

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
sns.set(style="ticks")
```

Загрузка и первичный анализ данных

```
from google.colab import files
```

```
uploaded = files.upload()
```

```
for filename in uploaded.keys():
    print('Uploaded file "{name}" with length
bytes'.format(name='Anaemia Prediction Dataset.csv',
length=len(uploaded['Anaemia Prediction Dataset.csv'])))
```

<IPython.core.display.HTML object>

Saving Anaemia Prediction Dataset.csv to Anaemia Prediction Dataset.csv

Uploaded file "Anaemia Prediction Dataset.csv" with length bytes

```
data = pd.read_csv('Anaemia Prediction Dataset.csv', sep=",")
```

data

```
{
  "summary": {
    "name": "data",
    "rows": 500,
    "fields": [
      {
        "column": "Number",
        "properties": {
          "dtype": "number",
          "std": 144,
          "min": 1,
          "max": 500,
          "num_unique_values": 500,
          "samples": [
            362, 74, 375
          ],
          "semantic_type": ""
        },
        "description": ""
      },
      {
        "column": "Sex",
        "properties": {
          "dtype": "category",
          "num_unique_values": 4,
          "samples": [
            "F", "F", "M"
          ],
          "semantic_type": ""
        },
        "description": ""
      },
      {
        "column": "%Red Pixel",
        "properties": {
          "dtype": "number",
          "std": 2.9500168444205075,
          "min": 36.8,
          "max": 56.85,
          "num_unique_values": 404,
          "samples": [
            45.92, 41.95, 50.47
          ],
          "semantic_type": ""
        },
        "description": ""
      },
      {
        "column": "%Green pixel",
        "properties": {
          "dtype": "number",
          "std": 1.8446234284497123,
          "min": 24.15,
          "max": 33.6,
          "num_unique_values": 358,
          "samples": [
            27.94, 28.22, 28.3
          ]
        },
        "description": ""
      }
    ]
  }
}
```

```

{"semantic_type": "\\",
  },
  {"column": "%Blue pixel",
    "properties": {
      "dtype": "number",
      "std": 2.2159130731812993,
      "min": 17.95,
      "max": 31.3,
      "num_unique_values": 371,
      "samples": [
        27.02,
        27.69,
        24.08
      ]
    },
    "semantic_type": "\\",
    "description": "\\",
  },
  {"column": "Hb",
    "properties": {
      "dtype": "number",
      "std": 3.04094282361787,
      "min": 4.0,
      "max": 18.55,
      "num_unique_values": 345,
      "samples": [
        4.83,
        7.13,
        15.56
      ]
    },
    "semantic_type": "\\",
    "description": "\\",
  },
  {"column": "Anaemic",
    "properties": {
      "dtype": "category",
      "num_unique_values": 2,
      "samples": [
        "No",
        "Yes"
      ]
    },
    "semantic_type": "\\",
    "description": "\\"
  }
],
"type": "dataframe",
"variable_name": "data"}

```

```

data = data.drop(columns='Number')
data

```

```

{"summary": {
  "name": "data",
  "rows": 500,
  "fields": [
    {
      "column": "Sex",
      "properties": {
        "dtype": "category",
        "num_unique_values": 4,
        "samples": [
          "F",
          "F ",
          "M"
        ],
        "semantic_type": "\\",
        "description": "\\"
      },
      "column": "%Red Pixel",
      "properties": {
        "dtype": "number",
        "std": 2.9500168444205075,
        "min": 36.8,
        "max": 56.85,
        "num_unique_values": 404,
        "samples": [
          45.92,
          41.95,
          50.47
        ]
      },
      "semantic_type": "\\",
      "description": "\\"
    },
    {
      "column": "%Green pixel",
      "properties": {
        "dtype": "number",
        "std": 1.8446234284497123,
        "min": 24.15,
        "max": 33.6,
        "num_unique_values": 358,
        "samples": [
          27.94,
          28.22,
          28.3
        ]
      },
      "semantic_type": "\\",
      "description": "\\"
    },
    {
      "column": "%Blue pixel",
      "properties": {
        "dtype": "number",
        "std": 2.2159130731812993,
        "min": 17.95,
        "max": 31.3,
        "num_unique_values": 371,
        "samples": [
          27.02,
          27.69,
          24.08
        ]
      },
      "semantic_type": "\\",
      "description": "\\"
    },
    {
      "column": "Hb",
      "properties": {
        "dtype": "number",
        "std": 3.04094282361787,
        "min": 4.0,
        "max": 18.55,
        "num_unique_values": 345,
        "samples": [
          4.83,
          7.13,
          15.56
        ]
      },
      "semantic_type": "\\",
      "description": "\\"
    }
  ]
}

```

