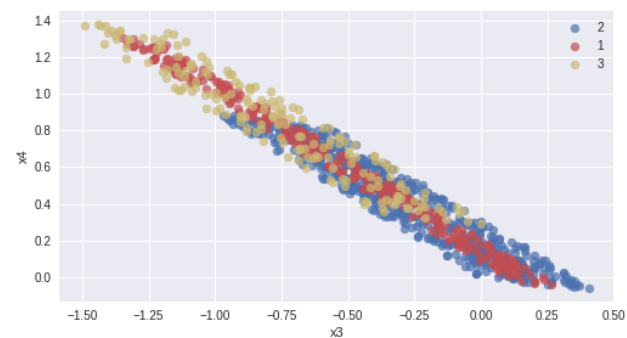
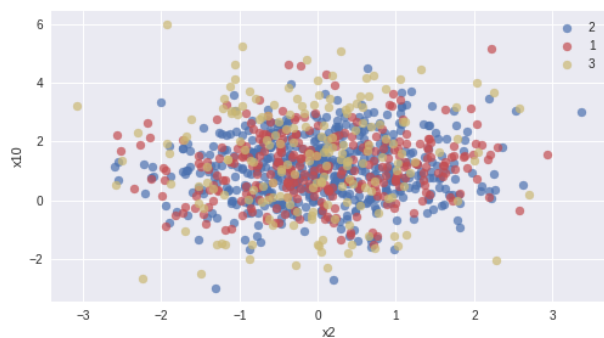
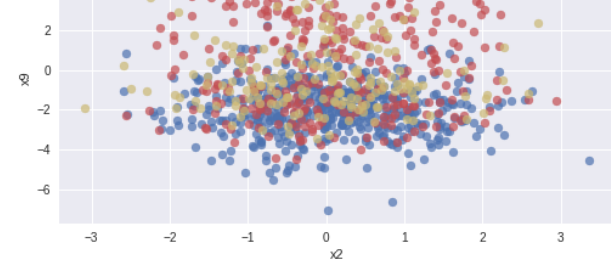
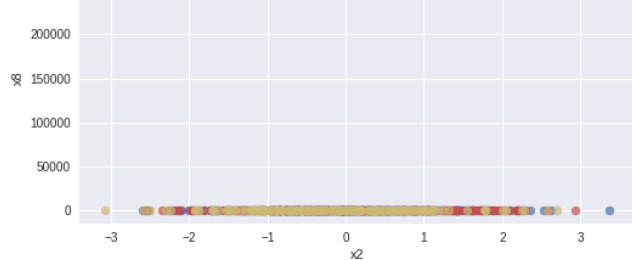
Time to run AutoViz (in seconds) = 5.221

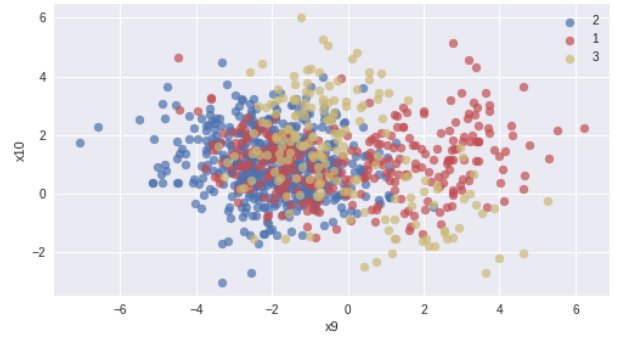
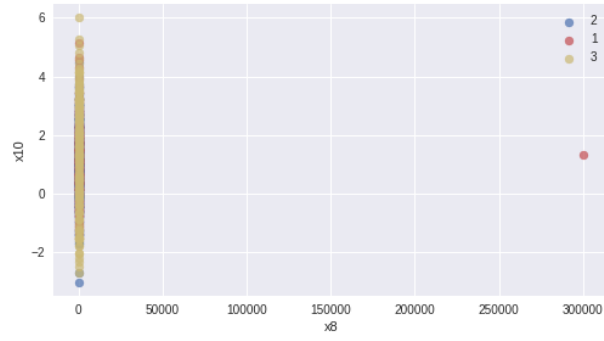
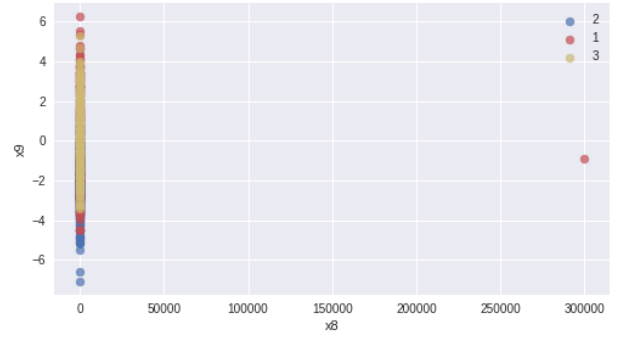
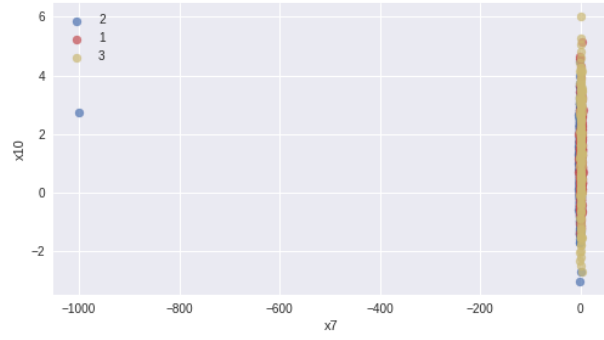
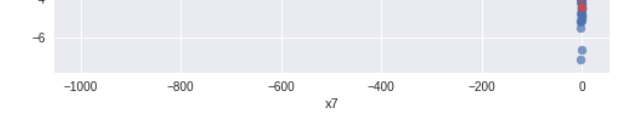
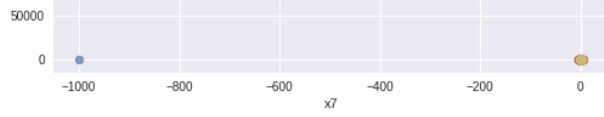
Scatter Plot of Continuous Variable vs Target (jitter=0.05)



