

```
In [1]: 1 import matplotlib.pyplot as plt  
2 import pandas as pd  
3 import numpy as np  
4 import seaborn as sns
```

```
In [2]: 1 df=pd.read_csv("heart_failure.csv")
```

```
In [3]: 1 df.head()
```

Out[3]:

	age	anaemia	creatinine_phosphokinase	diabetes	ejection_fraction	high_blood_pressure	platelets	serum_creatinine
0	75.0	0	582	0	20		1 265000.00	1.9
1	55.0	0	7861	0	38		0 263358.03	1.1
2	65.0	0	146	0	20		0 162000.00	1.3
3	50.0	1	111	0	20		0 210000.00	1.9
4	65.0	1	160	1	20		0 327000.00	2.7

```
In [4]: 1 df.shape
```

Out[4]: (299, 13)

```
In [5]: 1 df.describe()
```

Out[5]:

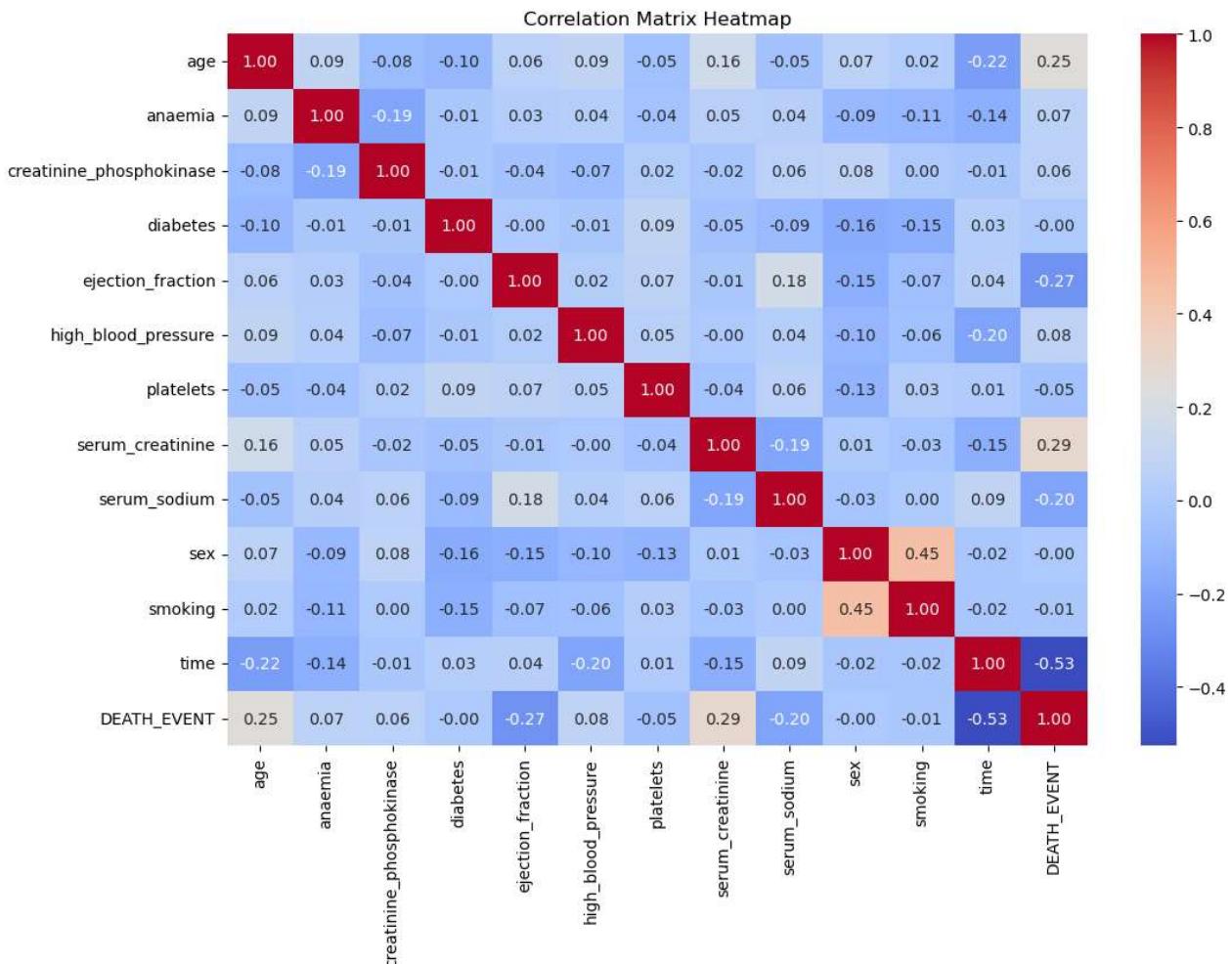
	age	anaemia	creatinine_phosphokinase	diabetes	ejection_fraction	high_blood_pressure	platelet
count	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000
mean	60.833893	0.431438	581.839465	0.418060	38.083612	0.351171	263358.02926
std	11.894809	0.496107	970.287881	0.494067	11.834841	0.478136	97804.23686
min	40.000000	0.000000	23.000000	0.000000	14.000000	0.000000	25100.00000
25%	51.000000	0.000000	116.500000	0.000000	30.000000	0.000000	212500.00000
50%	60.000000	0.000000	250.000000	0.000000	38.000000	0.000000	262000.00000
75%	70.000000	1.000000	582.000000	1.000000	45.000000	1.000000	303500.00000
max	95.000000	1.000000	7861.000000	1.000000	80.000000	1.000000	850000.00000

Data Exploration

```
In [246]: 1 # Data types and missing values
2 df.info()
```

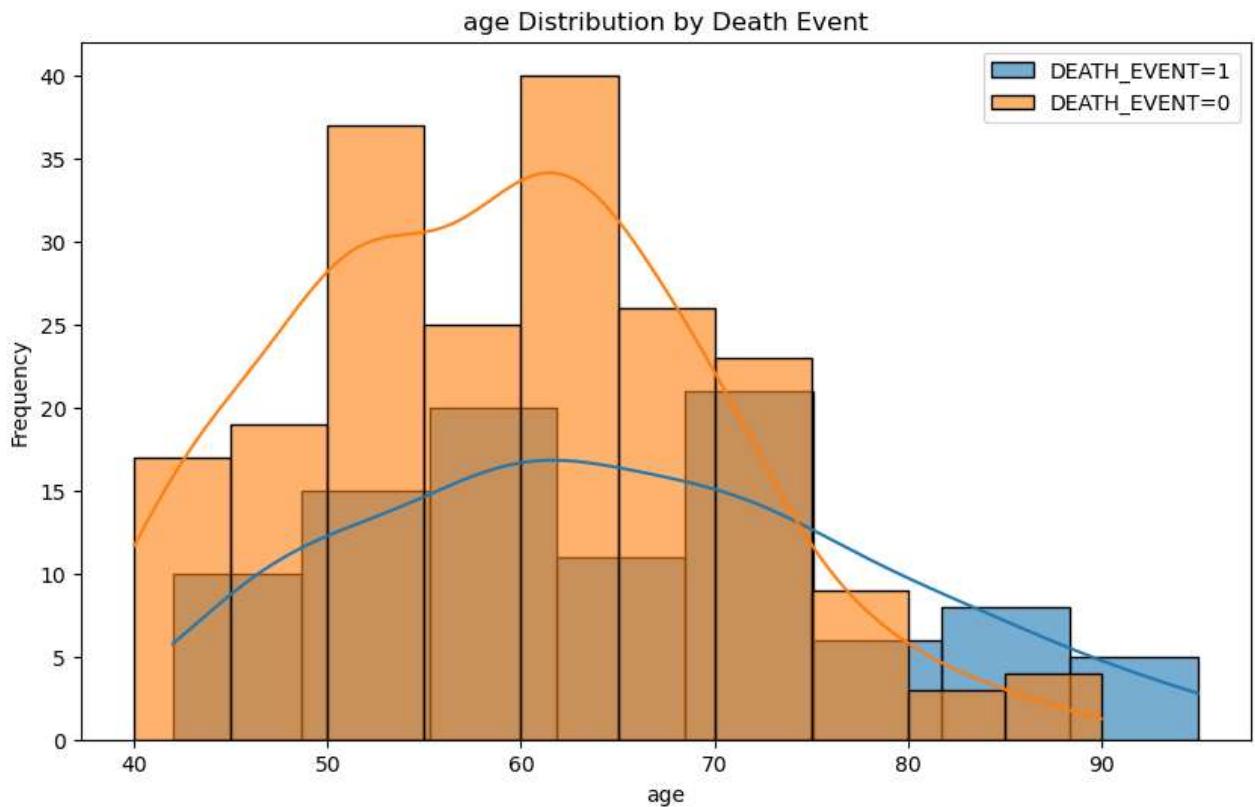
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 299 entries, 0 to 298
Data columns (total 13 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   age              299 non-null    float64
 1   anaemia          299 non-null    int64  
 2   creatinine_phosphokinase 299 non-null    int64  
 3   diabetes          299 non-null    int64  
 4   ejection_fraction 299 non-null    int64  
 5   high_blood_pressure 299 non-null    int64  
 6   platelets         299 non-null    float64
 7   serum_creatinine  299 non-null    float64
 8   serum_sodium      299 non-null    int64  
 9   sex               299 non-null    int64  
 10  smoking           299 non-null    int64  
 11  time              299 non-null    int64  
 12  DEATH_EVENT       299 non-null    int64  
dtypes: float64(3), int64(10)
memory usage: 30.5 KB
```