```

{"description\": \"\"\n        }\n    },\n    {\n        \"column\":  

\"Anaemic\", \n        \"properties\": {\n            \"dtype\":  

\"category\", \n            \"num_unique_values\": 2, \n            \"samples\":  

[\n                \"No\", \n                \"Yes\" \n            ], \n            \"semantic_type\": \"\", \n            \"description\": \"\" \n        } \n    } \n ] \n }\", \"type\": \"dataframe\", \"variable_name\": \"data\"}

```

```
data.Sex
```

```

0      M
1      F
2      F
3      F
4      M
..

```

```

495    F
496    F
497    F
498    F
499    F

```

```
Name: Sex, Length: 500, dtype: object
```

```
data.Sex.unique()
```

```
array(['M', 'F', 'M ', 'F '], dtype=object)
```

```
data['Sex'] = data['Sex'].map({'F': 0, 'F ': 0, 'M': 1, 'M ': 1})
```

```
data.Sex.unique()
```

```
array([1, 0])
```

```
data.Anaemic.unique()
```

```
array(['Yes', 'No'], dtype=object)
```

```
data['Anaemic'] = data['Anaemic'].map({'Yes': 1, 'No': 0})
```

```
data.dtypes
```

```

Sex                int64
%Red Pixel         float64
%Green pixel       float64
%Blue pixel        float64
Hb                 float64
Anaemic            int64
dtype: object

```

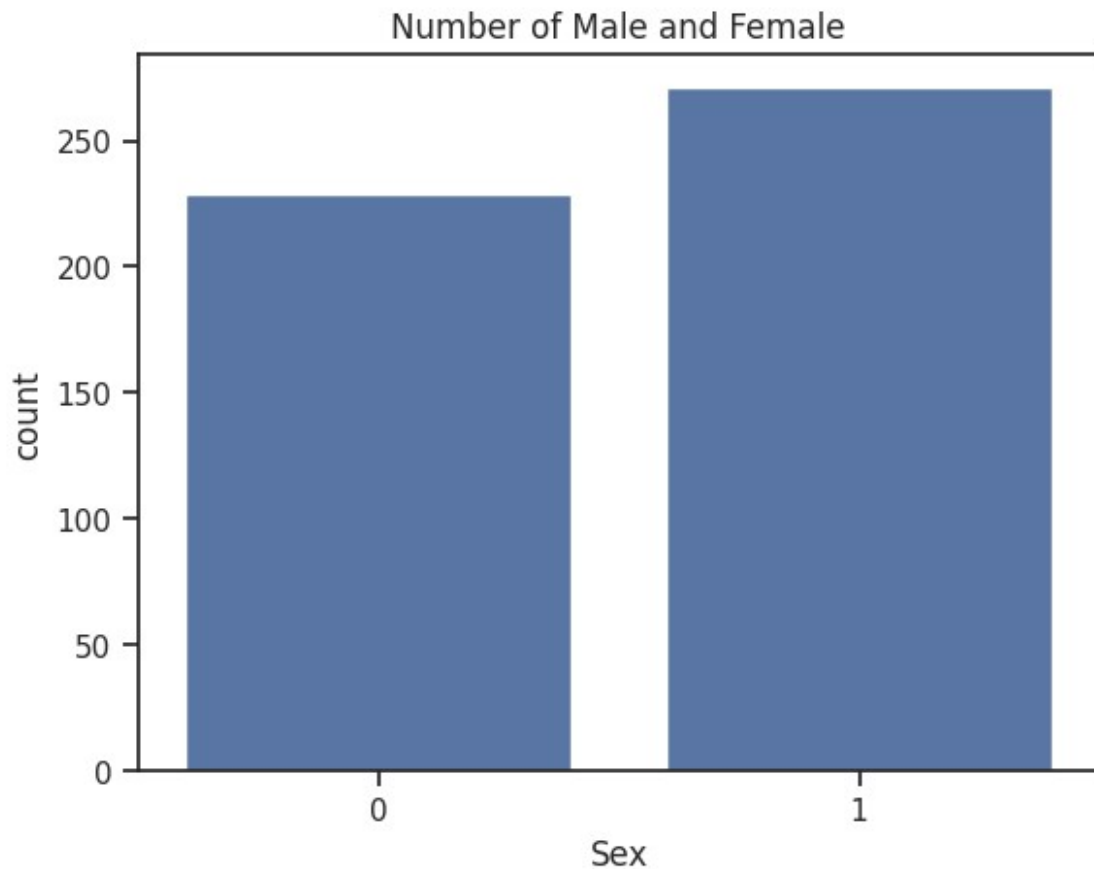
EDA