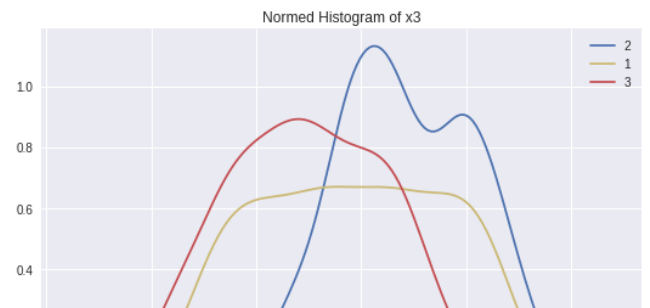
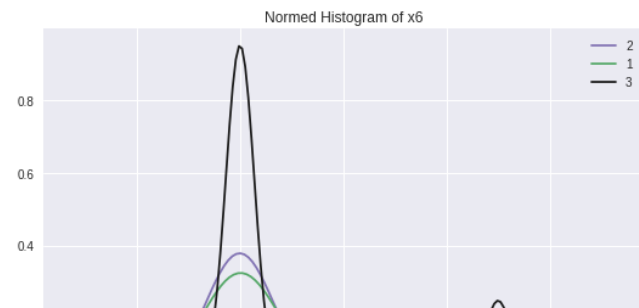
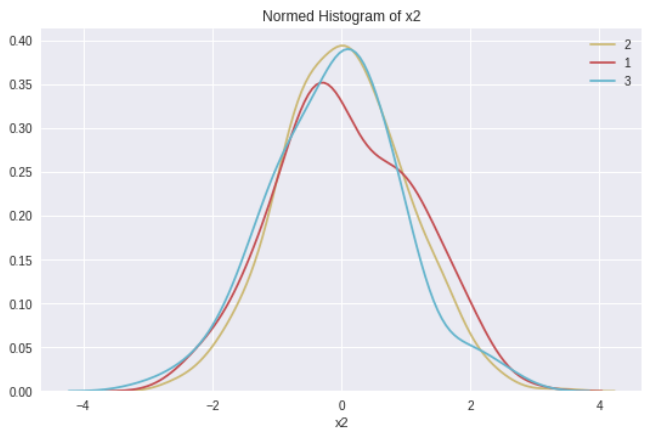
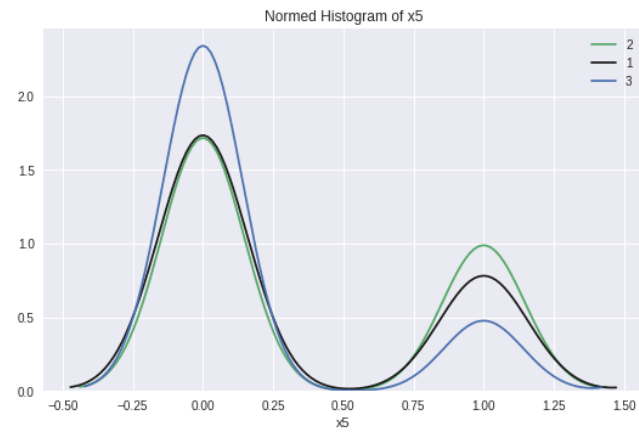
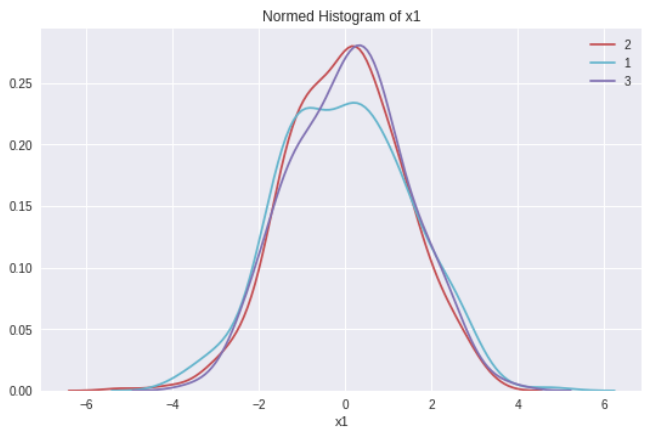
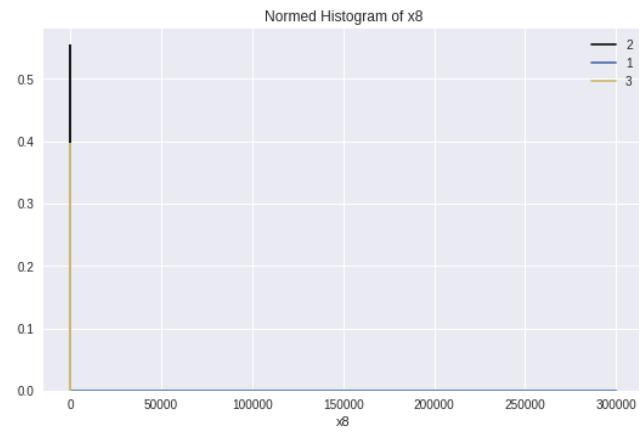