```
In [247]: 1 # Correlation matrix heatmap
2 correlation_matrix = df.corr()
3 plt.figure(figsize=(12, 8))
4 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
5 plt.title("Correlation Matrix Heatmap")
6 plt.show()
```

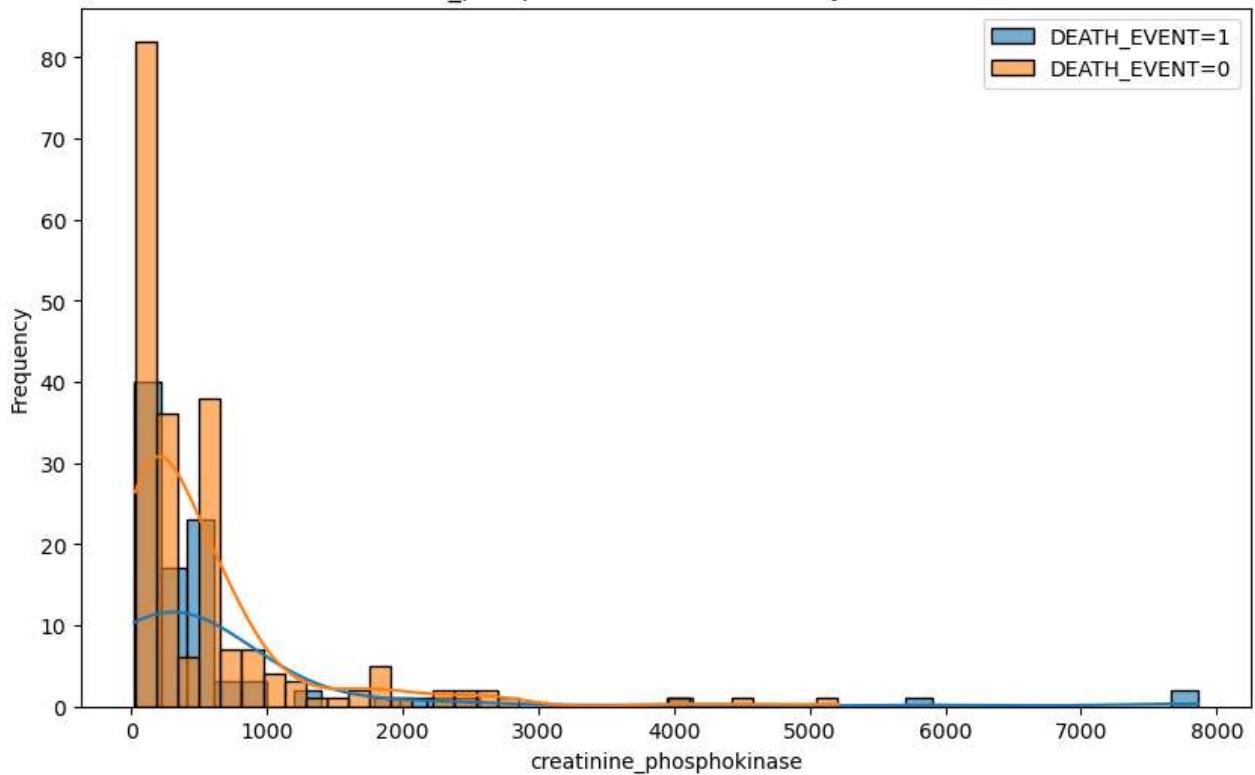


In [250]:

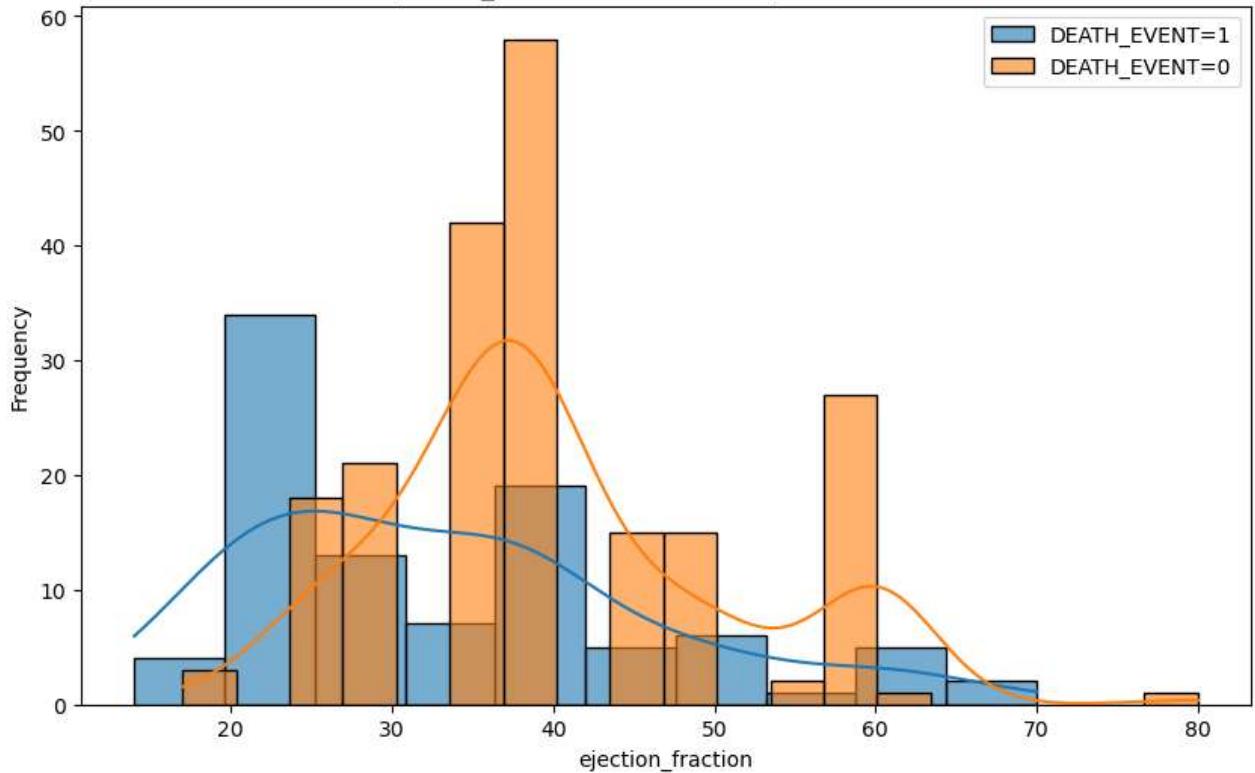
```
1 # Numerical features
2 numerical_features = ['age', 'creatinine_phosphokinase', 'ejection_fraction', 'platelets', 'serum_creatinine', 'serum_sodium']
3
4 # Group by 'DEATH_EVENT' and create histograms for numerical features
5 for feature in numerical_features:
6     plt.figure(figsize=(10, 6))
7
8     # Create separate histograms for each 'DEATH_EVENT' category
9     for death_event in df['DEATH_EVENT'].unique():