```

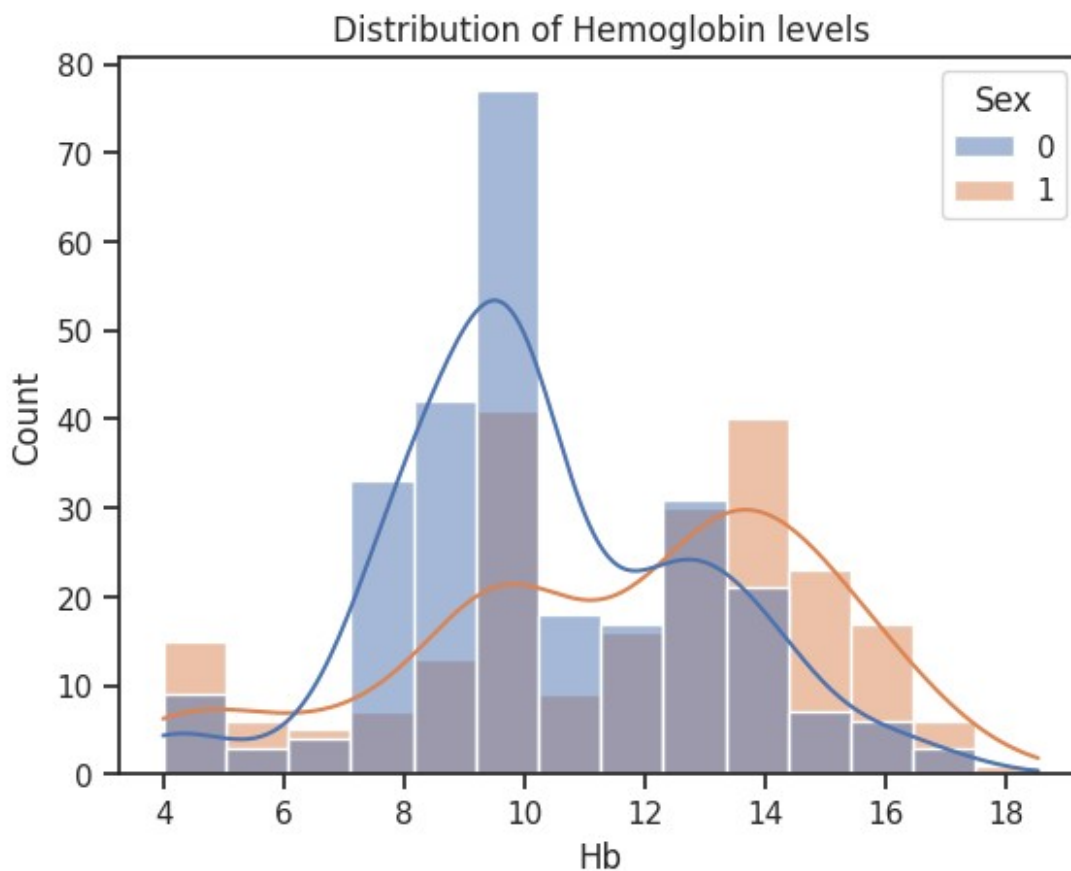
import matplotlib.pyplot as plt
import seaborn as sns

```