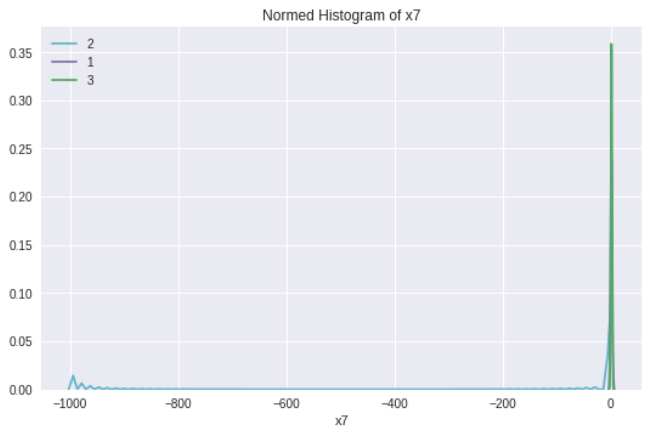
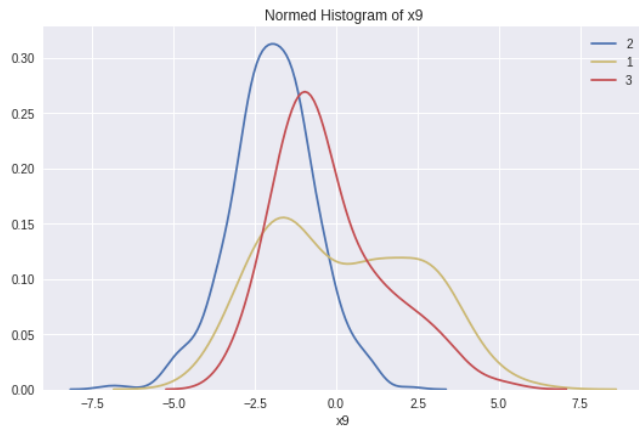
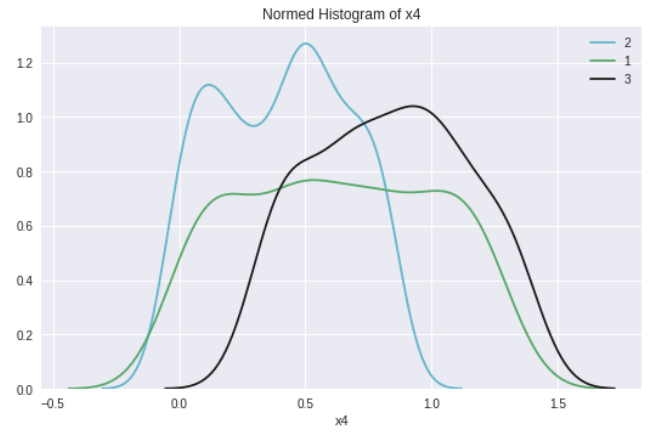
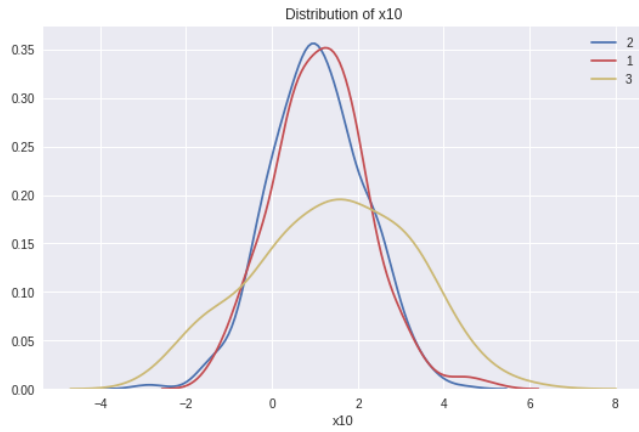
Scatter Plot of each Continuous Variable against Target Variable