10         sns.histplot(df[df['DEATH_EVENT'] == death_event][feature], kde=True, label=f'DEATH_EVENT={death_event}')
11
12     plt.title(f'{feature} Distribution by Death Event')
13     plt.xlabel(feature)
14     plt.ylabel('Frequency')
15     plt.legend()
16     plt.show()
17
```



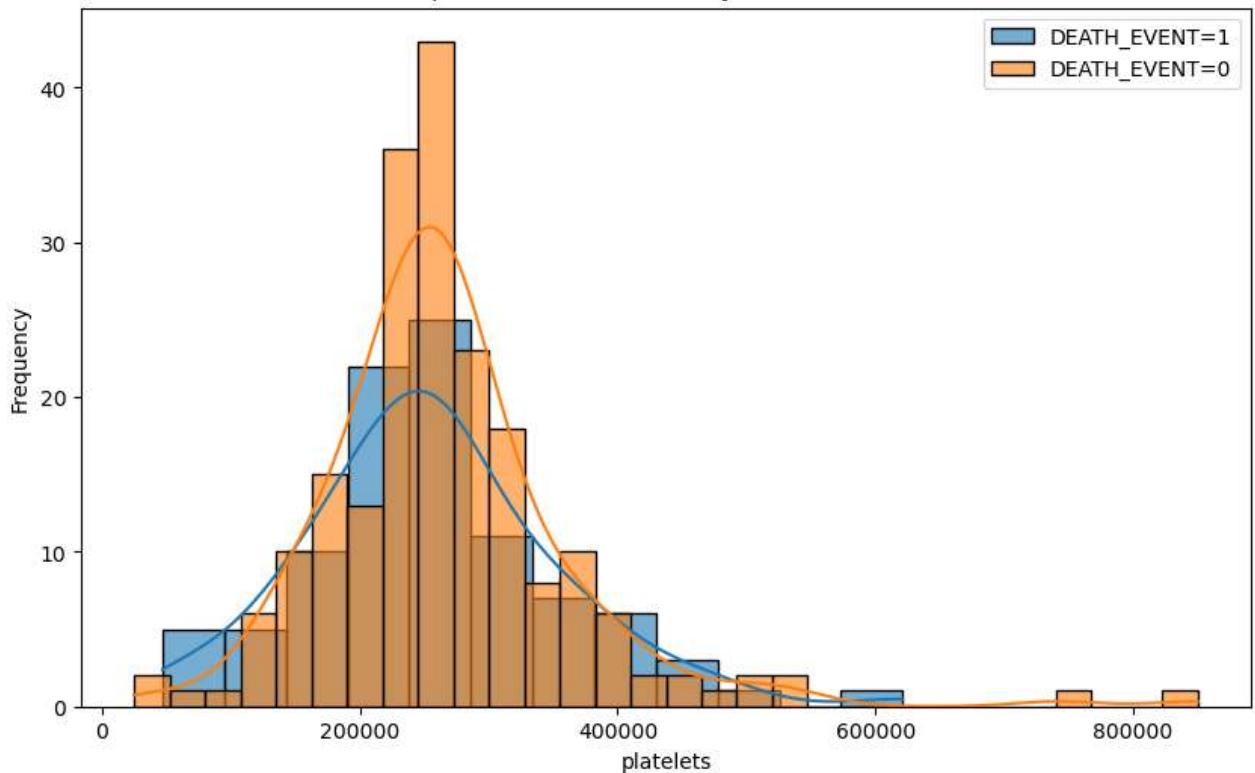
creatinine_phosphokinase Distribution by Death Event



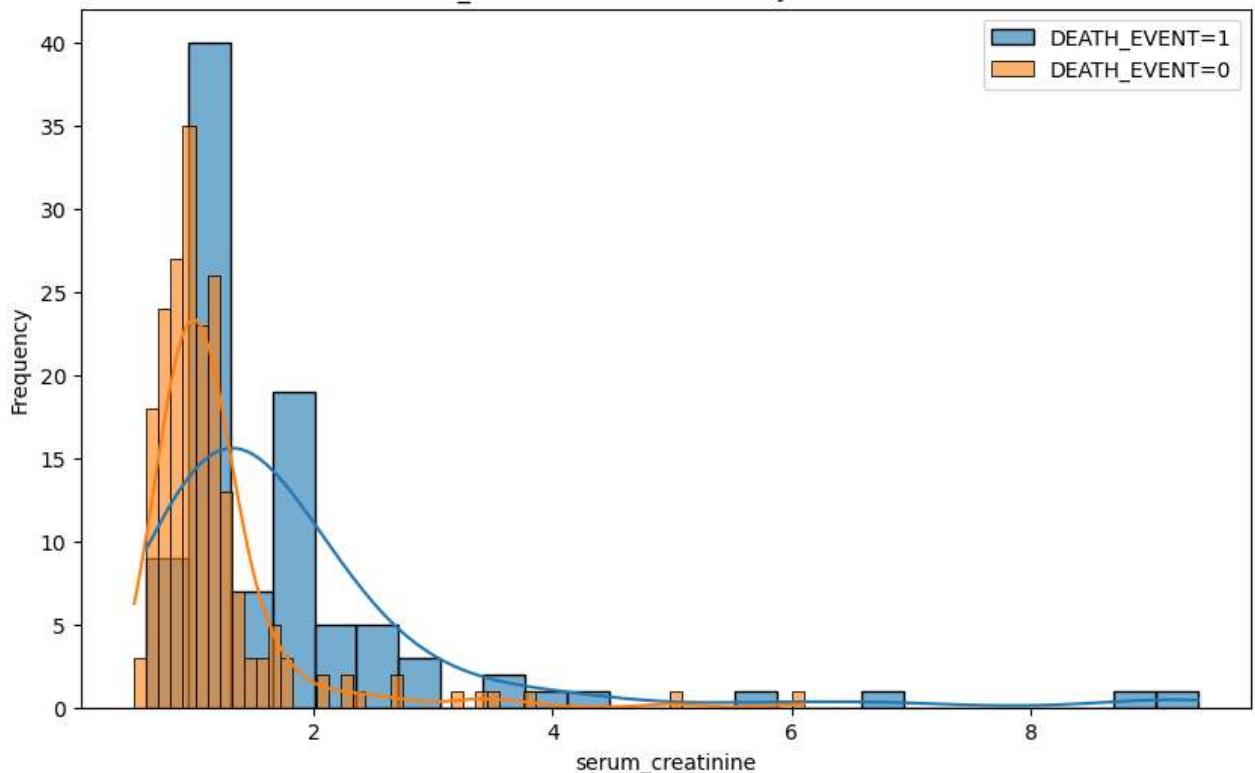
ejection_fraction Distribution by Death Event

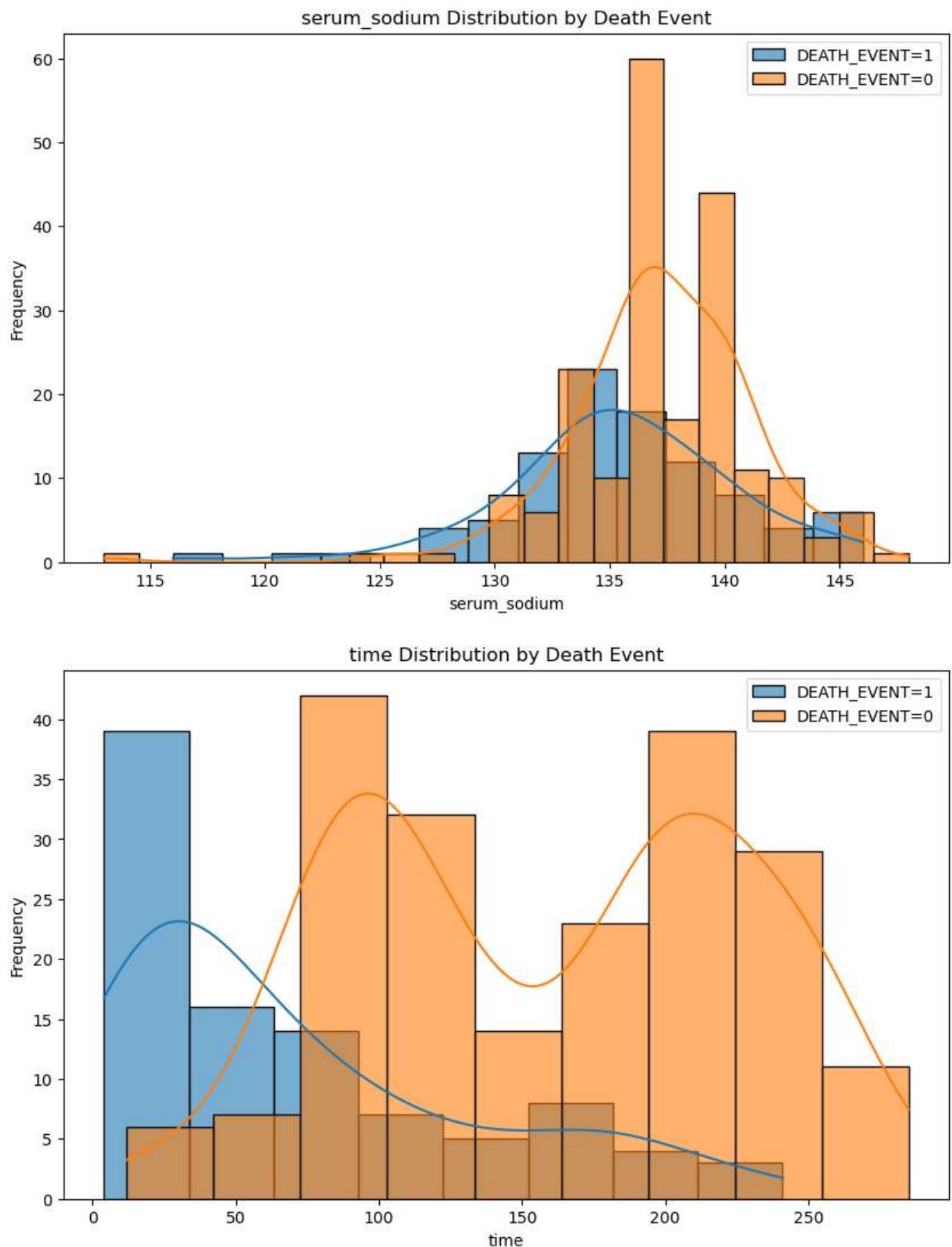


platelets Distribution by Death Event

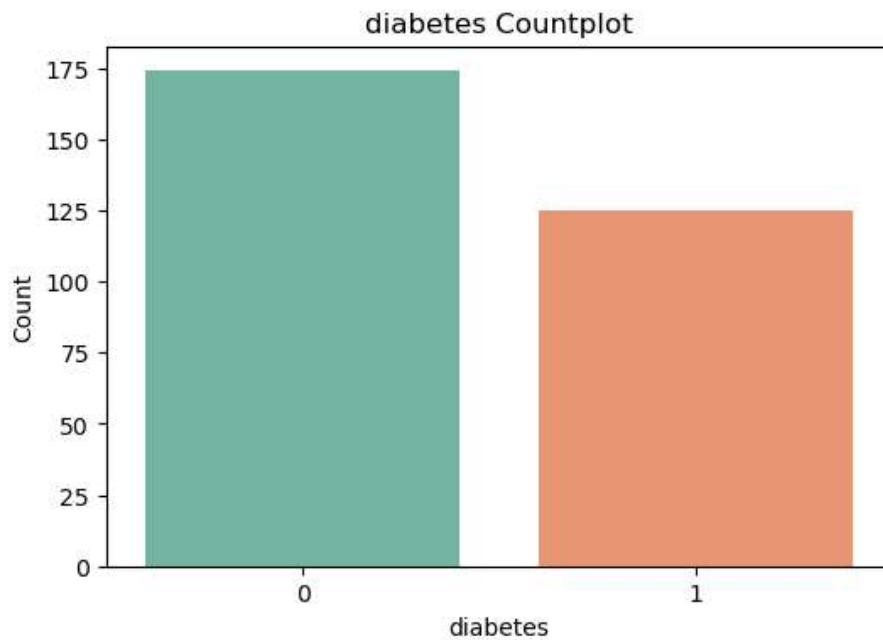
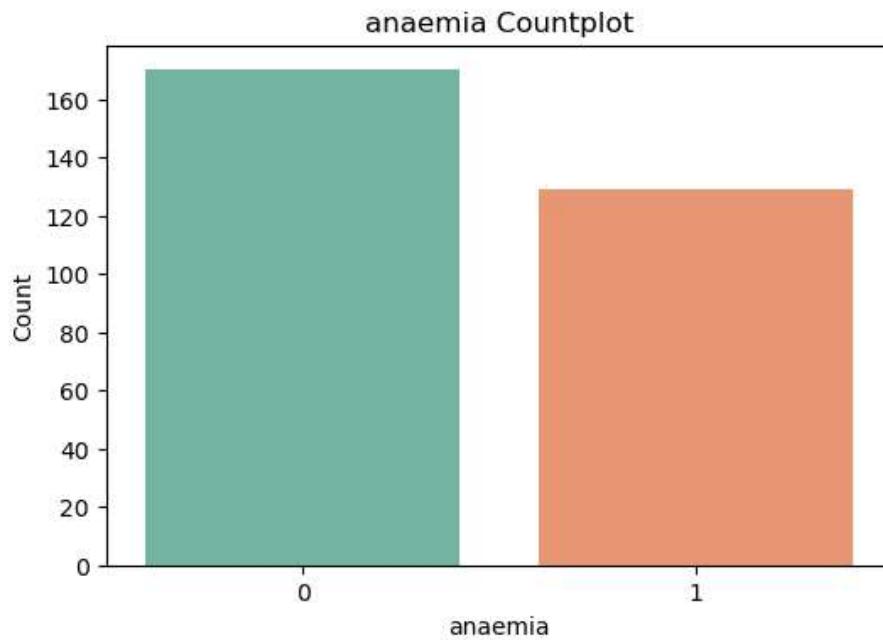


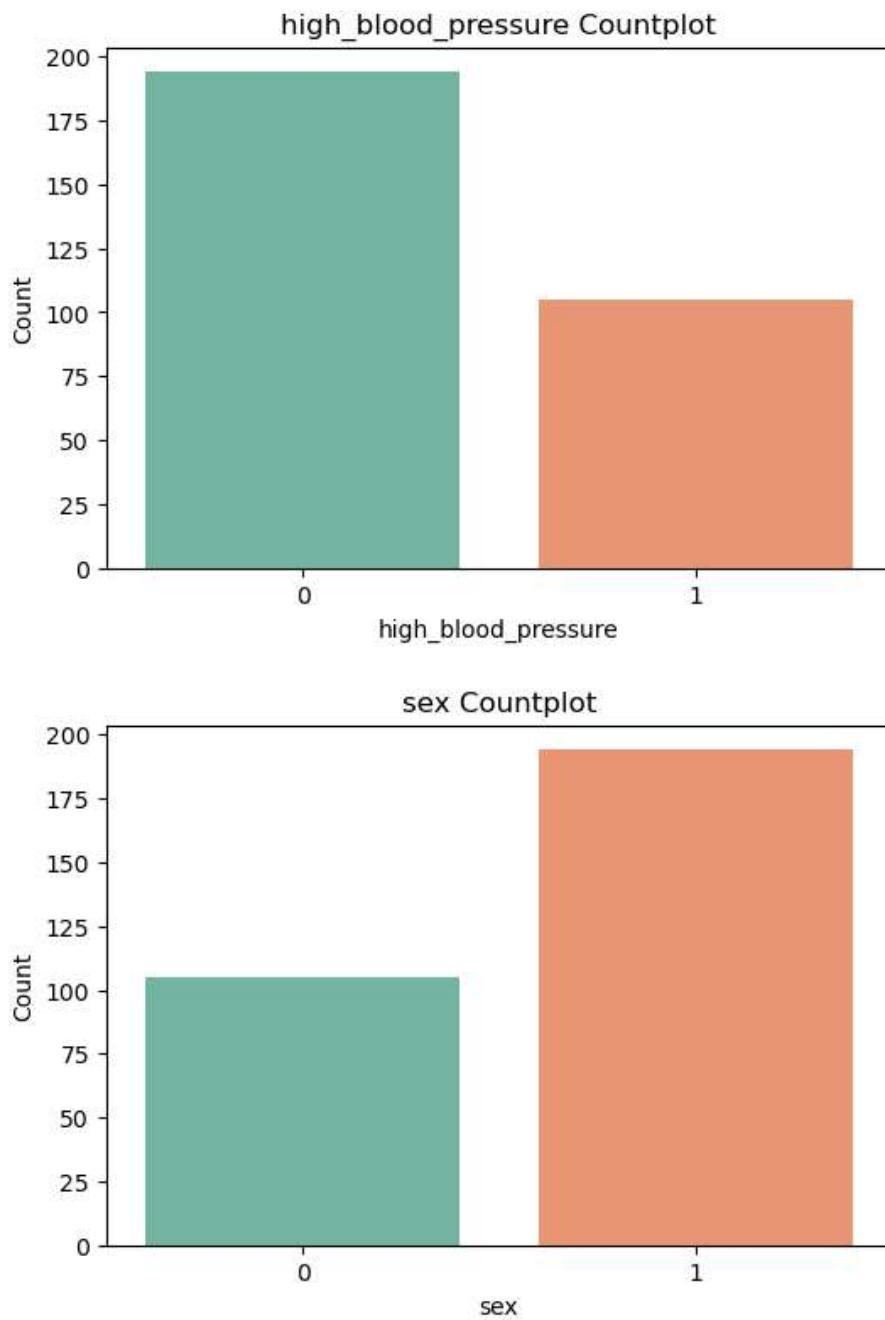
serum_creatinine Distribution by Death Event

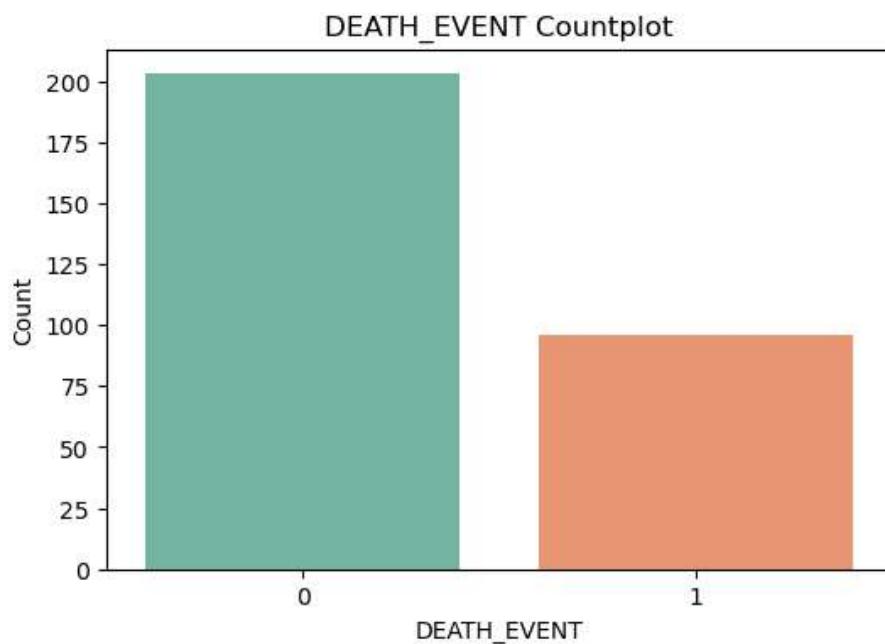
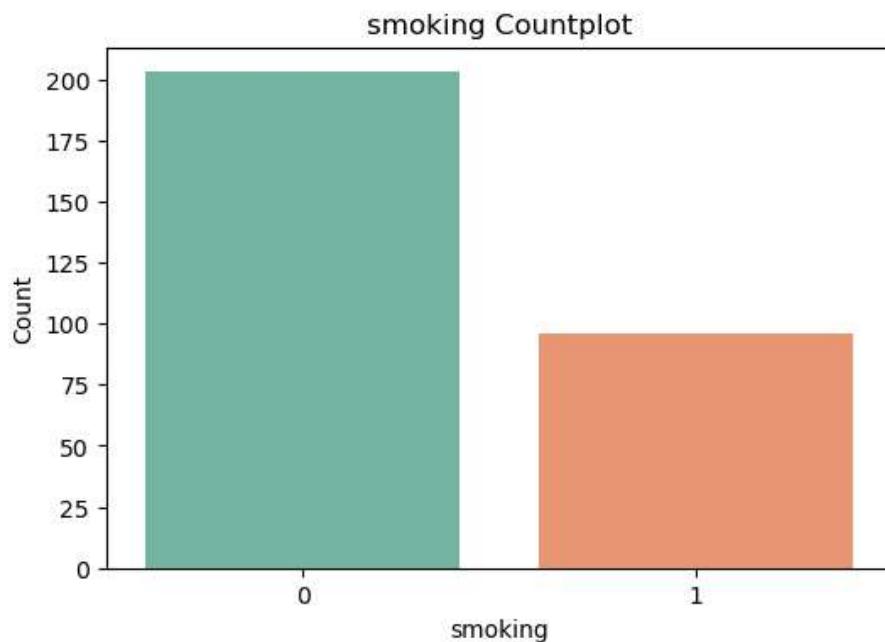