```
gender_count = data.Sex.value_counts() # data['Sex'].value_counts()
sns.barplot(data = data, x = 'Sex', y = gender_count)
plt.title('Number of Male and Female')
Text(0.5, 1.0, 'Number of Male and Female')
```



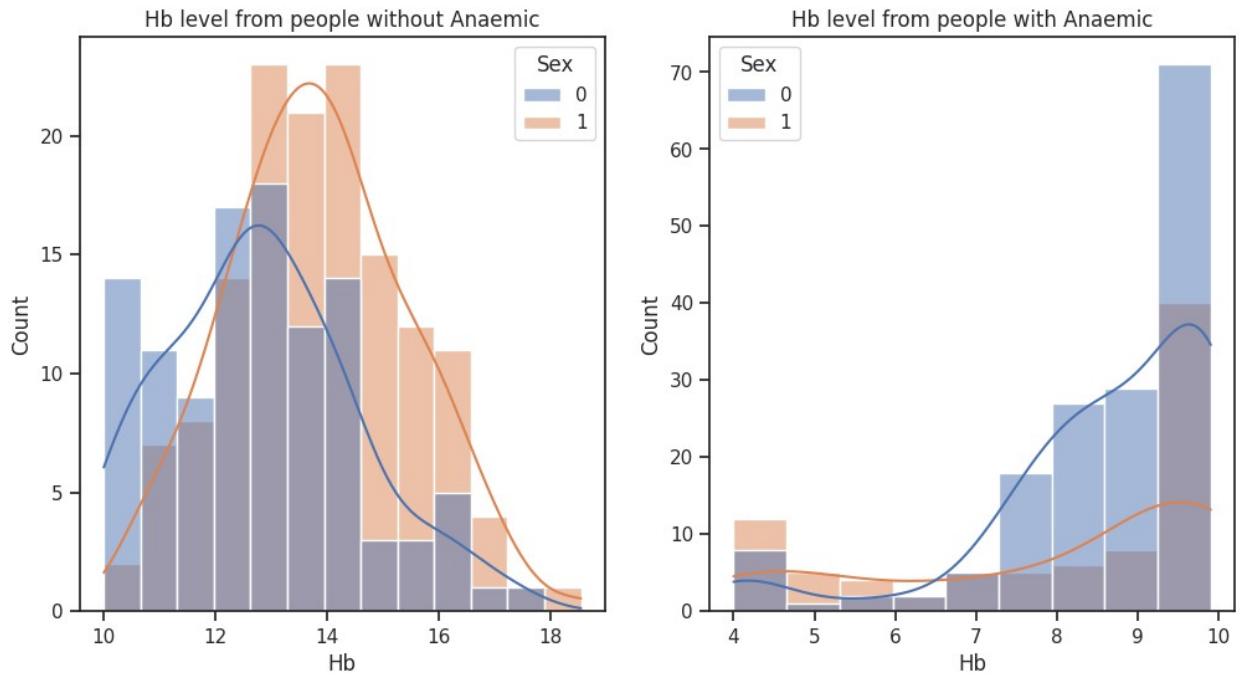
```
sns.histplot(data=data, x='Hb', hue='Sex', kde=True)
plt.title("Distribution of Hemoglobin levels")
Text(0.5, 1.0, 'Distribution of Hemoglobin levels')
```



```
anaemic_data = data[data['Anaemic'] == 1]
not_anaemic_data = data[data['Anaemic'] == 0]

