

```
In [3]: import matplotlib.pyplot as plt

In [4]: import pandas as pd

In [5]: import numpy as np

In [6]: import seaborn as sns

In [7]: data = pd.read_csv("D:/water_potability.csv")

In [8]: data.head()
```

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity	Potability
0	NaNa	129.422921	20791.10081	7.300212	368.516441	564.30654	10.376783	86.90070	2.963125	0
1	3.716080	129.422921	18630.057858	6.635246	NaN	592.885359	15.180013	66.320676	4.500566	0
2	8.099124	224.236259	19809.541732	9.275984	NaN	418.606213	16.868637	66.420093	3.055934	0
3	8.316766	214.373394	22018.417441	8.059332	358.886136	363.266516	18.436524	100.341674	4.628771	0
4	9.092223	181.101509	17978.986339	6.546600	310.135738	398.410813	11.586279	31.997993	4.075075	0

```
In [9]: data.isnull().sum()
```

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity	Potability
ph	0	491	0	0	0	0	0	0	0	0
Hardness	0	0	0	0	0	0	0	0	0	0
Solids	0	0	0	0	0	0	0	0	0	0
Chloramines	0	0	0	0	0	0	0	0	0	0
Sulfate	0	781	0	0	0	0	0	0	0	0
Conductivity	0	0	0	0	0	0	0	0	0	0
Organic_carbon	0	0	0	0	0	0	0	0	0	0
Trihalomethanes	0	0	0	0	0	0	0	0	0	0