```
In [251]: 1 categorical_features = ['anaemia', 'diabetes', 'high_blood_pressure', 'sex', 'smoking', 'DEATH']
2 for feature in categorical_features:
3     plt.figure(figsize=(6, 4))
4     sns.countplot(data=df, x=feature, palette='Set2')
5     plt.title(f'{feature} Countplot')
6     plt.xlabel(feature)
7     plt.ylabel('Count')
8     plt.show()
```

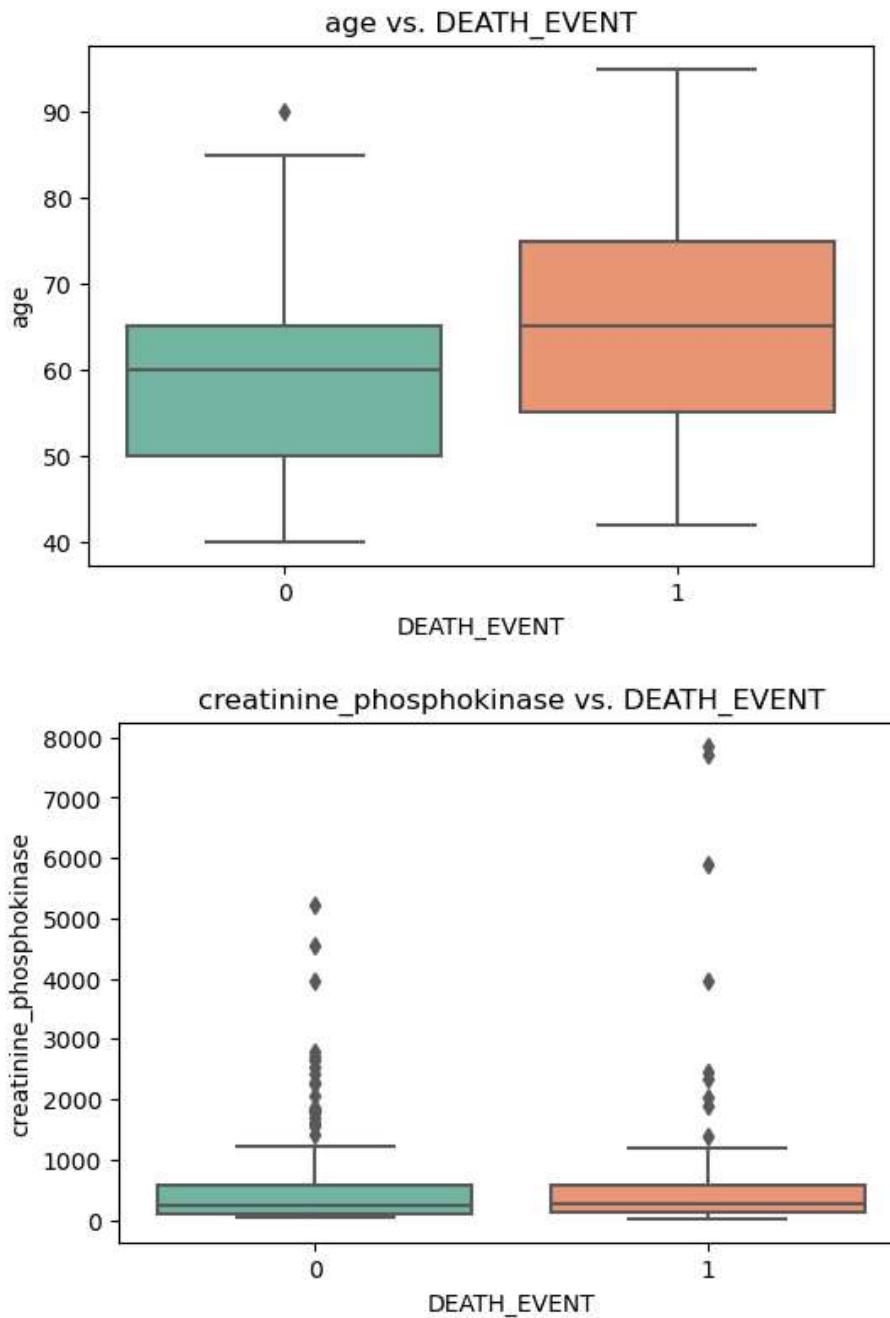


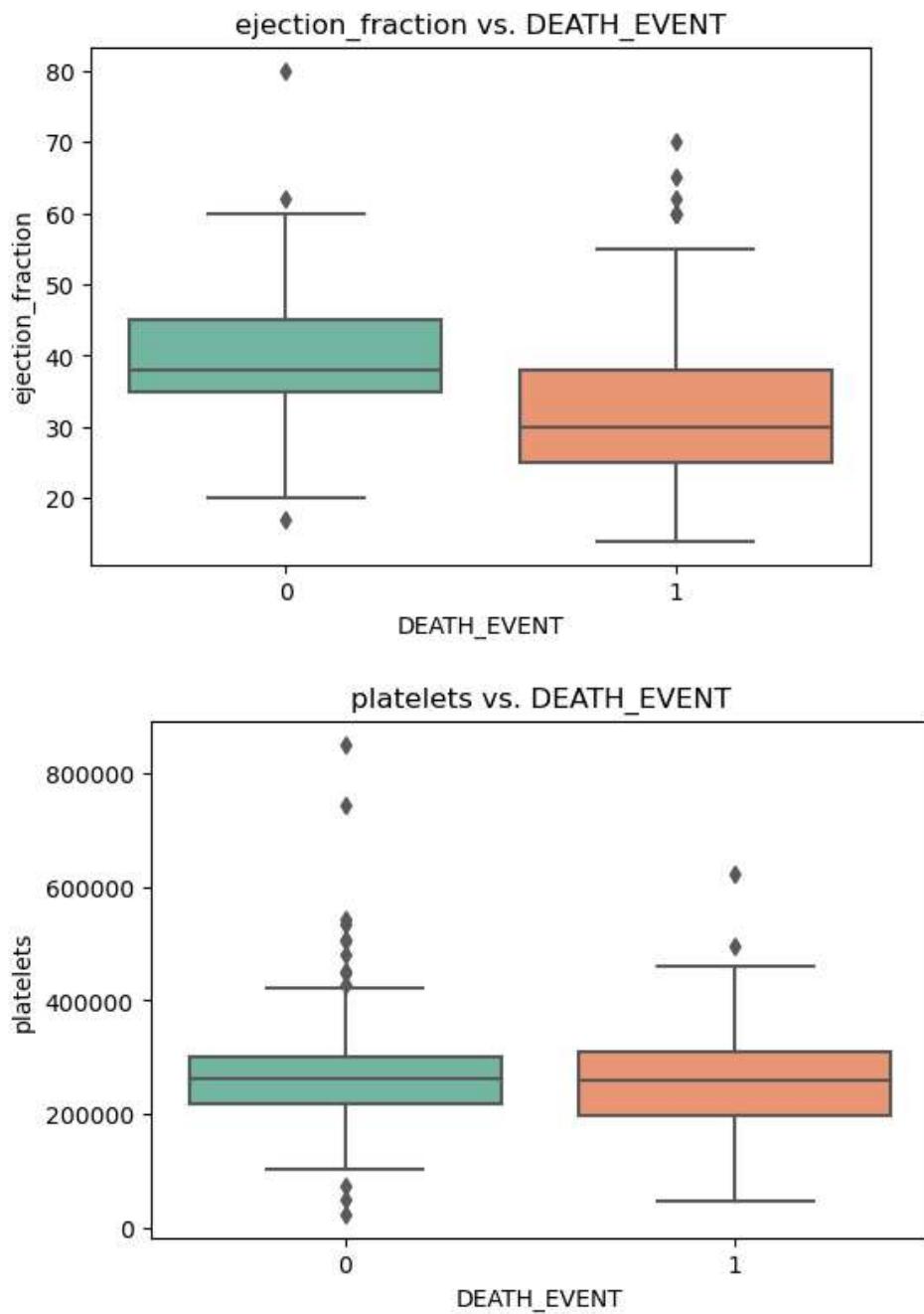


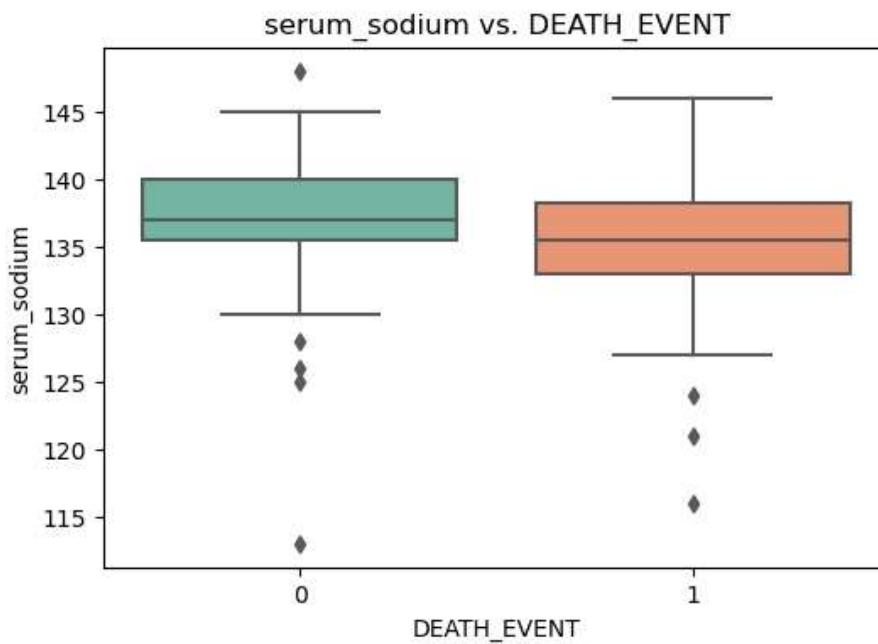
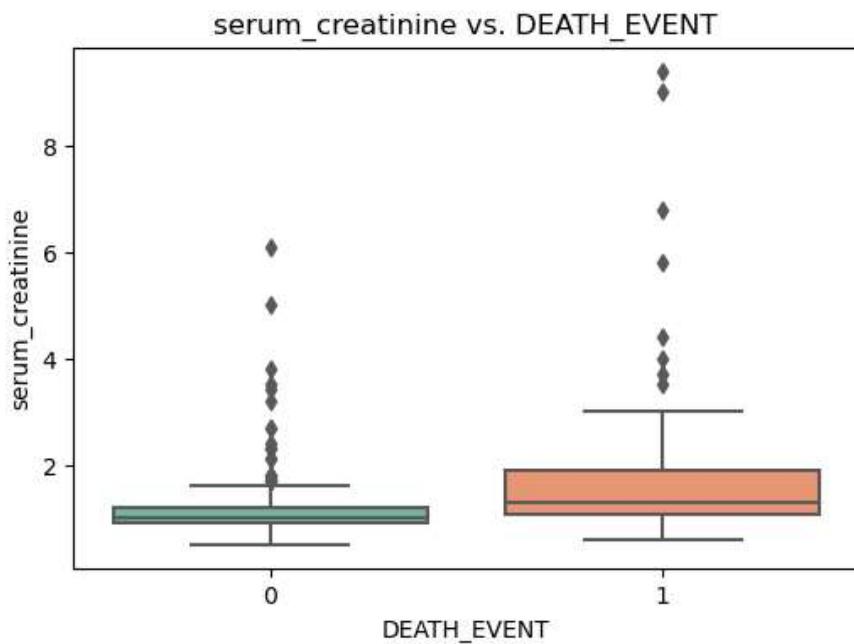


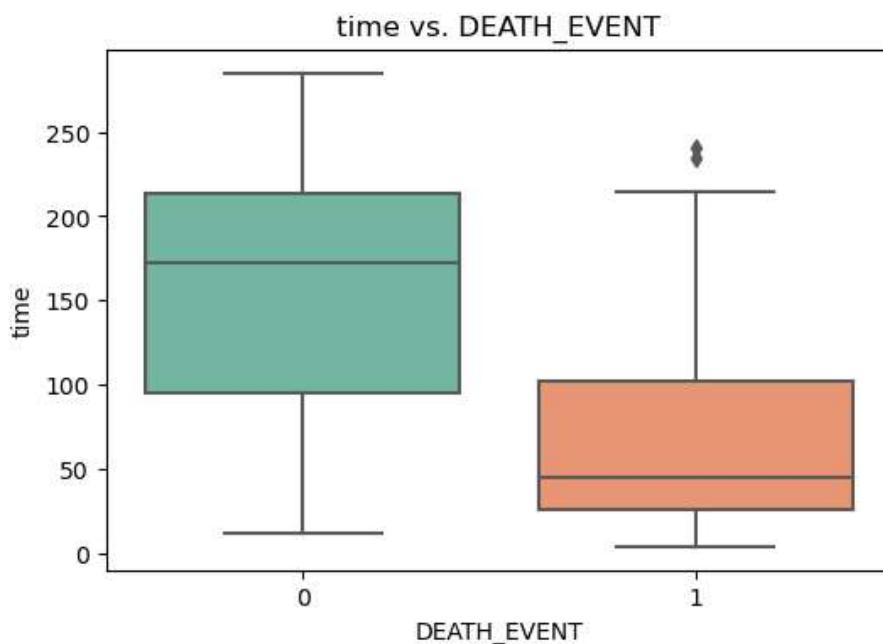
In [252]:

