

```
In [40]:  
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
import pandas_profiling  
from sklearn.preprocessing import LabelEncoder  
from sklearn.decomposition import PCA  
import xgboost as xgb  
import seaborn as sns  
from sklearn.metrics import r2_score  
from sklearn.model_selection import train_test_split  
import warnings  
warnings.filterwarnings('ignore')
```

## EDA

```
In [41]:  
df_train = pd.read_csv('mercedes_train.csv')  
  
print('Size of training set: {} rows and {} columns'.format(df_train.shape))  
# Print few rows and see how the data looks like  
df_train.head()
```

Size of training set: 4209 rows and 378 columns

```
Out[41]:  
   ID   y  X0  X1  X2  X3  X4  X5  X6  X8 ... X375  X376  X377  X378  X379  X380  X382  X383  X384  X385  
0  0  130.81  k  v  at  a  d  u  j  o ...  0  0  1  0  0  0  0  0  0  0  0  
1  6  88.53  k  t  av  e  d  y  l  o ...  1  0  0  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  0  0  0  
3  9  80.62  az  t  n  f  d  x  l  e ...  0  0  0  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  0  0  0
```

5 rows × 378 columns

```
In [42]:  
df_train.info()  
  
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 4209 entries, 0 to 4208  
Columns: 378 entries, ID to X385  
dtypes: float64(1), int64(369), object(8)  
memory usage: 12.1+ MB
```

```
In [43]:  
df_train.shape
```

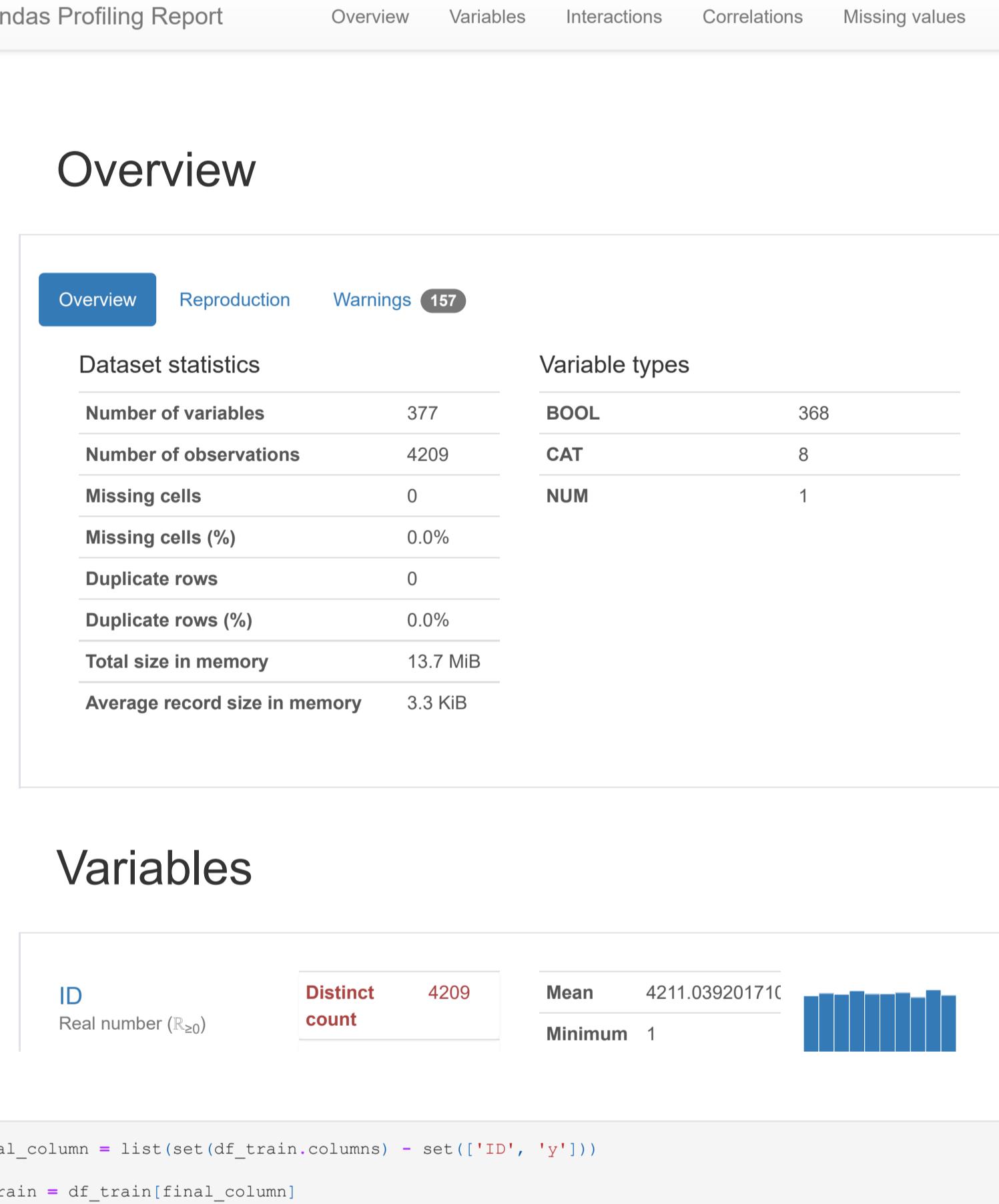
```
Out[43]: (4209, 378)
```

```
In [44]:  
pfr = pandas_profiling.ProfileReport(df_train)  
pfr.to_file("Descriptive_Analysis_train.html")
```