Turbidity	0	0	0	0	0	0	0	0	0	0
Potability	0	0	0	0	0	0	0	0	0	0
dtype:	int64									

```
In [10]: data=data.dropna()
```

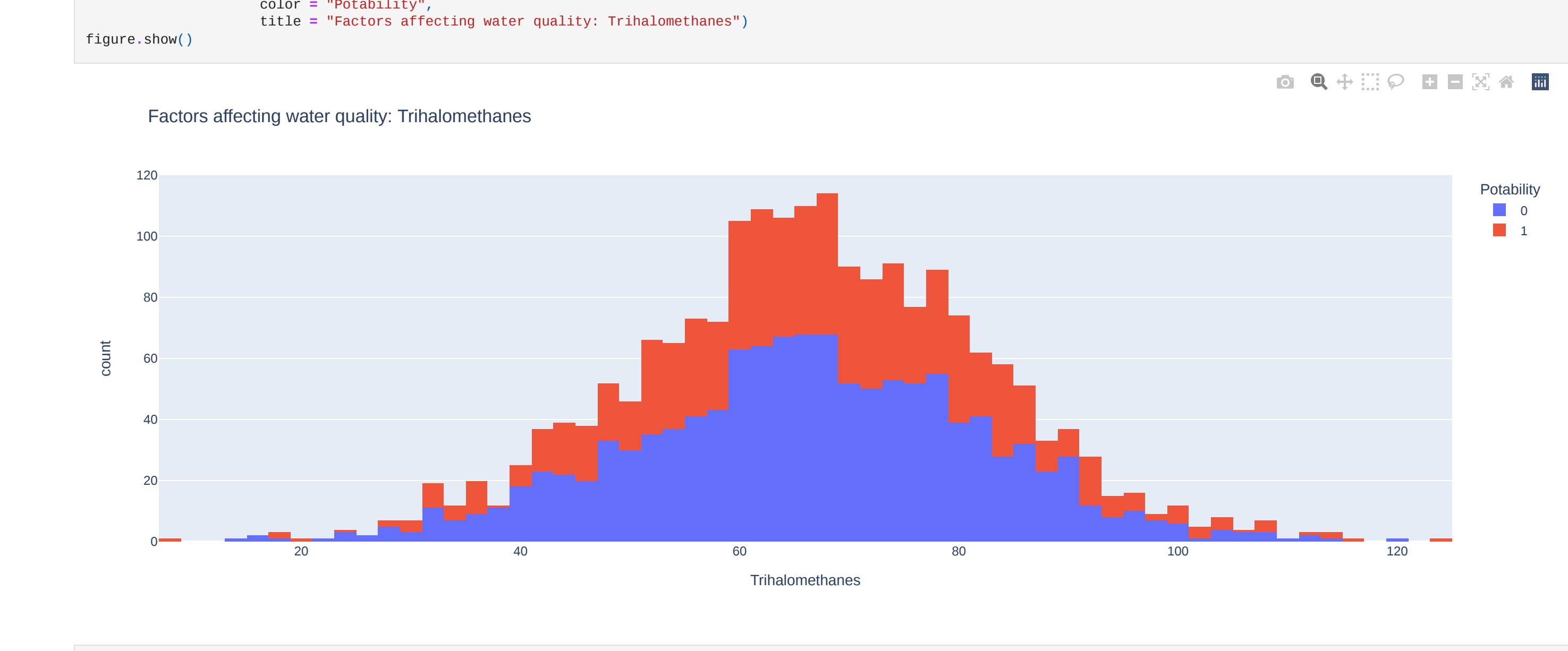
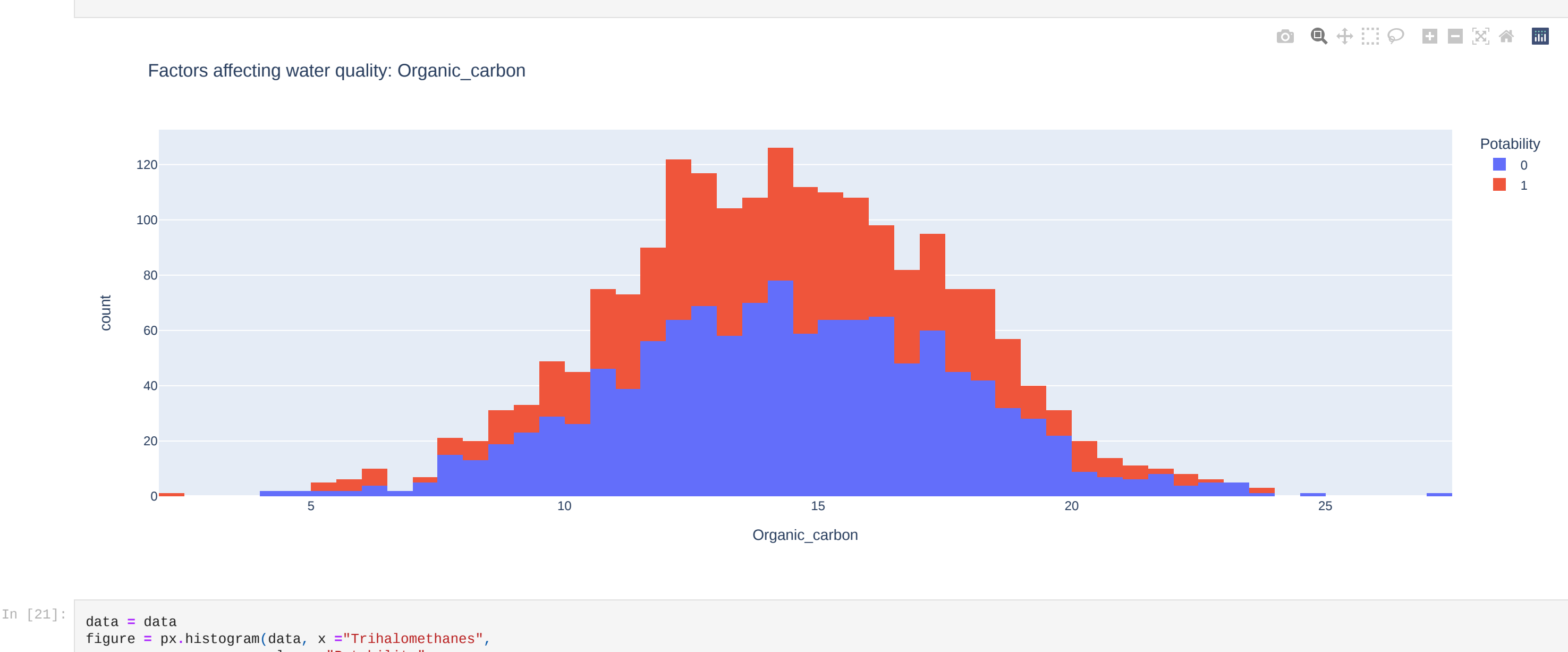
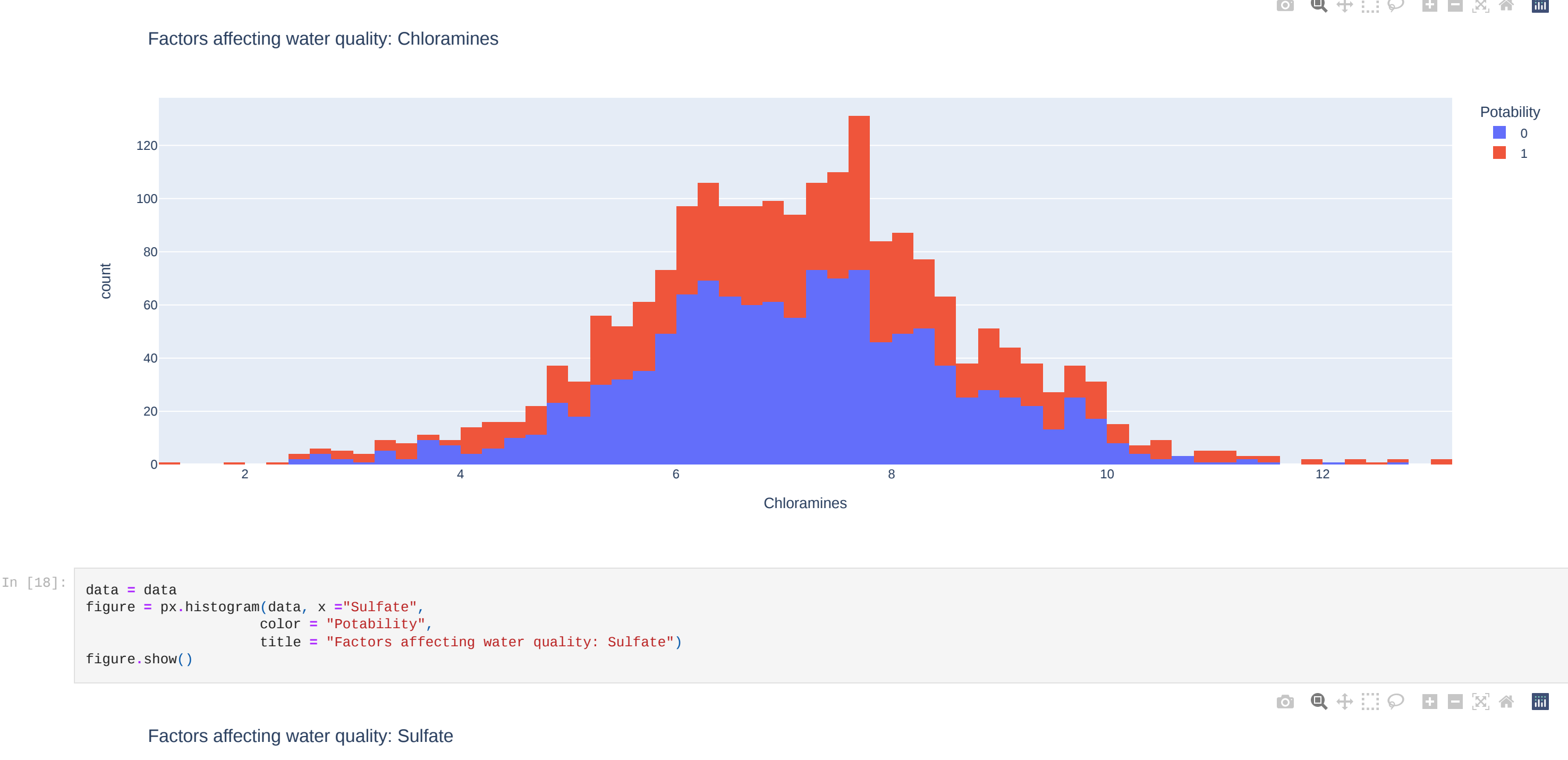
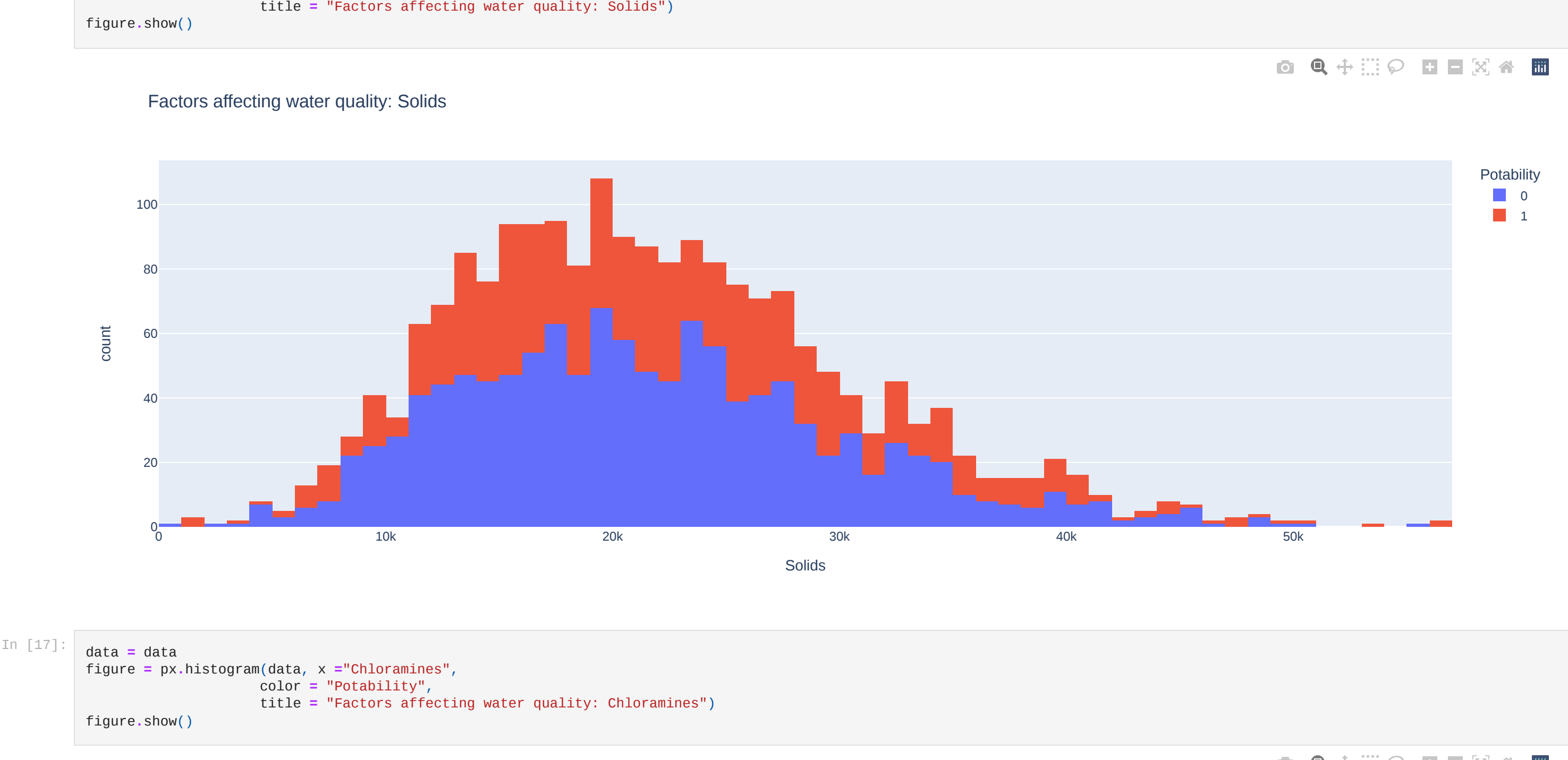
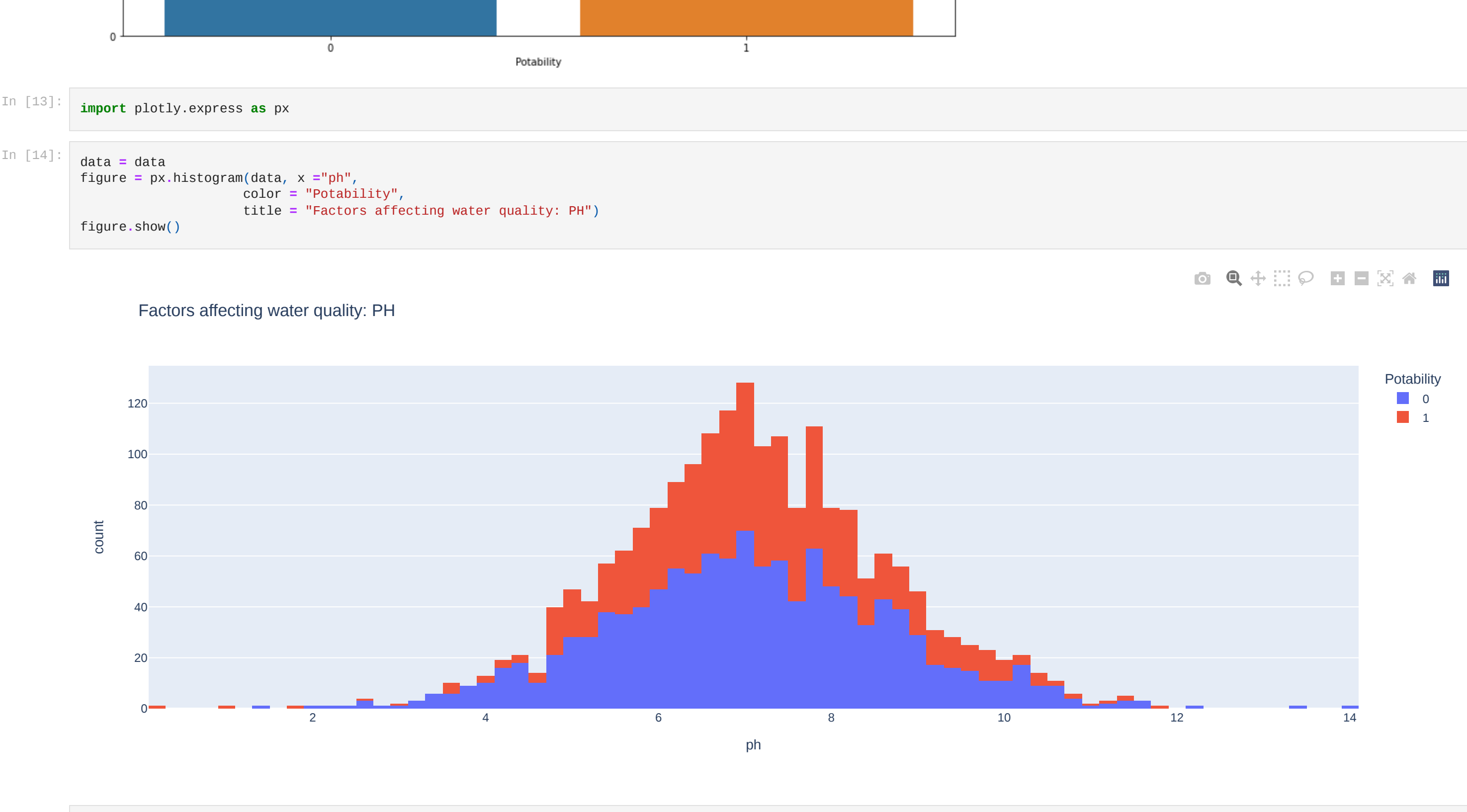
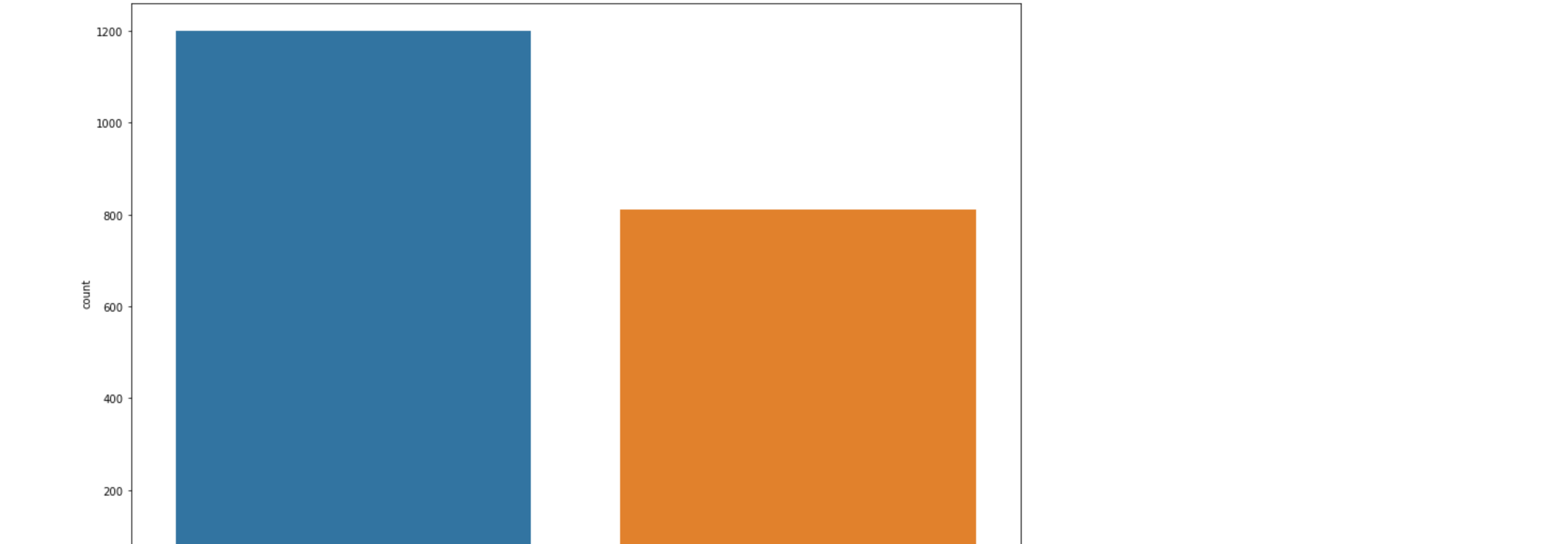
```
In [11]: data.isnull().sum()
```

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity	Potability
ph	0	0	0	0	0	0	0	0	0	0
Hardness	0	0	0	0	0	0	0	0	0	0
Solids	0	0	0	0	0	0	0	0	0	0
Chloramines	0	0	0	0	0	0	0	0	0	0
Sulfate	0	0	0	0	0	0	0	0	0	0
Conductivity	0	0	0	0	0	0	0	0	0	0
Organic_carbon	0	0	0	0	0	0	0	0	0	0
Trihalomethanes	0	0	0	0	0	0	0	0	0	0
Turbidity	0	0	0	0	0	0	0	0	0	0
Potability	0	0	0	0	0	0	0	0	0	0
dtype:	int64									

```
In [12]: plt.figure(figsize=(15,10))
sns.countplot(data.Potability)
plt.title("Distribution of unsafe and safe water")
plt.show()
```

C:\Users\lenovo\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be "data", and passing other arguments without an explicit keyword will result in an error or misinterpretation.