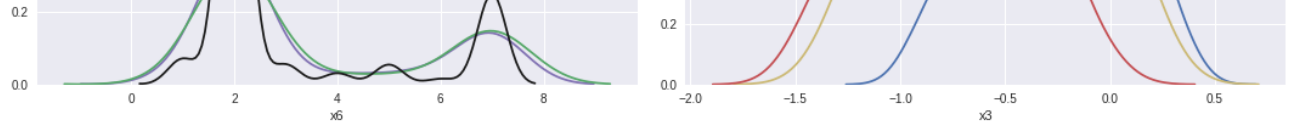




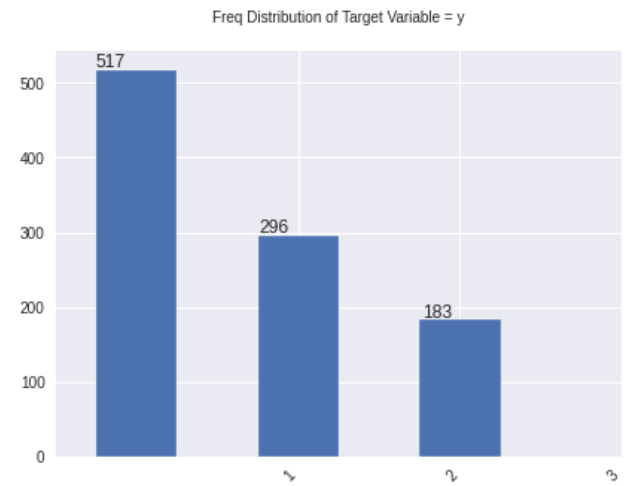
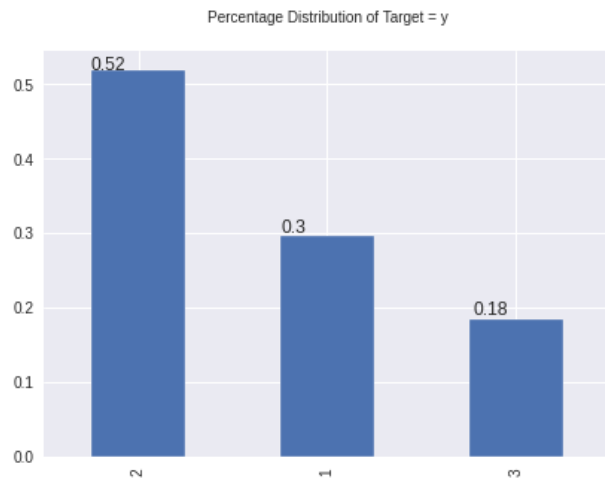


Histograms (KDE plots) of all Continuous Variables

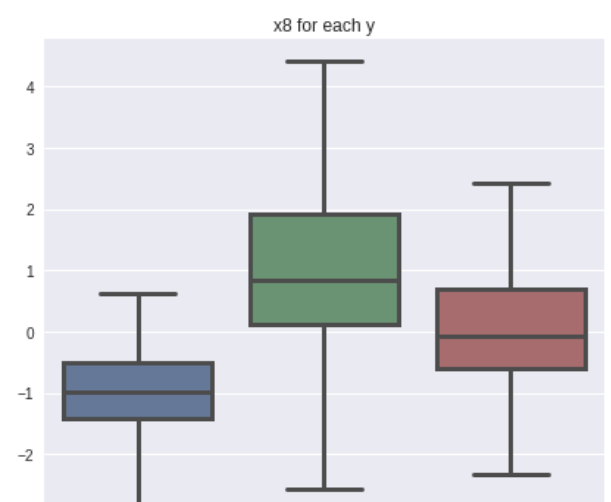