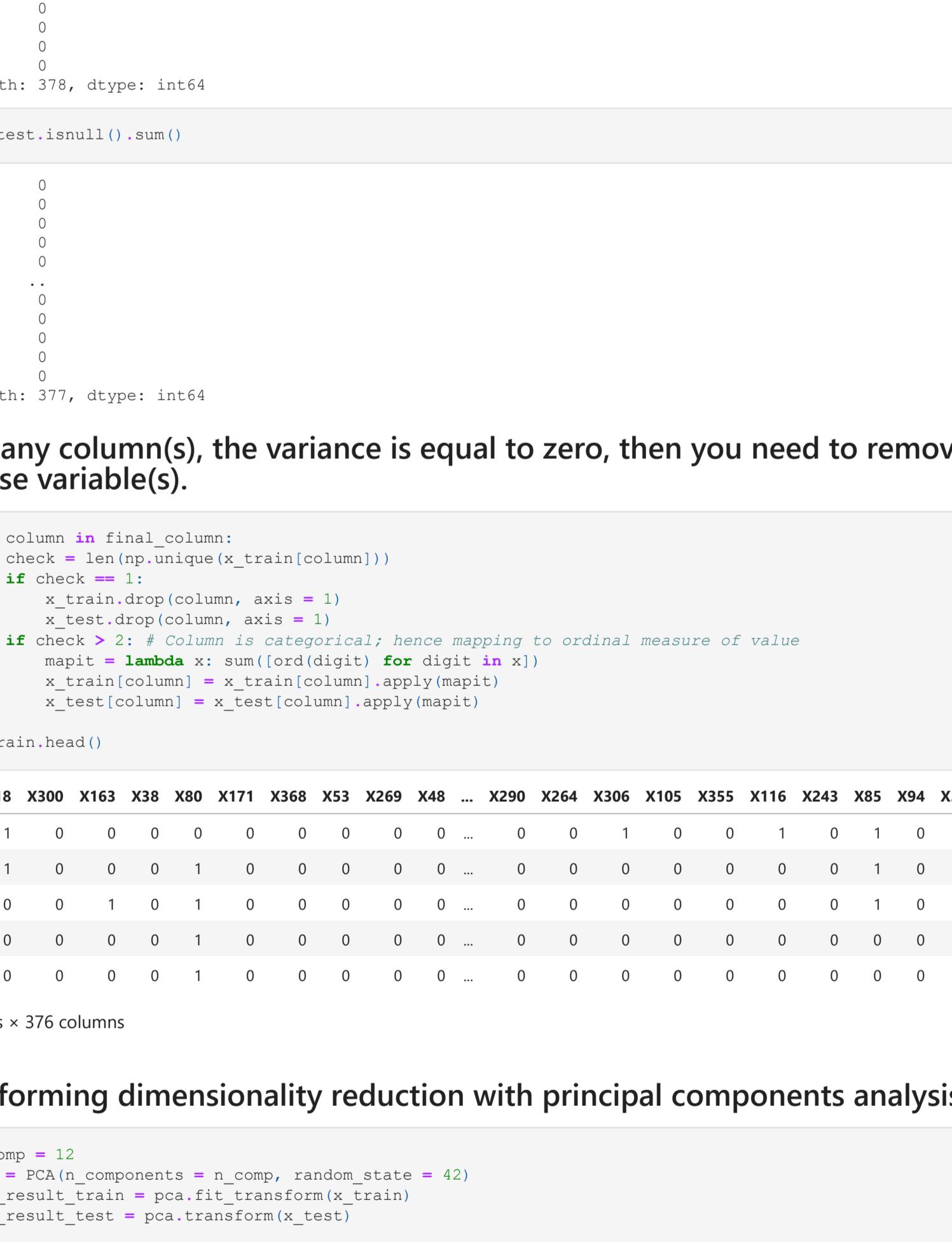
```
In [45]:  
pfr
```

Pandas Profiling Report      Overview      Variables      Interactions      Correlations      Missing values      Sample

## Overview



## Variables



```
Out[46]:  
y_train = df_train['y'].values  
y_train
```

```
Out[46]: array([130.81,  88.53,  76.26, ..., 109.22,  87.48, 110.85])
```

```
In [47]:  
var="X0"  
colu_order=np.sort(df_train[var].unique()).tolist()  
plt.figure(figsize=(12,6))  
sns.stripplot(x=var,y="y",data=df_train,order=colu_order)  
plt.xlabel(var,fontsize=12)  
plt.ylabel("y",fontsize=12)  
plt.title("Distribution of y variable with "+var, fontsize=15)  
plt.show()
```

Distribution of y variable with X0

250  
225  
200  
175  
150  
125  
100  
75

a a b a c a d a i a k a l a m a o a p a q a s a t a u a w a y a z b a b c d e f g h i j k l m n o q r s t u v w x y z

X0

250  
225  
200  
175  
150  
125  
100  
75

a a b a c a d a i a k a l a m a o a p a q a s a t a u a w a y a z b a b c d e f g h i j k l m n o q r s t u v w x y z

```
In [48]:  
plt.figure(figsize=(15,5))  
plt.subplot(121)  
sns.distplot(df_train.y.values, bins=20)  
plt.title('Target Value Distribution \n', fontsize=15)  