warnings.warn()



```
In [22]: correlation = data.corr()
```

```
In [23]: correlation.style.background_gradient('coolwarm')
```

```
Out[23]:
```

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity	Potability
ph	1.000000	0.108948	-0.087615	-0.024768	0.010524	0.014128	0.028378	0.018278	-0.035849	0.014530
Hardness	-0.108948	1.000000	-0.053209	-0.022085	-0.108521	0.011731	0.013224	-0.015400	-0.034831	-0.001505
Solids	-0.087615	-0.053209	1.000000	-0.051789	-0.162769	-0.005198	-0.005484	-0.015688	0.019409	0.001674
Chloramines	-0.024768	-0.022085	-0.051789	1.000000	0.008254	-0.028277	-0.023808	0.014990	0.013137	0.020784
Sulfate	0.010524	-0.108521	-0.162769	0.008254	1.000000	-0.016192	0.026776	-0.023347	-0.009834	-0.015303
Conductivity	0.014128	0.011731	-0.005198	-0.028277	-0.016192	1.000000	0.015647	-0.004888	-0.012495	-0.015486
Organic_carbon	0.028378	0.013224	-0.005484	-0.023808	0.026776	0.015647	1.000000	-0.005667	-0.015428	-0.015567
Trihalomethanes	0.018278	-0.015400	-0.015688	0.014990	-0.023347	0.004888	-0.005667	1.000000	0.020497	0.009244
Turbidity	-0.035849	-0.034831	0.019409	0.013137	-0.009834	0.012495	-0.015428	-0.020497	1.000000	0.022682
Potability	0.014530	-0.001505	0.001674	0.020784	-0.015303	-0.015486	-0.015567	0.009244	0.022682	1.000000

```
In [24]: y = data["Potability"]
X = data.drop(["Potability"], axis=1)
```

```
In [25]: from sklearn.preprocessing import StandardScaler
```

```
In [26]: scaler = StandardScaler()
scaler.fit(X)
```

```
Out[26]: StandardScaler()
```

```
In [27]: from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report, precision_score
from sklearn.linear_model import LogisticRegression
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.naive_bayes import GaussianNB
from sklearn.neural_network import MLPClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
```

```
In [28]: X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                         test_size=0.2,
                                                         random_state=42)
```

```
In [29]: models = [{"LogisticReg": LogisticRegression(max_iter=1000)}, {"SVC": SVC()},
{"DecisionTree": DecisionTreeClassifier()}, {"NaiveBayes": GaussianNB()},
{"RandomForest": RandomForestClassifier()}, {"GradientBoost": GradientBoostingClassifier()},
{"NeuralNetwork": MLPClassifier()}]
```

```
In [30]: results = []
names = []
accuracies = []
finalResults = []

for name,model in models:
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    score = precision_score(y_test, y_pred, average='macro')
    accuracyScore = accuracy_score(y_test, y_pred)
    accuracies.append(accuracyScore)
    results.append(score)
    names.append(name)
    finalResults.append((name,score, accuracyScore))
```

```
In [38]: finalResults_df = pd.DataFrame({'Model-Name': names, 'Accuracy-Score': accuracies, 'Precision-Score': results})
finalResults_df.sort_values(by='Accuracy-Score', ascending=False)
```

```
Out[38]:
```

	Model-Name	Accuracy-Score	Precision-Score
3	RandomForest	0.650124	0.647632
4	GradientBoost	0.645161	0.648925
2	DecisionTree	0.627792	0.616706
3	NaiveBayes	0.622829	0.636385
0	LogisticReg	0.573201	0.286600
1	SVC	0.573201	0.286600
6	NeuralNetwork	0.488983	0.514313

```
In [43]: rf_params = [{"max_depth": [12, 14, 15, 17, 19],
"max_features": [3, 5, 7, 10, 11],
"n_estimators": [100, 150, 180, 200, 300, 500],
"min_samples_split": [11, 13, 15, 16, 18]}]
```

```
rf = RandomForestClassifier()
rf_cv_model = GridSearchCV(rf, rf_params, cv=10, n_jobs=-1, verbose=2)
rf_cv_model.fit(X_train, y_train)
```

```
Out[43]: Fitting 10 folds for each of 750 candidates, totalling 7500 fits
GridSearchCV(cv=10, estimator=RandomForestClassifier(), n_jobs=-1,
param_grid={'max_depth': [12, 14, 15, 17, 19],
'max_features': [3, 5, 7, 10, 11],
'min_samples_split': [11, 13, 15, 16, 18],
'n_estimators': [100, 150, 180, 200, 300, 500]},
verbose=2)
```

```
In [42]: rf_cv_model.best_params_
```

```
Out[42]: {'max_depth': 12,
'max_features': 7,
'min_samples_split': 15,
'min_samples_split': 100}
```

```
In [39]: rf_tuned = RandomForestClassifier(max_depth = 12, max_features = 7, min_samples_split = 10, n_estimators = 100)
rf_tuned.fit(X_train, y_train)
```

```
Out[39]: RandomForestClassifier(max_depth=10, max_features=8, min_samples_split=10)
```

```
In [40]: y_pred = rf_tuned.predict(X_test)
accuracy_score(y_test, y_pred)
```

```
Out[40]: 0.674937965268546
```

```
In [ ] :
```