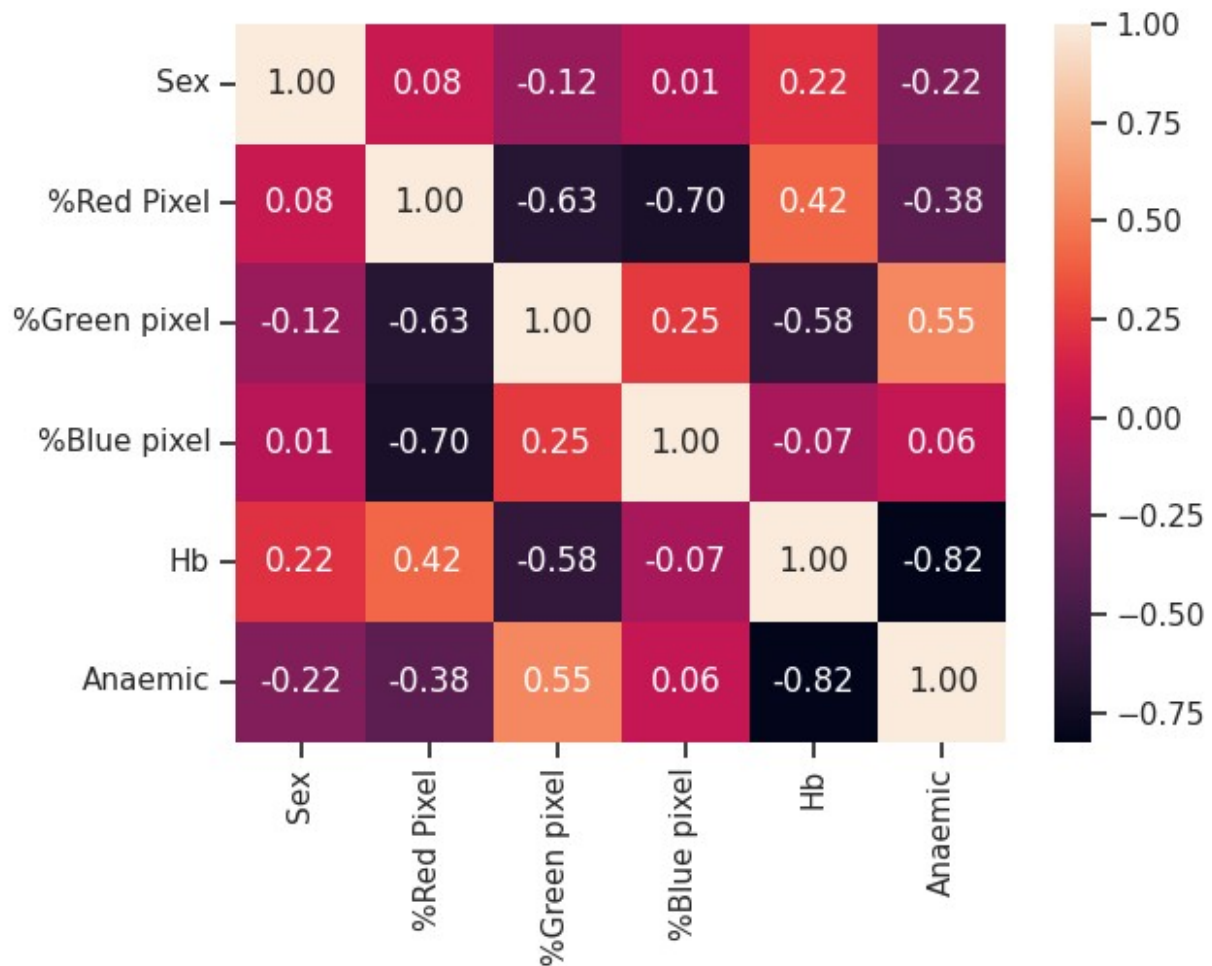
plt.figure(figsize=(12, 6))
plt.subplot(1,2,1)
sns.histplot(data=not_anaemic_data, x='Hb', hue='Sex', kde=True)
plt.title('Hb level from people without Anaemic')

plt.subplot(1,2,2)
sns.histplot(data=anaemic_data, x='Hb', hue='Sex', kde=True)
plt.title('Hb level from people with Anaemic')
plt.show()
```



people who didn't diagnose with anaemic have higher Hemoglobin level than people who diagnose with Anaemic.