plt.xlabel('Target Value in Seconds'); plt.ylabel('Occurrences');
```

Target Value Distribution

250  
225  
200  
175  
150  
125  
100  
75

0.000 0.005 0.010 0.015 0.020 0.025 0.030 0.035 0.040

```
plt.subplot(122)  
sns.boxplot(df_train.y.values)  
plt.title('Target Value Distribution \n', fontsize=15)  
plt.xlabel('Target Value in Seconds');
```

Target Value Distribution

250  
225  
200  
175  
150  
125  
100  
75

0.000 0.005 0.010 0.015 0.020 0.025 0.030 0.035 0.040

```
In [49]:  
# Looking at the test dataset for similar features  
df_test = pd.read_csv('mercedes_test.csv')  
df_test.head()
```

5 rows × 377 columns

```
Out[49]:  
df_test.info()
```

<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 4209 entries, 0 to 4208  
Columns: 377 entries, ID to X385  
dtypes: int64(369), object(8)  
memory usage: 12.1+ MB

```
In [50]:  
df_test.shape
```

```
Out[50]: (4209, 377)
```

```
In [52]:  
pfr1 = pandas_profiling.ProfileReport(df_test)  
pfr1.to_file("Descriptive_Analysis_test.html")
```

```
In [53]:  
pfr1
```

Pandas Profiling Report      Overview      Variables      Interactions      Correlations      Missing values      Sample

