

Importing the libraries

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
```

Load Dataset

```
url = 'https://docs.google.com/spreadsheets/d/1p_WuY33JZo00wRFvtI7kEAITRhrwG0OM/export?format=csv'
df = pd.read_csv(url)
df.head()
```

1 to 5 of 5 entries

Filter

?

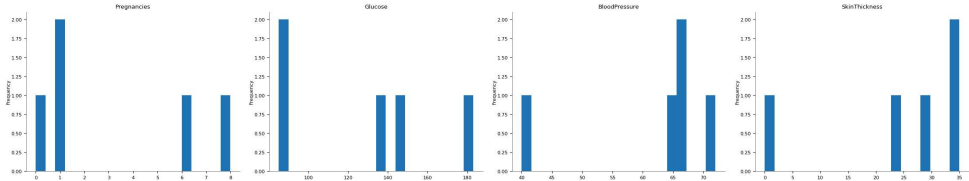
index	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1

Show 25 per page

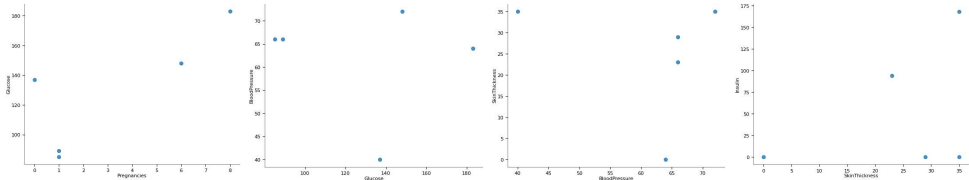
il

Like what you see? Visit the [data table notebook](#) to learn more about interactive tables.

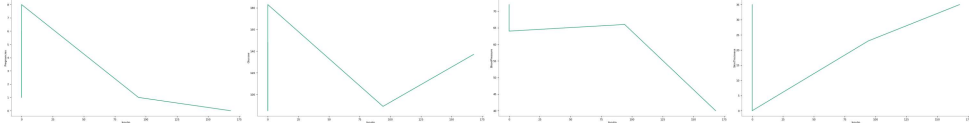
Distributions



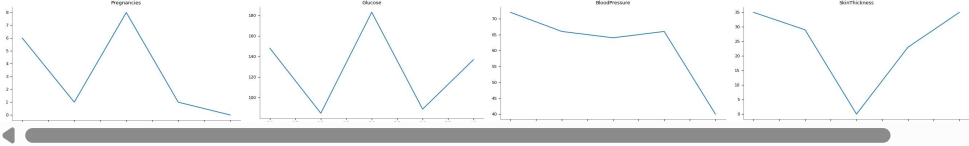
2-d distributions



Time series



Values



Next steps: [Generate code with df](#) [View recommended plots](#) [New interactive sheet](#)


Analyze the Dataset

```
print("Dataset info:")
df.info()
```

```
Dataset info:
<class 'pandas.core.frame.DataFrame'>
```

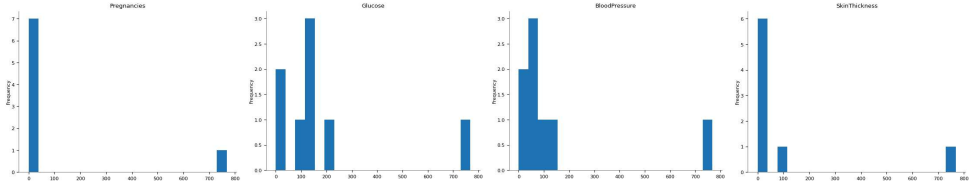
```
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
#   Column                                Non-Null Count  Dtype
---  ---                                ---
0   Pregnancies                          768 non-null    int64
1   Glucose                              768 non-null    int64
2   BloodPressure                       768 non-null    int64
3   SkinThickness                       768 non-null    int64
4   Insulin                             768 non-null    int64
5   BMI                                 768 non-null    float64
6   DiabetesPedigreeFunction            768 non-null    float64
7   Age                                 768 non-null    int64
8   Outcome                             768 non-null    int64
dtypes: float64(2), int64(7)
memory usage: 54.1 KB
```