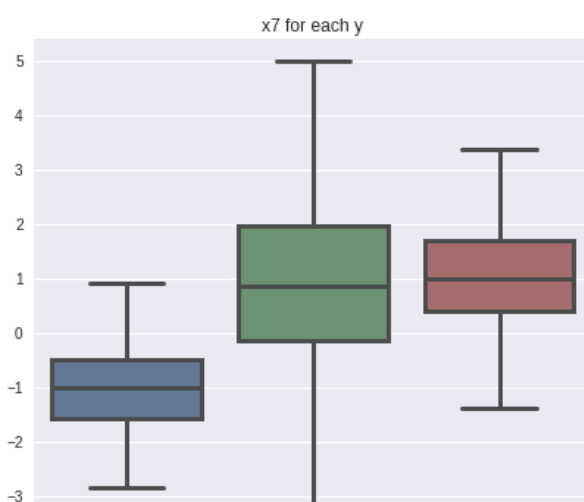
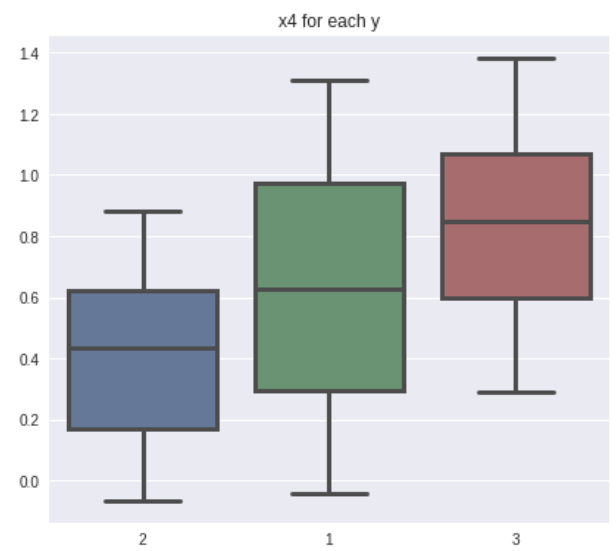
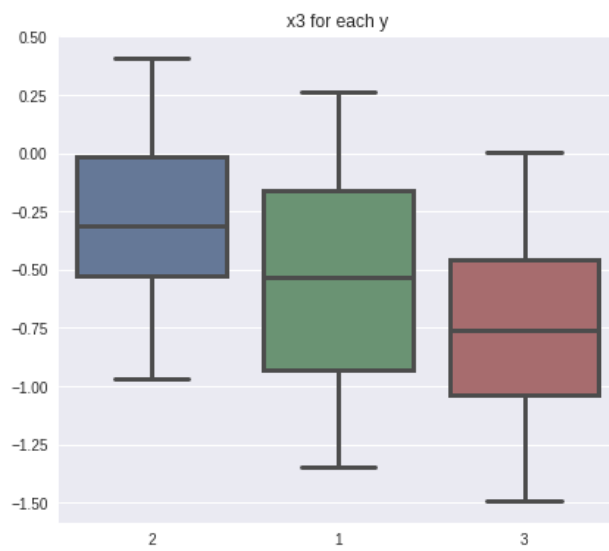
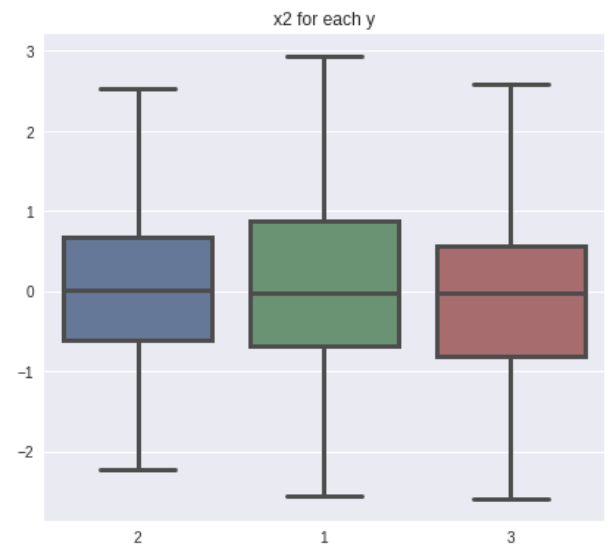
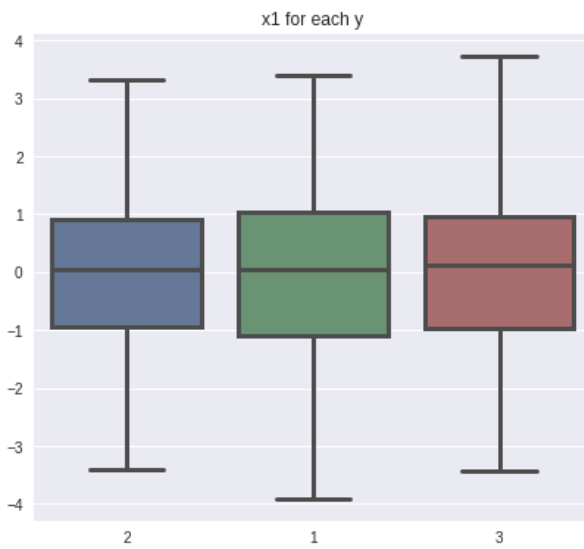