## Overview



## Variables



```
Out[53]:  
final_column = list(set(df_train.columns) - set(['ID', 'y']))
```

```
In [54]:  
x_train = df_train[final_column]  
y_train = df_train['y'].values
```

```
Out[54]: array([130.81,  88.53,  76.26, ..., 109.22,  87.48, 110.85])
```

```
In [47]:  
var="X0"  
colu_order=np.sort(df_train[var].unique()).tolist()  
plt.figure(figsize=(12,6))  
sns.stripplot(x=var,y="y",data=df_train,order=colu_order)  
plt.xlabel(var,fontsize=12)  
plt.ylabel("y",fontsize=12)  
plt.title("Distribution of y variable with "+var, fontsize=15)  
plt.show()
```

Distribution of y variable with X0

250  
225  
200  
175  
150  
125  
100  
75

a a b a c a d a i a k a l a m a o a p a q a s a t a u a w a y a z b a b c d e f g h i j k l m n o q r s t u v w x y z

X0

250  
225  
200  
175  
150  
125  
100  
75

0.000 0.005 0.010 0.015 0.020 0.025 0.030 0.035 0.040

```
plt.subplot(122)  
sns.boxplot(df_train.y.values)
```

Target Value Distribution

250  
225  
200  
175  
150  
125  
100  
75

0.000 0.005 0.010 0.015 0.020 0.025 0.030 0.035 0.040

```
In [48]:  
plt.figure(figsize=(15,5))  
plt.subplot(121)  
sns.distplot(df_train.y.values, bins=20)  
plt.title('Target Value Distribution \n', fontsize=15)  
plt.xlabel('Target Value in Seconds'); plt.ylabel('Occurrences');
```

Target Value Distribution

250  
225  
200  
175  
150  
125  
100  
75

0.000 0.005 0.010 0.015 0.020 0.025 0.030 0.035 0.040

```
plt.subplot(122)  
sns.boxplot(df_train.y.values)
```

Target Value Distribution

250  
225  
200  
175  
150  
125  
100  
75

0.000 0.005 0.010 0.015 0.020 0.025 0.030 0.035 0.040

```
In [49]:  
# Looking at the test dataset for similar features  
df_test = pd.read_csv('mercedes_test.csv')  
df_test.head()
```

5 rows × 377 columns

```
Out[49]:  
df_test.info()
```

<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 4209 entries, 0 to 4208  
Columns: 377 entries, ID to X385  
dtypes: int64(369), object(8)  
memory usage: 12.1+ MB

```
In [50]:  
df_test.shape
```

```
Out[50]: (4209, 377)
```

```
In [52]:  
pfr1 = pandas_profiling.ProfileReport(df_test)  
pfr1.to_file("Descriptive_Analysis_test.html")
```

```
In [53]:  
pfr1
```

Pandas Profiling Report      Overview      Variables      Interactions      Correlations      Missing values      Sample

## Overview



## Variables



```
Out[53]:  
final_column = list(set(df_train.columns) - set(['ID', 'y']))
```

```
In [54]:  
x_train = df_train[final_column]
```

```
Out[54]: array([130.81,  88.53,  76.26, ..., 109.22,  87.48, 110.85])
```

```
In [47]:  
var="X0"  
colu_order=np.sort(df_train[var].unique()).tolist()  
plt.figure(figsize=(12,6))  
sns.stripplot(x=var,y="y",data=df_train,order=colu_order)  
plt.xlabel(var,fontsize=12)  
plt.ylabel("y",fontsize=12)  
plt.title("Distribution of y variable with "+var, fontsize=15)  
plt.show()
```

Distribution of y variable with X0

250  
225  
200  
175  
150  
125  
100  
75

a a b a c a d a i a k a l a m a o a p a q a s a t a u a w a y a z b a b c d e f g h i j k l m n o q r s t u v w x y z

X0

250  
225  
200  
175  
150  
125  
100  
75

0.000 0.005 0.010 0.015 0.020 0.025 0.030 0.035 0.040

```
plt.subplot(122)  
sns.boxplot(df_train.y.values)
```

Target Value Distribution

250  
225  
200  
175  
150  
125  
100  
75

0.000 0.005 0.010 0