```
print("\nsummary Statistics:")
df.describe()
```

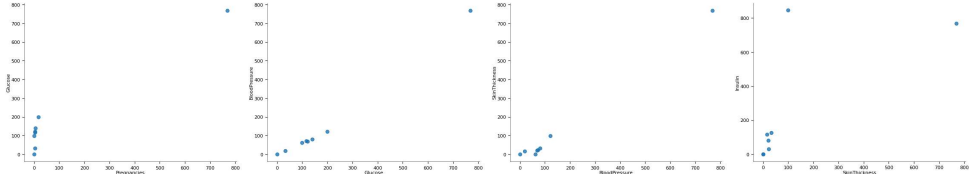
summary Statistics:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	0.471876	33.240885	0.348958
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	0.331329	11.760232	0.476951
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078000	21.000000	0.000000
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.243750	24.000000	0.000000
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.372500	29.000000	0.000000
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.626250	41.000000	1.000000
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	81.000000	1.000000

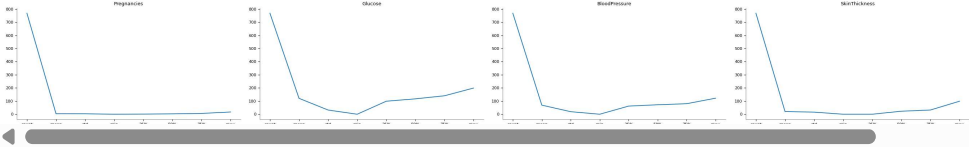
Distributions




2-d distributions



Values



```
print("\nMissing Values in Each Column:")
print(df.isnull().sum())
```

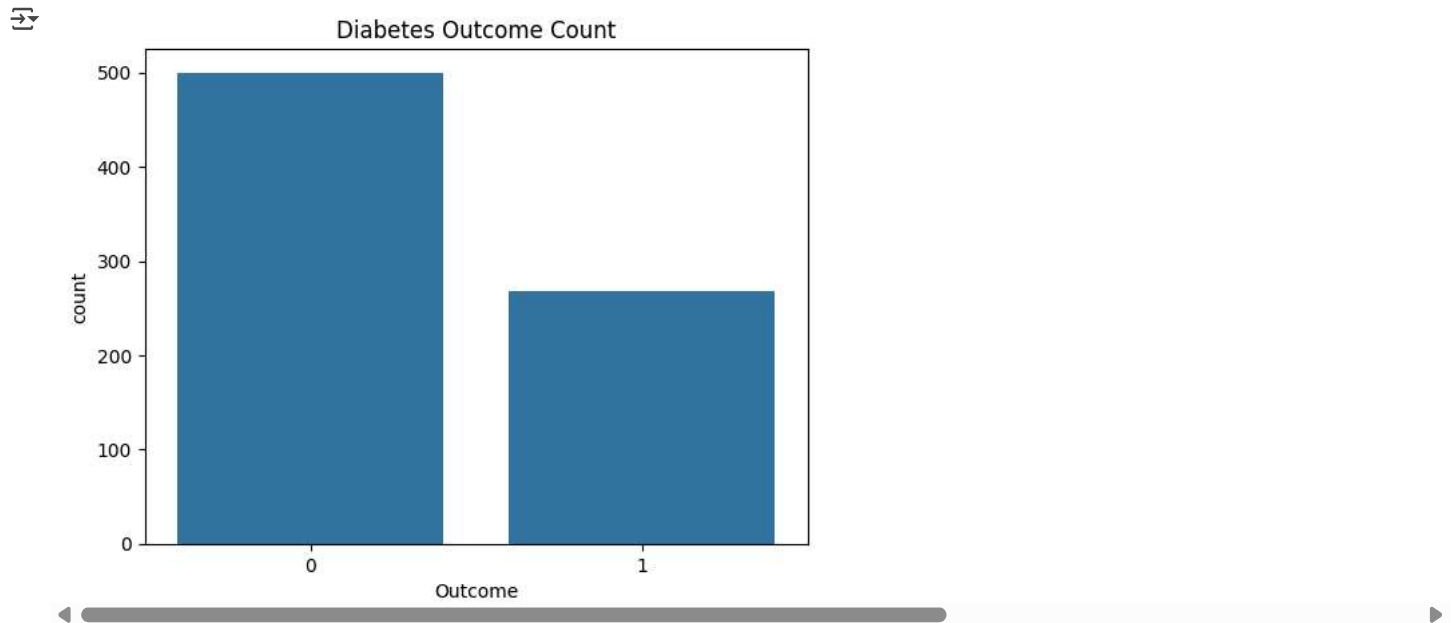
Missing Values in Each Column:

```
Pregnancies      0
Glucose          0
BloodPressure    0
SkinThickness    0
Insulin          0
BMI              0
DiabetesPedigreeFunction  0
Age              0
Outcome          0
dtype: int64
```

Visualize the Dataset

1)Check Diabetes outcome count(0 or 1)