```
1 # BoxPlots for numerical features vs. target variable (DEATH_EVENT)
2 for feature in numerical_features:
3     plt.figure(figsize=(6, 4))
4     sns.boxplot(data=df, x='DEATH_EVENT', y=feature, palette='Set2')
5     plt.title(f'{feature} vs. DEATH_EVENT')
6     plt.xlabel('DEATH_EVENT')
7     plt.ylabel(feature)
8     plt.show()
```









Polynomial Regression

In [331]:

```
1 from sklearn.model_selection import GridSearchCV
2 from sklearn.linear_model import LinearRegression
3 from sklearn.preprocessing import PolynomialFeatures
4 from sklearn.pipeline import Pipeline
5 from sklearn.metrics import confusion_matrix
6 from sklearn.model_selection import train_test_split
7 from sklearn.preprocessing import StandardScaler
8 from sklearn import metrics
9 from sklearn.metrics import roc_curve, auc
```

In [253]:

```

1 # Splitting Data into feature variables and target variables
2 df_X = df.drop(['DEATH_EVENT'], axis=1) # Features
3 df_Y = df["DEATH_EVENT"] # Target variable
4
5 # Creating Feature array and Label array
6 X = df.drop(['DEATH_EVENT'], axis=1).values # Features
7 Y = df["DEATH_EVENT"].values.reshape(-1, 1) # Target variable, reshaped to a 2D array
8
9 # Splitting the data into training and testing sets
10 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.3, random_state=42)
11
12 # Define the hyperparameter grid for polynomial degree
13 grid_param = {'poly_degree': np.arange(1, 10)}
14
15 # Create a Pipeline for polynomial regression
16 PIPELINE = Pipeline([
17     ('poly', PolynomialFeatures(include_bias=False)),
18     ('regression', LinearRegression())
19 ])
20
21 # Perform GridSearchCV with cross-validation
22 gridsearch = GridSearchCV(PIPELINE, grid_param, cv=10, scoring="neg_mean_squared_error", verbose=0)
23
24 # Fit the grid search to the training data
25 gridsearch.fit(X_train, y_train)
26
27 # Make predictions on the test data
28 pred = gridsearch.predict(X_test)
29
30 # Convert predictions to binary Labels (0 or 1) based on a threshold of 0.5
31 pred = pred > 0.5
32 pred = pred.astype(int)
33
34 # Calculate evaluation metrics
35 poly_conf = confusion_matrix(y_test, pred)
36 poly_acc = metrics.accuracy_score(y_test, pred)
37 poly_recall = metrics.recall_score(y_test, pred)
38 poly_f1 = metrics.f1_score(y_test, pred)
39
40 # Print confusion matrix and evaluation metrics
41 print('Confusion Matrix \n {}'.format(poly_conf))
42 print('Accuracy : {:.2f}\nRecall : {:.2f}\nf1 score : {:.2f}'.format(poly_acc, poly_recall, poly_f1))
43

```

Confusion Matrix

```

[[50  3]
 [17 20]]

```

Accuracy : 0.78

Recall : 0.54

f1 score : 0.67

Lasso Regression

In [254]:

```
1 from sklearn.linear_model import Lasso
2 from sklearn.metrics import mean_squared_error
3
4 # Define a hyperparameter grid for Lasso alpha
5 grid_param = {'lasso_alpha': np.arange(0.01, 1, 0.01)}
6
7 # Create a Pipeline for Lasso regression
8 PIPELINE = Pipeline([('lasso', Lasso())])
9
10 # Perform GridSearchCV with cross-validation to find the optimal Lasso alpha
11 gridsearch = GridSearchCV(PIPELINE, grid_param, cv=10, scoring="neg_mean_squared_error", verbose=1)
12
13 # Fit the grid search to the training data
14 gridsearch.fit(X_train, y_train)
15
16 # Make predictions on the test data using the best Lasso model from grid search
17 pred = gridsearch.predict(X_test)
18
19 # Convert predictions to binary Labels (0 or 1) based on a threshold of 0.5
20 pred = pred > 0.5
21 pred = pred.astype(int)
22
23 # Calculate the confusion matrix and evaluation metrics for Lasso regression
24 lasso_conf = confusion_matrix(y_test, pred)
25 lasso_acc = metrics.accuracy_score(y_test, pred)
26 lasso_recall = metrics.recall_score(y_test, pred)
27 lasso_f1 = metrics.f1_score(y_test, pred)
28
29 # Print the confusion matrix and evaluation metrics for Lasso regression
30 print('Confusion Matrix \n{}\n'.format(lasso_conf))
31 print('Accuracy : {:.2f}\nRecall : {:.2f}\nf1 score : {:.2f}'.format(lasso_acc, lasso_recall, lasso_f1))
```

Confusion Matrix

```
[[50  3]
 [16 21]]
```

Accuracy : 0.79

Recall : 0.57

f1 score : 0.69

Ridge Regression