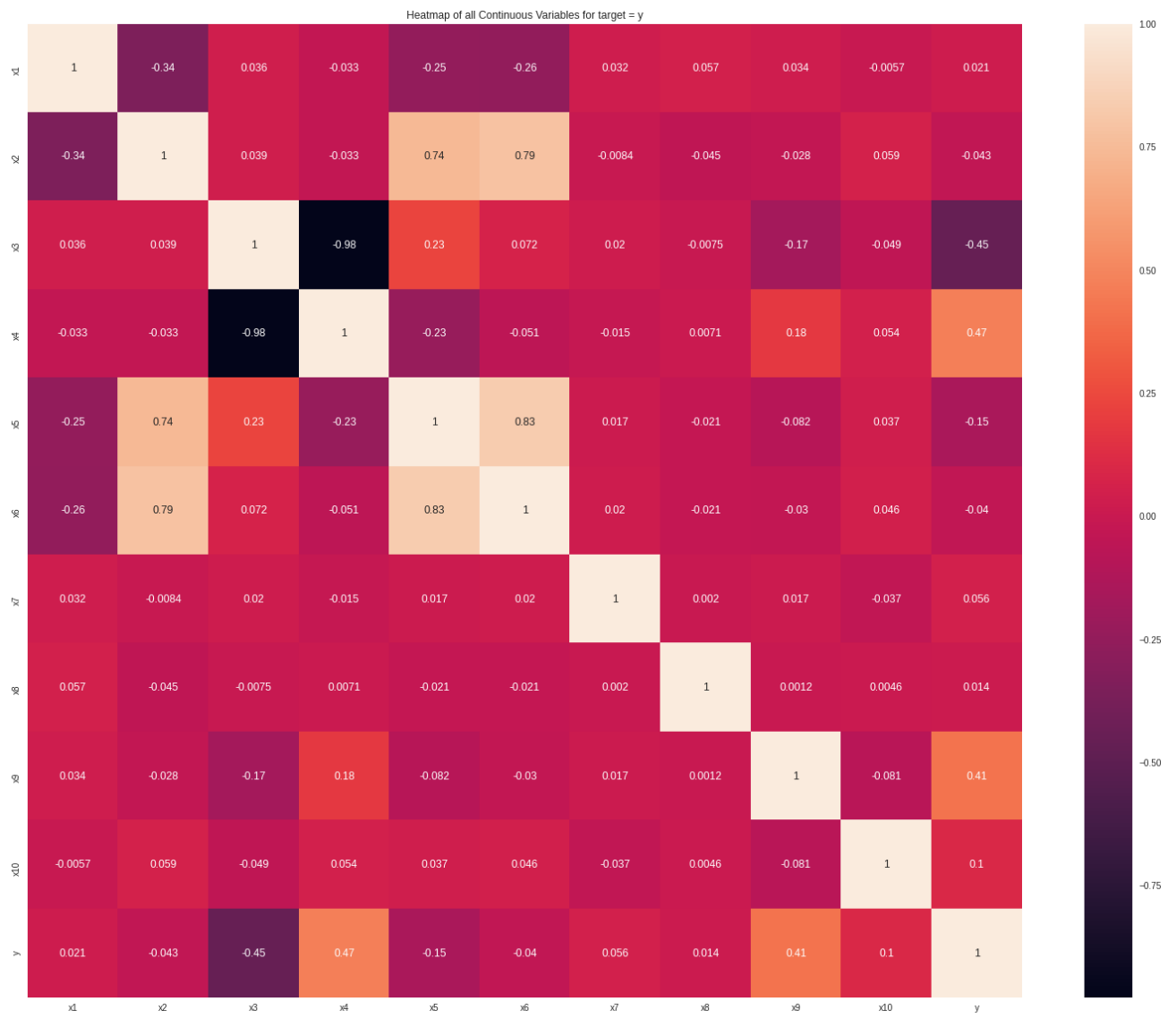
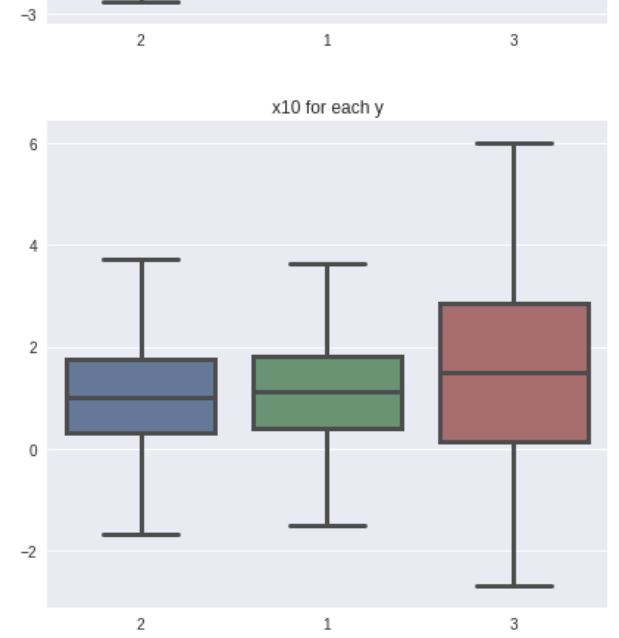
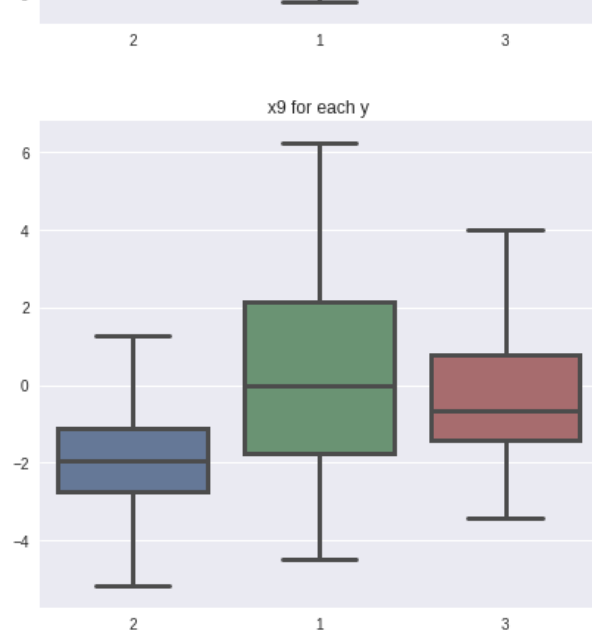