```
# Countplot of diabetic vs non-diabetic
sns.countplot(x='Outcome', data=df)
plt.title("Diabetes Outcome Count")
plt.show()
```

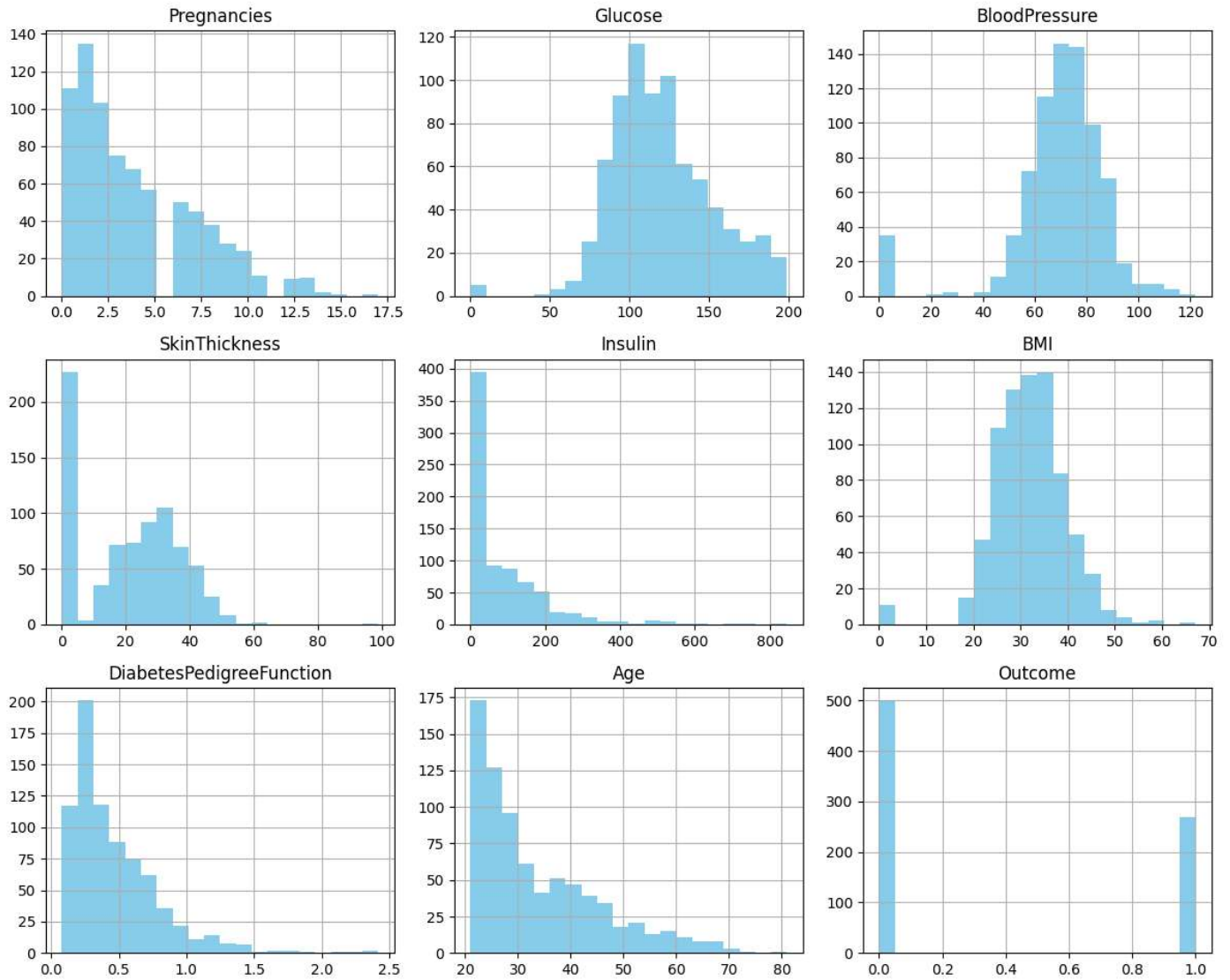


2)Show Histogram for all features