```
In [255]: 1 from sklearn.linear_model import Ridge
2 from sklearn.metrics import mean_squared_error
3
4 # Define a hyperparameter grid for Ridge alpha
5 grid_param = {'ridge_alpha': np.arange(0.01, 1, 0.01)}
6
7 # Create a Pipeline for Ridge regression
8 PIPELINE = Pipeline([('ridge', Ridge())])
9
10 # Perform GridSearchCV with cross-validation to find the optimal Ridge alpha
11 gridsearch = GridSearchCV(PIPELINE, grid_param, cv=10, scoring="neg_mean_squared_error", verbose=1)
12
13 # Fit the grid search to the training data
14 gridsearch.fit(X_train, y_train)
15
16 # Make predictions on the test data using the best Ridge model from grid search
17 pred = gridsearch.predict(X_test)
18
19 # Convert predictions to binary Labels (0 or 1) based on a threshold of 0.5
20 pred = pred > 0.5
21 pred = pred.astype(int)
22
23 # Calculate the confusion matrix and evaluation metrics for Ridge regression
24 ridge_conf = confusion_matrix(y_test, pred)
25 ridge_acc = metrics.accuracy_score(y_test, pred)
26 ridge_recall = metrics.recall_score(y_test, pred)
27 ridge_f1 = metrics.f1_score(y_test, pred)
28
29 # Print the confusion matrix and evaluation metrics for Ridge regression
30 print('Confusion Matrix \n {} \n'.format(ridge_conf))
31 print('Accuracy : {:.2f}\nRecall : {:.2f}\nf1 score : {:.2f}'.format(ridge_acc, ridge_recall, ridge_f1))
```

Confusion Matrix
[[50 3]
[17 20]]

Accuracy : 0.78
Recall : 0.54
f1 score : 0.67

Feature Selection

```
In [202]: 1 from sklearn.linear_model import LassoCV
2 from sklearn.feature_selection import SelectFromModel
3
4 lasso_model = LassoCV(cv=5)
5 lasso_model.fit(X_train, y_train_1d)
6 selected_features = SelectFromModel(lasso_model, prefit=True).transform(X_train)
7 selected_test_features = SelectFromModel(lasso_model, prefit=True).transform(X_test)
8
```

```
In [237]: 1 important_feature_names = df_X.columns[lasso_model.coef_ != 0].tolist()
2 print(important_feature_names)
```

['age', 'creatinine_phosphokinase', 'ejection_fraction', 'platelets', 'time']

K Nearest Neighbour

```
In [51]: 1 from sklearn.neighbors import KNeighborsClassifier
```

```
In [256]: 1 from sklearn.preprocessing import StandardScaler
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.pipeline import Pipeline
4
5 # Standardize the selected features for both training and test data
6 scaler = StandardScaler()
7 selected_features_norm = scaler.fit_transform(selected_features)
8 selected_test_features_norm = scaler.transform(selected_test_features)
9
10 # Define a hyperparameter grid for the number of neighbors in KNeighborsClassifier
11 param_grid = {'knn_n_neighbors': np.arange(1, 20)}
12
13 # Create a Pipeline for KNeighborsClassifier
14 est = Pipeline([('knn', KNeighborsClassifier())])
15
16 # Perform GridSearchCV with cross-validation to find the optimal number of neighbors
17 gridsearch = GridSearchCV(estimator=est, param_grid=param_grid, cv=10, scoring="recall_micro")
18
19 # Fit the grid search to the standardized training data
20 model = gridsearch.fit(selected_features_norm, y_train_1d)
21
22 # Print the best hyperparameters and corresponding cross-validation score
23 print(model.best_params_, model.best_score_, '\n')
24
25 # Make predictions on the standardized test data using the best model
26 pred = model.predict(selected_test_features_norm)
27
28 # Calculate the confusion matrix and evaluation metrics for KNeighborsClassifier
29 knn_conf = confusion_matrix(y_test, pred)
30 knn_acc = metrics.accuracy_score(y_test, pred)
31 knn_recall = metrics.recall_score(y_test, pred)
32 knn_f1 = metrics.f1_score(y_test, pred)
33
34 # Print the confusion matrix and evaluation metrics for KNeighborsClassifier
35 print('Confusion Matrix \n {} \n'.format(knn_conf))
36 print('Accuracy : {:.2f}\nRecall : {:.2f}\nf1 score : {:.2f}'.format(knn_acc, knn_recall, knn_f1))
37
```

```
{'knn_n_neighbors': 3} 0.8421428571428571
```

```
Confusion Matrix
```

```
[[48  5]
 [16 21]]
```

```
Accuracy : 0.77
Recall : 0.57
f1 score : 0.67
```

Decision Tree

```
In [78]: 1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn import metrics
3 import sklearn.tree as tree
```

```

In [257]: 1 # Splitting the data into training and testing sets using a 70-30 split and a fixed random seed
2 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.3, random_state=42)
3
4 # Define hyperparameter grid for DecisionTreeClassifier
5 param = {'tree_max_depth': np.arange(1, 20),
6           'tree_criterion': ["entropy", "gini"]}
7
8 # Create a Pipeline for DecisionTreeClassifier
9 pipe = Pipeline([('tree', DecisionTreeClassifier())])
10
11 # Perform GridSearchCV with cross-validation to find the optimal hyperparameters
12 gridsearch = GridSearchCV(estimator=pipe, param_grid=param, cv=10, scoring="accuracy", n_jobs=-1)
13
14 # Fit the grid search to the selected features and training labels
15 gridsearch.fit(selected_features, y_train)
16
17 # Print the best cross-validation score and corresponding hyperparameters
18 print(gridsearch.best_score_, gridsearch.best_params_, '\n')
19
20 # Get the best estimator from grid search
21 model = gridsearch.best_estimator_
22
23 # Make predictions on the selected test features
24 pred = model.predict(selected_test_features)
25
26 # Calculate the confusion matrix for DecisionTreeClassifier
27 dec_conf = confusion_matrix(y_test, pred)
28 dec_acc = metrics.accuracy_score(y_test, pred)
29 dec_recall = metrics.recall_score(y_test, pred)
30 dec_f1 = metrics.f1_score(y_test, pred)
31
32 # Print the confusion matrix and evaluation metrics for DecisionTreeClassifier
33 print('Confusion Matrix \n {} \n'.format(dec_conf))
34 print('Accuracy : {:.2f}\nRecall : {:.2f}\nf1 score : {:.2f}'.format(dec_acc, dec_recall, dec_f1))
35

```

```
0.865952380952381 {'tree_criterion': 'entropy', 'tree_max_depth': 1}
```

```
Confusion Matrix
[[50  3]
 [18 19]]
```

```
Accuracy : 0.77
Recall : 0.51
f1 score : 0.64
```

Logistic Regression

```

In [103]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import confusion_matrix

```

In [258]:

```
1 # Define hyperparameter grid for LogisticRegression
2 param_grid = {'lr_C': np.arange(0.01, 4, 0.5),
3                 'lr_solver': ['liblinear']}
4
5 # Create a Pipeline for LogisticRegression
6 est = Pipeline([('lr', LogisticRegression())])
7
8 # Perform GridSearchCV with cross-validation to find the optimal hyperparameters
9 gridsearch = GridSearchCV(estimator=est, param_grid=param_grid, cv=10, scoring="accuracy", n_jobs=-1)
10
11 # Fit the grid search to the selected features and training labels
12 gridsearch.fit(selected_features, y_train_1d)
13
14 # Make predictions on the selected test features
15 pred = gridsearch.predict(selected_test_features)
16
17 # Calculate the confusion matrix for LogisticRegression
18 lr_conf = confusion_matrix(y_test, pred)
19 lr_acc = metrics.accuracy_score(y_test, pred)
20 lr_recall = metrics.recall_score(y_test, pred)
21 lr_f1 = metrics.f1_score(y_test, pred)
22
23 # Print the best hyperparameters from grid search
24 print(gridsearch.best_params_, '\n')
25
26 # Print the confusion matrix and evaluation metrics for LogisticRegression
27 print('Confusion Matrix \n {} \n'.format(lr_conf))
28 print('Accuracy : {:.2f}\nRecall : {:.2f}\nf1 score : {:.2f}'.format(lr_acc, lr_recall, lr_f1))
```

```
{'lr_C': 0.01, 'lr_solver': 'liblinear'}
```

```
Confusion Matrix
```

```
[[50  3]
 [16 21]]
```

```
Accuracy : 0.79
```

```
Recall : 0.57
```

```
f1 score : 0.69
```

SVM

```
In [261]: 1 from sklearn.preprocessing import StandardScaler
2 from sklearn.svm import SVC
3 from sklearn.preprocessing import PolynomialFeatures
4
5 # Standardize the selected features for both training and test data
6 selected_features_norm = stdscaler.fit_transform(selected_features)
7 selected_test_features_norm = stdscaler.transform(selected_test_features)
8
9 # Define a hyperparameter grid for the SVM classifier
10 param_grid = {'svm_degree': np.arange(1, 10, 1),
11             'svm_C': np.arange(0.01, 10, 0.5),
12             'svm_gamma': np.linspace(0.01, 5, 10)}
13
14 # Create a Pipeline for Polynomial SVM classifier
15 pipe = Pipeline([('polyfeature', PolynomialFeatures()), # Add polynomial features
16                  ('svm', SVC(kernel='rbf'))]) # Use the radial basis function kernel
17
18 # Perform GridSearchCV with cross-validation to find the optimal hyperparameters
19 gridsearch = GridSearchCV(estimator=pipe, param_grid=param_grid, cv=10, scoring='f1', verbose=1)
20
21 # Fit the grid search to the standardized training features and labels
22 gridsearch.fit(selected_features_norm, y_train)
23
24 # Print the best hyperparameters from grid search
25 print(gridsearch.best_params_, '\n')
26
27 svm_conf = confusion_matrix(y_test, pred)
28 svm_acc = metrics.accuracy_score(y_test, pred)
29 svm_recall = metrics.recall_score(y_test, pred)
30 svm_f1 = metrics.f1_score(y_test, pred)
31
32
33 print('Confusion Matrix \n {} \n'.format(svm_conf))
34 print('Accuracy : {:.2f}\nRecall : {:.2f}\nf1 score : {:.2f}'.format(svm_acc, svm_recall, svm_f1))
35
```

{'svm_C': 4.01, 'svm_degree': 1, 'svm_gamma': 0.01}

Confusion Matrix

50	3
16	21

Accuracy : 0.79
 Recall : 0.57
 f1 score : 0.69

E:\Apps\Anaconda\Lib\site-packages\sklearn\utils\validation.py:1143: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
 y = column_or_1d(y, warn=True)

Random Forest

```
In [118]: 1 from sklearn.ensemble import RandomForestClassifier
```

```

In [262]: 1 # Define a hyperparameter grid for the number of estimators in RandomForestClassifier
2 param_grid = {'rf__n_estimators': np.arange(1, 100, 2)}
3
4 # Create a RandomForestClassifier with a fixed random state for reproducibility
5 rf_classifier = RandomForestClassifier(random_state=42)
6
7 # Create a Pipeline for RandomForestClassifier
8 pipe = Pipeline([('rf', rf_classifier)])
9
10 # Perform GridSearchCV with cross-validation to find the optimal number of estimators
11 gridsearch = GridSearchCV(estimator=pipe, param_grid=param_grid, cv=10, scoring='accuracy', ve
12
13 # Fit the grid search to the selected features and training labels
14 gridsearch.fit(selected_features, y_train)
15
16 # Print the best hyperparameters from grid search
17 print(gridsearch.best_params_, '\n')
18
19 # Make predictions on the selected test features using the best model
20 pred = gridsearch.predict(selected_test_features)
21
22 # Calculate the confusion matrix and evaluation metrics for RandomForestClassifier
23 rf_conf = confusion_matrix(y_test, pred)
24 rf_acc = metrics.accuracy_score(y_test, pred)
25 rf_recall = metrics.recall_score(y_test, pred)
26 rf_f1 = metrics.f1_score(y_test, pred)
27
28 # Print the confusion matrix and evaluation metrics for RandomForestClassifier
29 print('Confusion Matrix \n {} \n'.format(rf_conf))
30 print('Accuracy : {:.2f}\nRecall : {:.2f}\nf1 score : {:.2f}'.format(rf_acc, rf_recall, rf_f1)
31

```