y : Distribution of Target Variable

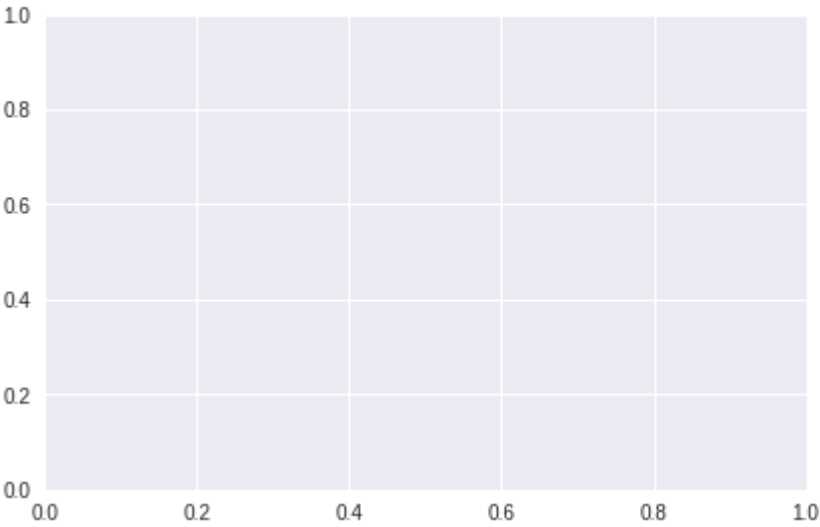


Box Plots without Outliers shown

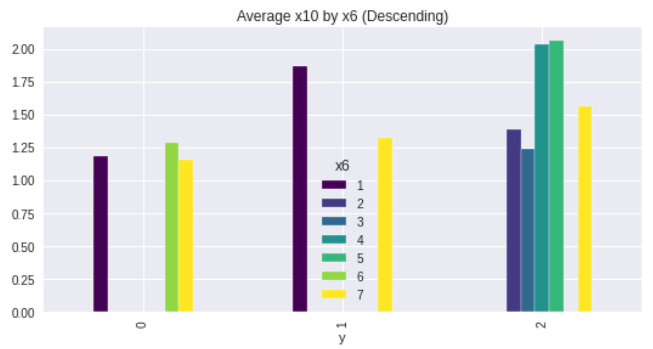
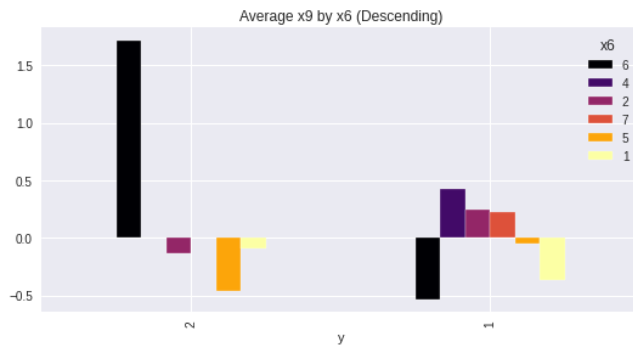
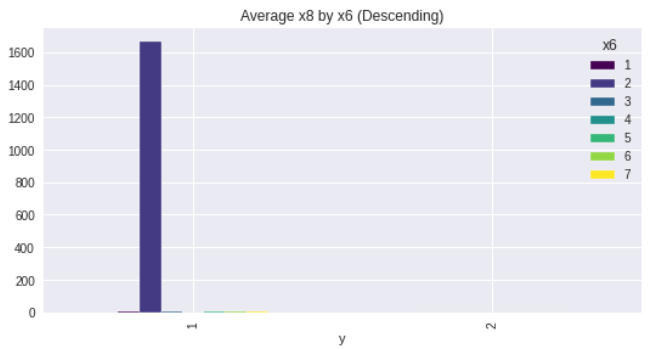
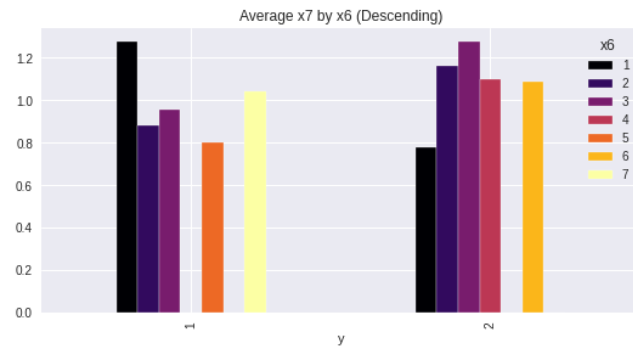
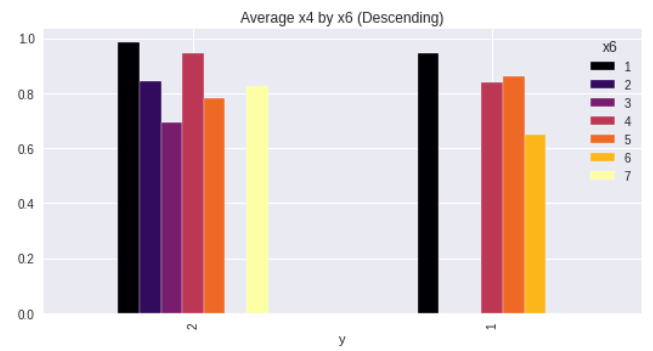
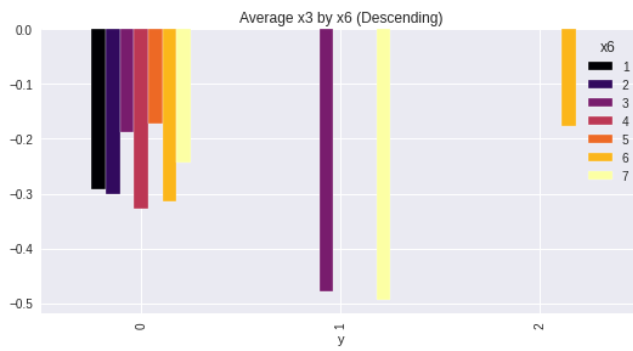
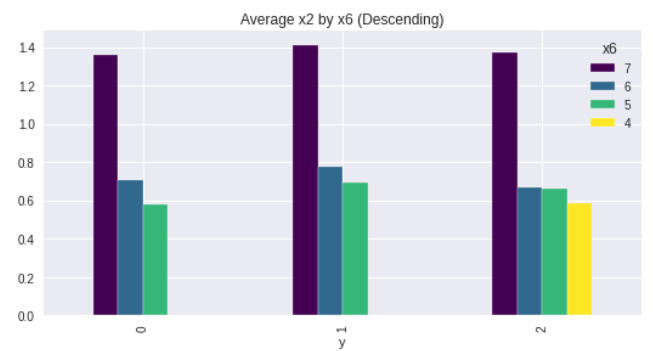
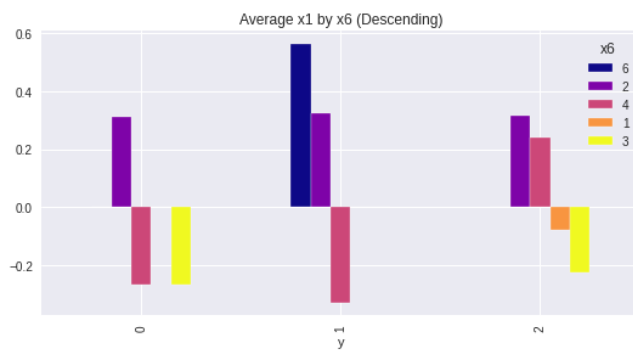