```
# Histogram of features
df.hist(bins=20, figsize=(12,10), color='skyblue')
plt.suptitle("Feature Distributions")
plt.tight_layout()
plt.show()
```

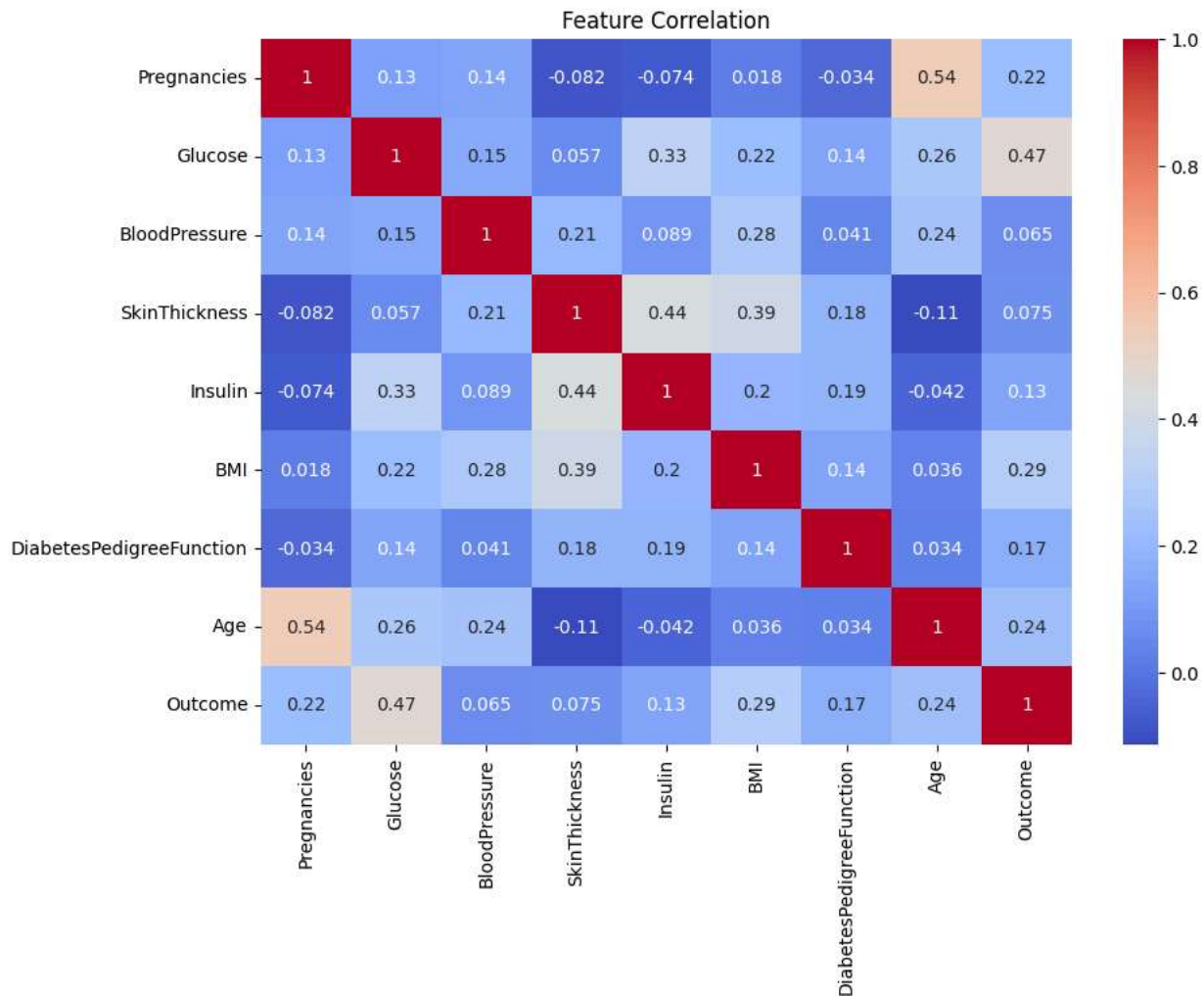


Feature Distributions



3)Correlation heatmap

```
# Correlation heatmap
plt.figure(figsize=(10,7))
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
plt.title("Feature Correlation")
plt.show()
```



Split features and labels

```
X = df.drop('Outcome', axis=1)
y = df['Outcome']

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42)
```

Feature scaling

```
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

Train SVM model

```
model = SVC(kernel='linear')
model.fit(X_train, y_train)
```



SVC

```
SVC(kernel='linear')
```

Evaluate the model

```
y_pred = model.predict(X_test)

print("Accuracy Score:", accuracy_score(y_test, y_pred))
print("\nConfusion Matrix:\n", confusion_matrix(y_test, y_pred))
```

```
print("\nClassification Report:\n", classification_report(y_test,y_pred))
```

➦ Accuracy Score: 0.7597402597402597

Confusion Matrix:
[[81 18]
[19 36]]

Classification Report:

	precision	recall	f1-score	support
0	0.81	0.82	0.81	99
1	0.67	0.65	0.66	55
accuracy			0.76	154
macro avg	0.74	0.74	0.74	154
weighted avg	0.76	0.76	0.76	154