```
corr = data.corr()
sns.heatmap(corr, fmt='.2f', annot=True)
<Axes: >
```



Split Dataset

```
X = data.drop(columns='Anaemic')
y = data['Anaemic']

X.shape, y.shape

((500, 5), (500,))

from sklearn.model_selection import train_test_split

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

Define the Model

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.linear_model import LogisticRegression
```

```

from sklearn.metrics import accuracy_score, classification_report,
confusion_matrix

models = {
    'RF': RandomForestClassifier(),
    'DT': DecisionTreeClassifier(),
    'LogR': LogisticRegression()
}

for model_name, model in models.items():
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    accuracy = accuracy_score(y_pred, y_test)
    print(f"Model name: {model_name}\nAccuracy: {accuracy}\n
n{classification_report(y_pred, y_test)}")

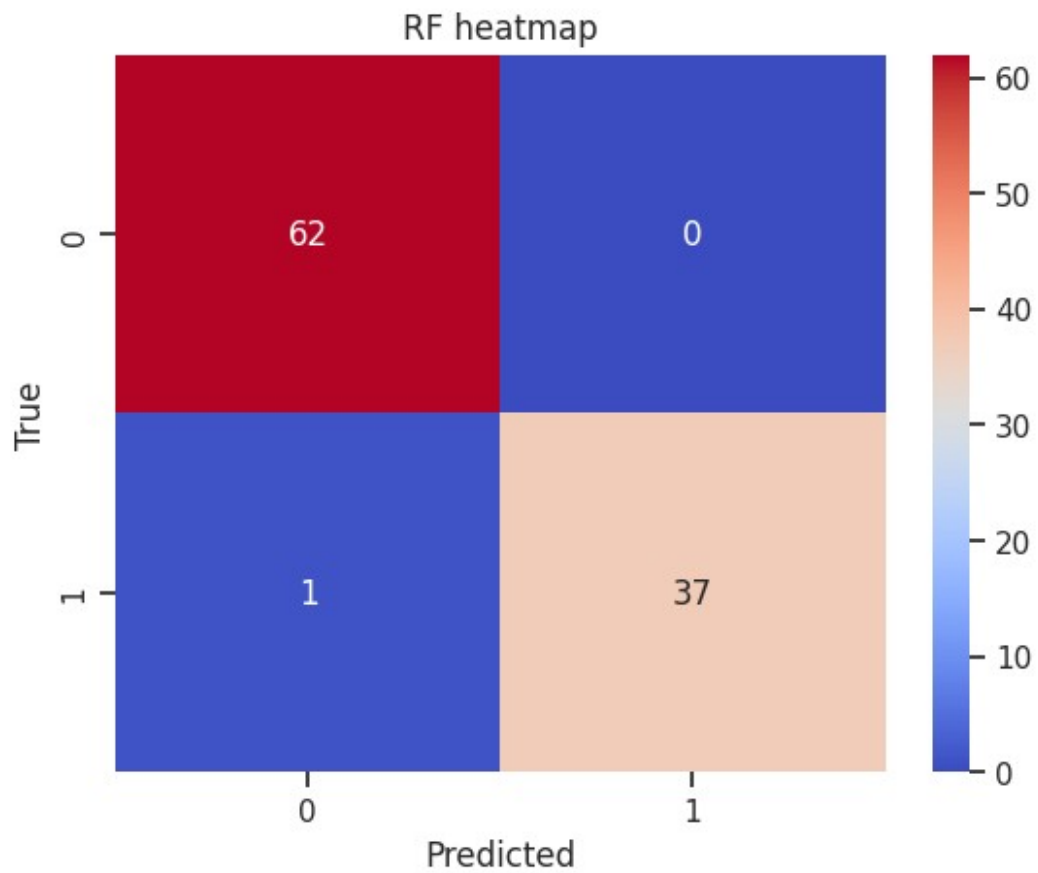
    cm = confusion_matrix(y_pred, y_test)
    sns.heatmap(cm, annot=True, fmt='d', cmap='coolwarm')
    plt.xlabel('Predicted')
    plt.ylabel('True')
    plt.title(f"{model_name} heatmap")
    plt.show()

```

Model name: RF

Accuracy: 0.99

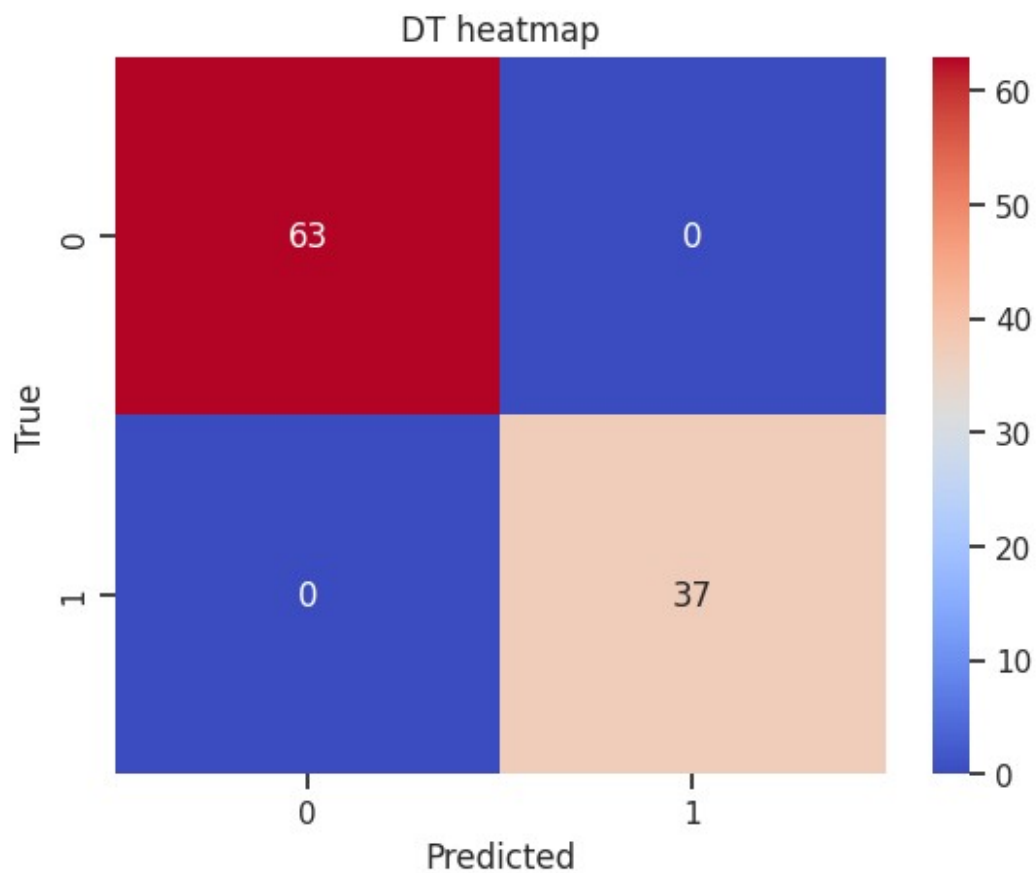
	precision	recall	f1-score	support
0	0.98	1.00	0.99	62
1	1.00	0.97	0.99	38
accuracy			0.99	100
macro avg	0.99	0.99	0.99	100
weighted avg	0.99	0.99	0.99	100



Model name: DT

Accuracy: 1.0

	precision	recall	f1-score	support
0	1.00	1.00	1.00	63
1	1.00	1.00	1.00	37
accuracy			1.00	100
macro avg	1.00	1.00	1.00	100
weighted avg	1.00	1.00	1.00	100



Model name: LogR

Accuracy: 0.98

	precision	recall	f1-score	support
0	0.97	1.00	0.98	61
1	1.00	0.95	0.97	39
accuracy			0.98	100
macro avg	0.98	0.97	0.98	100
weighted avg	0.98	0.98	0.98	100