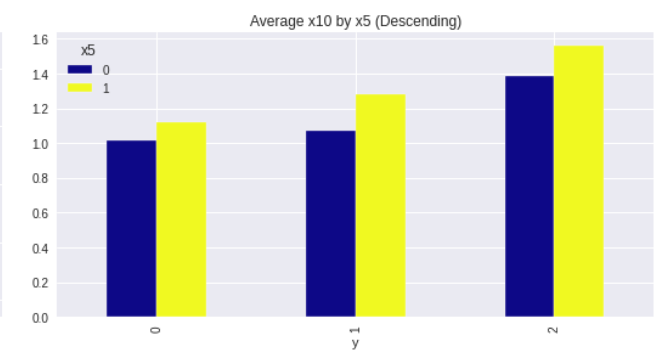
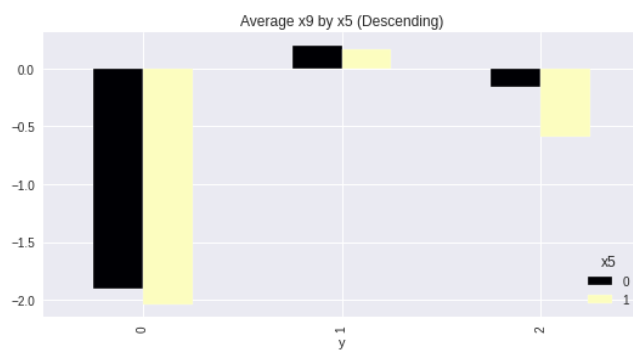
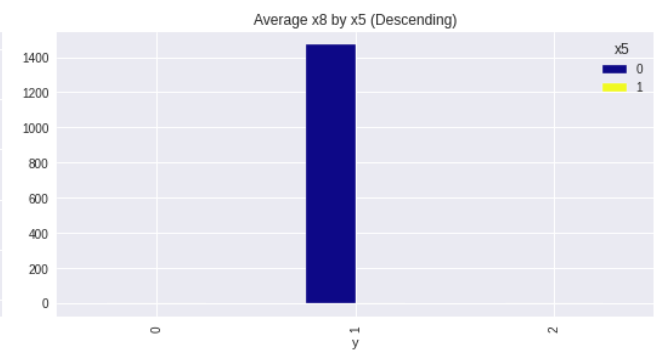
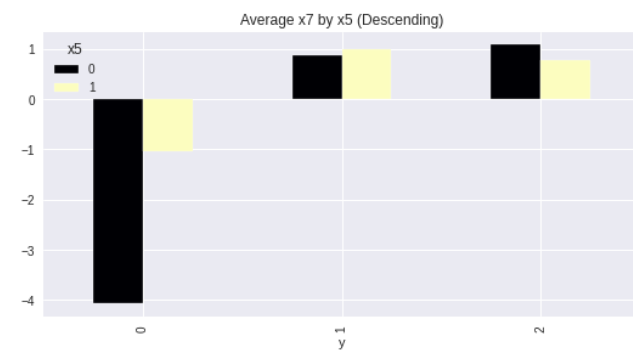
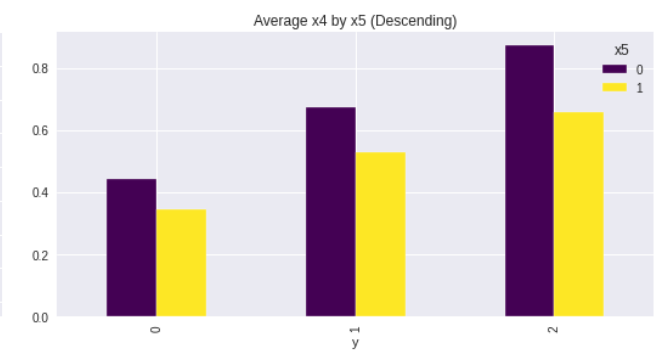
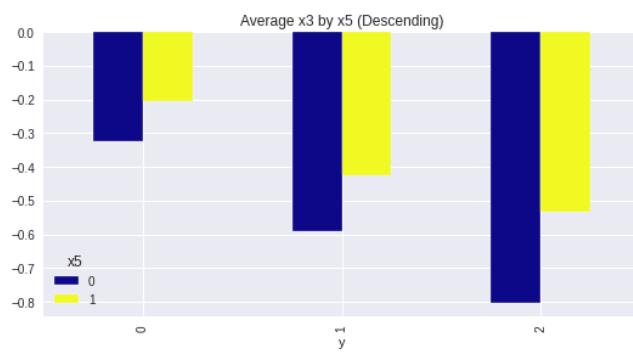
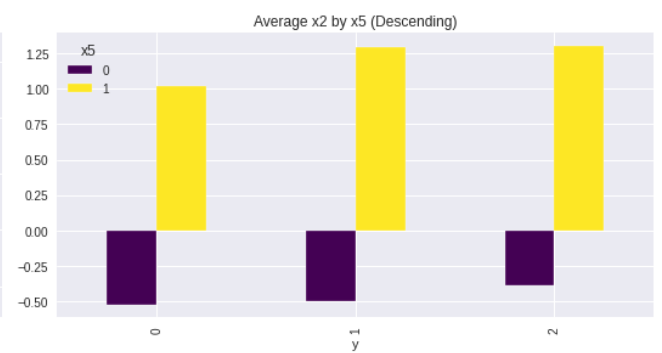
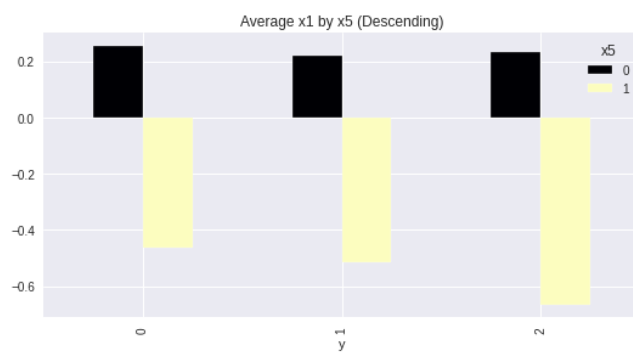
Plots of each Continuous Var by y



Bar Plots of Continuous Variables by x6



Bar Plots of Continuous Variables by x5



Bar Plots of Continuous Variables by x5