```
{'rf__n_estimators': 89}
```

```
Confusion Matrix
[[48  5]
 [18 19]]
```

```
Accuracy : 0.74
Recall : 0.51
f1 score : 0.62
```

```
E:\Apps\Anaconda\Lib\site-packages\sklearn\pipeline.py:405: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
  self._final_estimator.fit(Xt, y, **fit_params_last_step)
```

```

In [264]: 1 # Define a hyperparameter grid for the number of estimators in RandomForestClassifier
2 param_grid = {'rf__n_estimators': np.arange(1, 100, 2)}
3
4 # Create a RandomForestClassifier with a fixed random state for reproducibility
5 rf_classifier = RandomForestClassifier(random_state=42)
6
7 # Create a Pipeline for RandomForestClassifier
8 pipe = Pipeline([('rf', rf_classifier)])
9
10 # Perform GridSearchCV with cross-validation to find the optimal number of estimators
11 gridsearch = GridSearchCV(estimator=pipe, param_grid=param_grid, cv=10, scoring='accuracy', ve
12
13 y_train_1d = y_train.ravel()
14 # Fit the grid search to the training features and Labels
15 gridsearch.fit(X_train, y_train_1d)
16
17 # Print the best hyperparameters from grid search
18 print(gridsearch.best_params_, '\n')
19
20 # Make predictions on the test features using the best model
21 pred = gridsearch.predict(X_test)
22
23 # Calculate the confusion matrix and evaluation metrics for RandomForestClassifier
24 rf2_conf = confusion_matrix(y_test, pred)
25 rf2_acc = metrics.accuracy_score(y_test, pred)
26 rf2_recall = metrics.recall_score(y_test, pred)
27 rf2_f1 = metrics.f1_score(y_test, pred)
28
29 # Print the confusion matrix and evaluation metrics for RandomForestClassifier
30 print('Confusion Matrix \n {} \n'.format(rf2_conf))
31 print('Accuracy : {:.2f}\nRecall : {:.2f}\nf1 score : {:.2f}'.format(rf2_acc, rf2_recall, rf2_
32

```

```
{'rf__n_estimators': 63}
```

```

Confusion Matrix
[[49  4]
 [17 20]]

```

```

Accuracy : 0.77
Recall : 0.54
f1 score : 0.66

```

VotingClassifier

```

In [142]: 1 from sklearn.ensemble import VotingClassifier

```

In [269]:

```

1 # Create individual classifiers with specified hyperparameters
2 log_clf = LogisticRegression(C=0.01)
3 rnd_clf = RandomForestClassifier(n_estimators=89)
4 tree_clf = DecisionTreeClassifier(criterion='entropy', max_depth=1)
5 svm_clf = SVC(C=4.01, degree=1, gamma=0.01, probability=True)
6 knn_clf = KNeighborsClassifier(n_neighbors=3)
7
8 # Create a VotingClassifier that combines the predictions of individual classifiers using hard
9 voting_clf = VotingClassifier(estimators=[('log', log_clf), ('rfc', rnd_clf), ('tree', tree_clf)])
10
11 # Fit the VotingClassifier to the standardized training features and labels
12 voting_clf.fit(selected_features_norm, y_train_1d)
13
14 # Make predictions on the standardized test features
15 pred = voting_clf.predict(selected_test_features_norm)
16
17 # Calculate the confusion matrix and evaluation metrics for the VotingClassifier
18 vot_conf = confusion_matrix(y_test, pred)
19 vot_acc = metrics.accuracy_score(y_test, pred)
20 vot_recall = metrics.recall_score(y_test, pred)
21 vot_f1 = metrics.f1_score(y_test, pred)
22
23 # Print the confusion matrix and evaluation metrics for the VotingClassifier
24 print('Confusion Matrix \n{} \n'.format(vot_conf))
25 print('Accuracy : {:.2f}\nRecall : {:.2f}\nf1 score : {:.2f}'.format(vot_acc, vot_recall, vot_f1))
26

```

Confusion Matrix

```

[[50  3]
 [20 17]]

```

Accuracy : 0.74

Recall : 0.46

f1 score : 0.60

Neural Network

In [149]:

```

1 from sklearn.preprocessing import StandardScaler
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.layers import Dense
4 from tensorflow.keras.utils import plot_model
5 import tensorflow as tf

```

In [334]:

```

1 # Define a Sequential neural network model
2 model = Sequential([
3     Dense(4, activation='relu', input_shape=(selected_features_norm.shape[1], )),
4     Dense(2, activation='relu'),
5     Dense(1, activation='sigmoid')
6 ])
7
8 # Compile the model with optimizer, loss function, and metrics
9 model.compile(optimizer='adam', loss=tf.keras.losses.binary_crossentropy, metrics=['accuracy'])
10
11 # Train the model on your training data
12 h = model.fit(selected_features_norm, y_train, epochs=200, batch_size=8, validation_split=0.2,
13
14 # Evaluate the model on your test data
15 model.evaluate(selected_test_features_norm, y_test)
16
17 # Make predictions using the model and threshold them to binary values (0 or 1)
18 pred = model.predict(selected_test_features_norm)
19 pred = (pred > 0.5).astype(int)
20
21 # Calculate the confusion matrix, accuracy, recall, and F1 score
22 nn_conf = confusion_matrix(y_test, pred)
23 nn_acc = metrics.accuracy_score(y_test, pred)
24 nn_recall = metrics.recall_score(y_test, pred)
25 nn_f1 = metrics.f1_score(y_test, pred)
26
27 # Print the confusion matrix and evaluation metrics
28 print('Confusion Matrix \n{}\n'.format(nn_conf))
29 print('Accuracy : {:.2f}\nRecall : {:.2f}\nf1 score : {:.2f}'.format(nn_acc, nn_recall, nn_f1))

```

3/3 [=====] - 0s 2ms/step - loss: 0.5154 - accuracy: 0.8222

3/3 [=====] - 0s 2ms/step

Confusion Matrix

```

[[50  3]
 [13 24]]

```

Accuracy : 0.82

Recall : 0.65

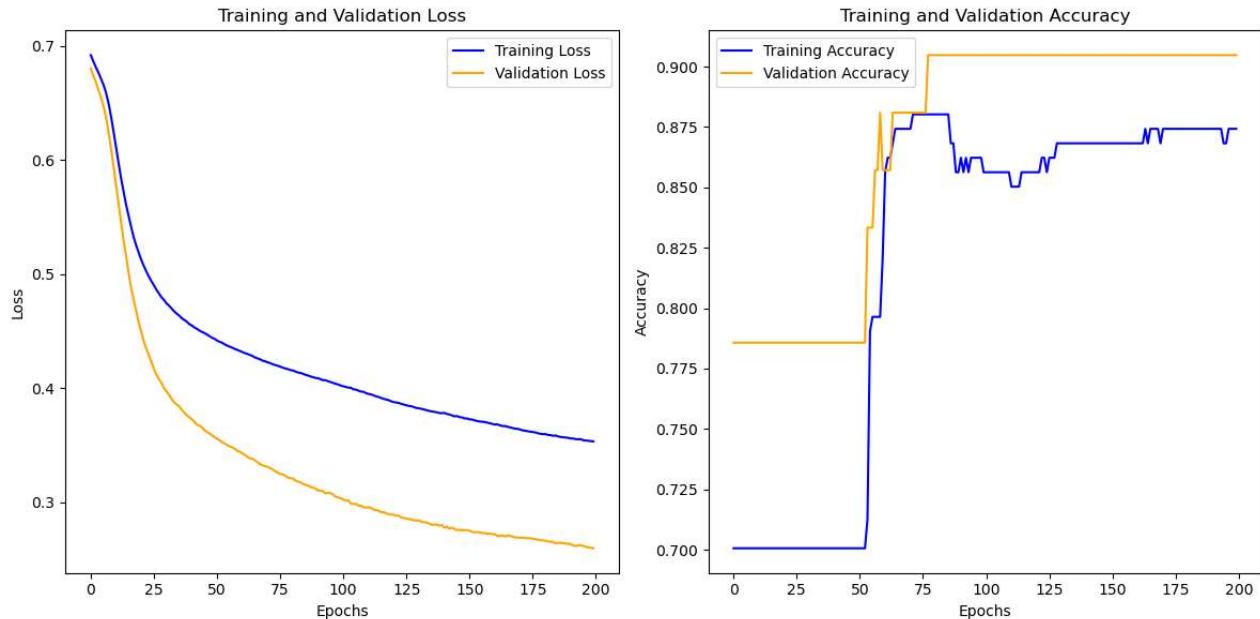
f1 score : 0.75

In [335]:

```

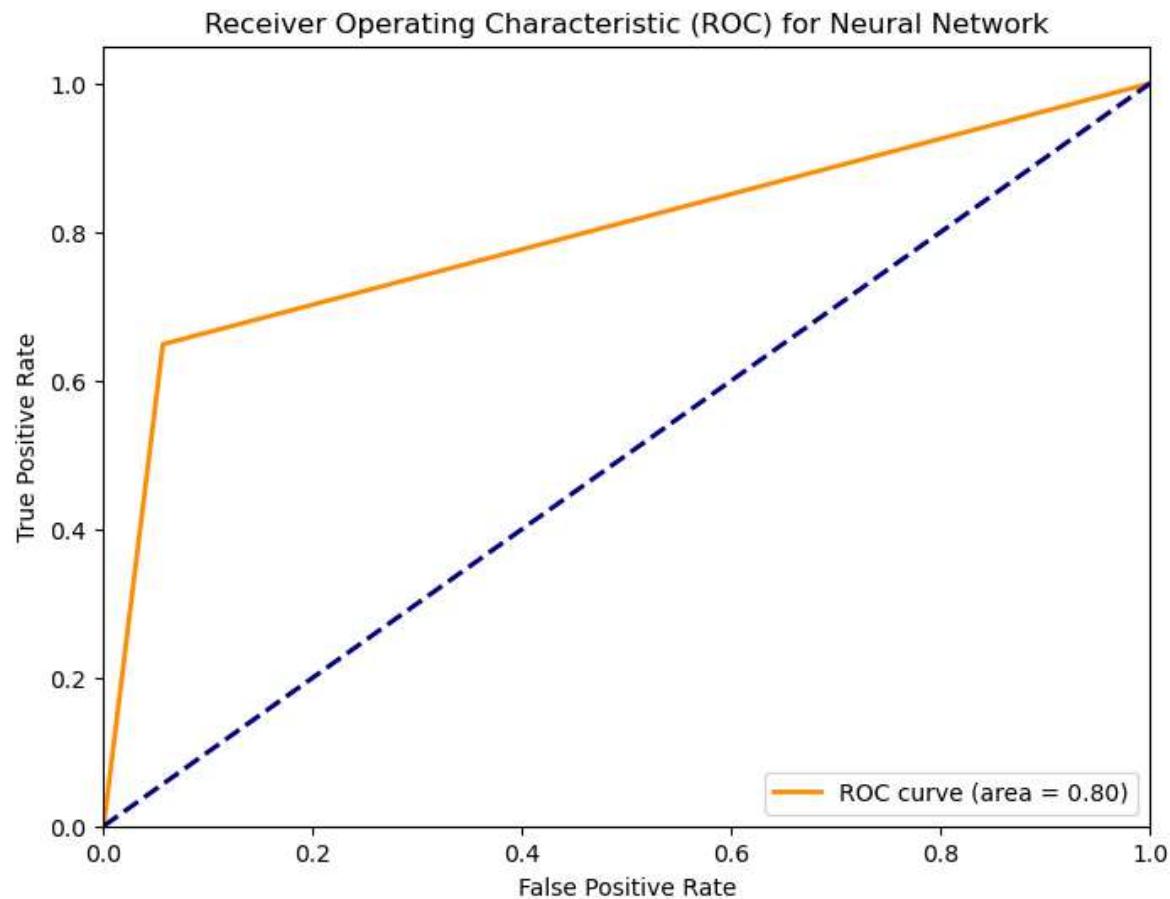
1 # Access training history
2 training_loss = h.history['loss']
3 validation_loss = h.history['val_loss']
4 training_accuracy = h.history['accuracy']
5 validation_accuracy = h.history['val_accuracy']
6
7 # Create subplots for loss and accuracy
8 plt.figure(figsize=(12, 6))
9
10 # Loss subplot
11 plt.subplot(1, 2, 1)
12 plt.plot(training_loss, label='Training Loss', color='blue')
13 plt.plot(validation_loss, label='Validation Loss', color='orange')
14 plt.title('Training and Validation Loss')
15 plt.xlabel('Epochs')
16 plt.ylabel('Loss')
17 plt.legend()
18
19 # Accuracy subplot
20 plt.subplot(1, 2, 2)
21 plt.plot(training_accuracy, label='Training Accuracy', color='blue')
22 plt.plot(validation_accuracy, label='Validation Accuracy', color='orange')
23 plt.title('Training and Validation Accuracy')
24 plt.xlabel('Epochs')
25 plt.ylabel('Accuracy')
26 plt.legend()
27
28 plt.tight_layout()
29 plt.show()
30

```



In [339]:

```
1 fpr, tpr, thresholds = roc_curve(y_test, pred)
2 roc_auc=auc(fpr,tpr)
3
4 plt.figure(figsize=(8, 6))
5 plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = {:.2f})'.format(roc_auc))
6 plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
7 plt.xlim([0.0, 1.0])
8 plt.ylim([0.0, 1.05])
9 plt.xlabel('False Positive Rate')
10 plt.ylabel('True Positive Rate')
11 plt.title('Receiver Operating Characteristic (ROC) for Neural Network')
12 plt.legend(loc='lower right')
13 plt.show()
```



AdaBoost

In [329]:

```

1  from sklearn.ensemble import AdaBoostClassifier
2  from sklearn.linear_model import LogisticRegression
3  from sklearn.model_selection import GridSearchCV
4  from sklearn.datasets import load_iris
5
6  # Load a dataset (e.g., Iris dataset)
7  data = load_iris()
8  X = data.data
9  y = data.target
10
11 # Create a LogisticRegression estimator
12 base_estimator = LogisticRegression(C=0.01)
13
14 # Create an AdaBoostClassifier with the Logistic regression estimator
15 ada_classifier = AdaBoostClassifier(base_estimator=base_estimator, algorithm="SAMME.R")
16
17 # Define the hyperparameters and their possible values
18 param_grid = {
19     'n_estimators': np.arange(1, 50, 1), # Number of boosting rounds
20     'learning_rate': np.arange(0.01, 0.1, 0.01) # Learning rate for each round
21 }
22
23 # Create a GridSearchCV instance with cross-validation (e.g., 5-fold)
24 grid_search = GridSearchCV(ada_classifier, param_grid, cv=20, scoring='accuracy', verbose=2, n_jobs=-1)
25
26 # Fit the GridSearchCV to the data
27 grid_search.fit(selected_features, y_train)
28
29 pred = ada_classifier.predict(selected_test_features)
30
31 print(grid_search.best_estimator_, '\n')
32
33 ada_conf = confusion_matrix(y_test, pred)
34 ada_acc = metrics.accuracy_score(y_test, pred)
35 ada_recall = metrics.recall_score(y_test, pred)
36 ada_f1 = metrics.f1_score(y_test, pred)
37
38 print('Confusion Matrix \n {} \n'.format(ada_conf))
39 print('Accuracy : {:.3f}\nRecall : {:.2f}\nf1 score : {:.2f}'.format(ada_acc, ada_recall, ada_f1))
40

```

Fitting 20 folds for each of 441 candidates, totalling 8820 fits
 AdaBoostClassifier(base_estimator=LogisticRegression(C=0.01), learning_rate=0.01,
 n_estimators=1)

Confusion Matrix
 [[49 4]
 [15 22]]

Accuracy : 0.789
 Recall : 0.59
 f1 score : 0.70

E:\Apps\Anaconda\Lib\site-packages\sklearn\utils\validation.py:1143: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
 y = column_or_1d(y, warn=True)

Gradient Boosting

In [197]:

```

1 #For Most Important Features
2
3 from sklearn.ensemble import GradientBoostingClassifier
4 from sklearn.model_selection import GridSearchCV
5 from sklearn.datasets import load_iris
6 from sklearn.metrics import confusion_matrix, accuracy_score, recall_score, f1_score
7 import numpy as np
8
9 # Load a dataset (e.g., Iris dataset)
10 data = load_iris()
11 X = data.data
12 y = data.target
13
14 # Define the hyperparameters and their possible values
15 param_grid = {
16     'n_estimators': np.arange(1, 200, 1),           # Number of boosting rounds
17     'learning_rate': np.arange(0.01, 0.1, 0.01)    # Learning rate for each round
18 }
19
20 # Create a GridSearchCV instance with cross-validation (e.g., 5-fold)
21 grid_search = GridSearchCV(GradientBoostingClassifier(), param_grid, cv=10, scoring='accuracy')
22
23 # Fit the GridSearchCV to the data
24 grid_search.fit(selected_features, y_train)
25
26 # Get the best estimator (GradientBoostingClassifier with the best hyperparameters)
27 best_gbc = grid_search.best_estimator_
28
29 # Make predictions on the test data
30 pred = best_gbc.predict(selected_test_features)
31
32 print(grid_search.best_estimator_, '\n')
33
34 # Calculate evaluation metrics
35 gbc_conf = confusion_matrix(y_test, pred)
36 gbc_acc = accuracy_score(y_test, pred)
37 gbc_recall = recall_score(y_test, pred)
38 gbc_f1 = f1_score(y_test, pred)
39
40 print('Confusion Matrix \n {} \n'.format(gbc_conf))
41 print('Accuracy : {:.3f}\nRecall : {:.3f}\nf1 score : {:.3f}'.format(gbc_acc, gbc_recall, gbc_f1))
42

```

Fitting 10 folds for each of 1791 candidates, totalling 17910 fits
 GradientBoostingClassifier(learning_rate=0.09, n_estimators=15)

Confusion Matrix
 [[49 4]
 [18 19]]

Accuracy : 0.756
 Recall : 0.514
 f1 score : 0.633

In [312]:

```

1 # For All Features
2
3 from sklearn.ensemble import GradientBoostingClassifier
4 from sklearn.model_selection import GridSearchCV
5 from sklearn.datasets import load_iris
6 from sklearn.metrics import confusion_matrix, accuracy_score, recall_score, f1_score
7 import numpy as np
8
9 # Load a dataset (e.g., Iris dataset)
10 data = load_iris()
11 X = data.data
12 y = data.target
13
14 # Define the hyperparameters and their possible values
15 param_grid = {
16     'n_estimators': np.arange(1, 200, 1),           # Number of boosting rounds
17     'learning_rate': np.arange(0.01, 0.1, 0.01)    # Learning rate for each round
18 }
19
20 # Create a GridSearchCV instance with cross-validation (e.g., 5-fold)
21 grid_search = GridSearchCV(GradientBoostingClassifier(), param_grid, cv=10, scoring='accuracy')
22
23 # Fit the GridSearchCV to the data
24 grid_search.fit(X_train, y_train)
25
26 # Get the best estimator (GradientBoostingClassifier with the best hyperparameters)
27 best_gbc = grid_search.best_estimator_
28
29 # Make predictions on the test data
30 pred = best_gbc.predict(X_test)
31
32 print(grid_search.best_estimator_, '\n')
33
34 # Calculate evaluation metrics
35 gbc2_conf = confusion_matrix(y_test, pred)
36 gbc2_acc = accuracy_score(y_test, pred)
37 gbc2_recall = recall_score(y_test, pred)
38 gbc2_f1 = f1_score(y_test, pred)
39
40 print('Confusion Matrix \n{} \n'.format(gbc2_conf))
41 print('Accuracy : {:.3f}\nRecall : {:.3f}\nf1 score : {:.3f}'.format(gbc2_acc, gbc2_recall, gbc2_f1))
42

```

Fitting 10 folds for each of 1791 candidates, totalling 17910 fits
 GradientBoostingClassifier(learning_rate=0.06000000000000005, n_estimators=87)

Confusion Matrix

```

[[47  6]
 [18 19]]

```

Accuracy : 0.733
 Recall : 0.514
 f1 score : 0.613

E:\Apps\Anaconda\Lib\site-packages\sklearn\ensemble_gb.py:437: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
 y = column_or_1d(y, warn=True)

```
In [338]: 1 var=['poly', 'knn', 'dec', 'svm', 'svm2', 'lr', 'vot', 'nn', 'ada', 'gbc', 'gbc2']  
2 param=['conf', 'acc', 'recall', 'f1']  
3 for j in range(4):  
4     for i in range(9):  
5         variable = f'{var[i]}_{param[j]}'  
6         if variable in locals():  
7             value=locals()[variable]  
8  
9         print(variable, "\n", value)  
10    print("\n")  
11
```

```
poly_conf
[[50 3]
 [17 20]]
knn_conf
[[48 5]
 [16 21]]
dec_conf
[[50 3]
 [18 19]]
svm_conf
[[50 3]
 [16 21]]
lr_conf
[[50 3]
 [16 21]]
vot_conf
[[50 3]
 [20 17]]
nn_conf
[[50 3]
 [13 24]]
ada_conf
[[49 4]
 [15 22]]

poly_acc
0.7777777777777778
knn_acc
0.7666666666666667
dec_acc
0.7666666666666667
svm_acc
0.7888888888888889
lr_acc
0.7888888888888889
vot_acc
0.7444444444444445
nn_acc
0.8222222222222222
ada_acc
0.7888888888888889

poly_recall
0.5405405405405406
knn_recall
0.5675675675675675
dec_recall
0.5135135135135135
svm_recall
0.5675675675675675
lr_recall
0.5675675675675675
vot_recall
0.4594594594594595
nn_recall
0.6486486486486487
ada_recall
0.5945945945945946

poly_f1
0.6666666666666666
knn_f1
0.6666666666666666
dec_f1
0.6440677966101694
```

```
svm_f1
0.6885245901639344
lr_f1
0.6885245901639344
vot_f1
0.5964912280701754
nn_f1
0.75
ada_f1
0.6984126984126985
```

Results

We trained 9 machine learning model to predict the heart attack based on 12 features. We Used L1 regularization for feature selection and trained the ML models on these features.

We found that Out of all the ML models:

- Neural Network with structure 4-->2-->1 had the highest accuracy, recall and f1 score of 0.82, 0.64 and 0.75 respectively with an auc of 0.8
- The second Best performing model is adaptive boost with accuracy, recall and f1 score of 0.78, 0.59 and 0.69 respectively

```
In [ ]: 1
```

```
In [ ]: 1
```