# 🏥 Heart Failure Clinical Records Analysis

**Objective**

Analyze survival of patients with heart failure using 13 clinical features.

Dataset size: 299 patients.

**Steps in Analysis**

1. Dataset loading & inspection
2. Data quality check (missing, duplicates)
3. Outlier detection & handling (IQR method)
4. Exploratory Data Analysis (EDA)
5. Grouped analysis (age bins, comorbidities)
6. Key presentation visuals
7. Feature engineering (risk score)
8. Baseline ML models (Logistic Regression & Random Forest)
9. ROC & feature importance
10. Save cleaned dataset + portfolio recommendations

```
In [19]:  # Imports and plotting settings (compatible with your versions)
          import pandas as pd
          import numpy as np
          import matplotlib.pyplot as plt
          import seaborn as sns
          import os
          from sklearn.model_selection import train_test_split
          from sklearn.preprocessing import StandardScaler
          from sklearn.linear_model import LogisticRegression
          from sklearn.ensemble import RandomForestClassifier
          from sklearn.metrics import roc_auc_score, classification_report, confusion_matr

          # plotting defaults for presentation
          plt.style.use('seaborn-v0_8')
          sns.set_palette("husl")
          plt.rcParams['figure.figsize'] = (10,6)
          plt.rcParams['savefig.dpi'] = 300

          # Print versions (for reproducibility)
          import sys
          print("Python:", sys.version.split()[0])
          print("pandas:", pd.__version__, "numpy:", np.__version__)
          import sklearn
          print("scikit-learn:", sklearn.__version__)
```

```
Python: 3.13.7
pandas: 2.3.2 numpy: 2.2.2
scikit-learn: 1.7.2
```

# 1 — Load dataset

Ensure `heart_failure_clinical_records_dataset.csv` is in the project folder.
Load into a DataFrame and inspect.

In [20]:
```python
csv_path = "heart_failure_clinical_records_dataset.csv"
assert os.path.exists(csv_path), f"File not found: {csv_path}"

df = pd.read_csv(csv_path)
print("Shape:", df.shape)

display(df.head())
display(df.info())
display(df.describe().T)
```

Shape: (299, 13)

| | age | anaemia | creatinine_phosphokinase | diabetes | ejection_fraction | high_blood_pressu |
|---|---|---|---|---|---|---|
| **0** | 75.0 | 0 | 582 | 0 | 20 | |
| **1** | 55.0 | 0 | 7861 | 0 | 38 | |
| **2** | 65.0 | 0 | 146 | 0 | 20 | |
| **3** | 50.0 | 1 | 111 | 0 | 20 | |
| **4** | 65.0 | 1 | 160 | 1 | 20 | |

```
<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
None
```

|  | count | mean | std | min | 25% | 50% |
|---|---|---|---|---|---|---|
| **age** | 299.0 | 60.833893 | 11.894809 | 40.0 | 51.0 | 60.0 |
| **anaemia** | 299.0 | 0.431438 | 0.496107 | 0.0 | 0.0 | 0.0 |
| **creatinine_phosphokinase** | 299.0 | 581.839465 | 970.287881 | 23.0 | 116.5 | 250.0 |
| **diabetes** | 299.0 | 0.418060 | 0.494067 | 0.0 | 0.0 | 0.0 |
| **ejection_fraction** | 299.0 | 38.083612 | 11.834841 | 14.0 | 30.0 | 38.0 |
| **high_blood_pressure** | 299.0 | 0.351171 | 0.478136 | 0.0 | 0.0 | 0.0 |
| **platelets** | 299.0 | 263358.029264 | 97804.236869 | 25100.0 | 212500.0 | 262000.0 |
| **serum_creatinine** | 299.0 | 1.393880 | 1.034510 | 0.5 | 0.9 | 1.1 |
| **serum_sodium** | 299.0 | 136.625418 | 4.412477 | 113.0 | 134.0 | 137.0 |
| **sex** | 299.0 | 0.648829 | 0.478136 | 0.0 | 0.0 | 1.0 |
| **smoking** | 299.0 | 0.321070 | 0.467670 | 0.0 | 0.0 | 0.0 |
| **time** | 299.0 | 130.260870 | 77.614208 | 4.0 | 73.0 | 115.0 |
| **DEATH_EVENT** | 299.0 | 0.321070 | 0.467670 | 0.0 | 0.0 | 0.0 |

```
In [21]:   # Missing values
           print("Missing values per column:\n", df.isnull().sum())

           # Duplicates
           dups = df.duplicated().sum()
           print("\nDuplicate rows:", dups)

           # Data types
           print("\nData types:\n", df.dtypes)
```

```
Missing values per column:
 age                           0
anaemia                       0
creatinine_phosphokinase      0
diabetes                      0
ejection_fraction             0
high_blood_pressure           0
platelets                     0
serum_creatinine              0
serum_sodium                  0
sex                           0
smoking                       0
time                          0
DEATH_EVENT                   0
dtype: int64

Duplicate rows: 0

Data types:
 age                           float64
anaemia                         int64
creatinine_phosphokinase        int64
diabetes                        int64
ejection_fraction               int64
high_blood_pressure             int64
platelets                     float64
serum_creatinine              float64
serum_sodium                    int64
sex                             int64
smoking                         int64
time                            int64
DEATH_EVENT                     int64
dtype: object
```

# Target variable overview (DEATH_EVENT)

Check survival rate and balance of the target classes.

In [22]:
```python
target = "DEATH_EVENT"
counts = df[target].value_counts()
print(counts)
print(f"Survival rate: {(1 - df[target].mean())*100:.2f}%")

plt.figure(figsize=(5,5))
colors = ["#2ecc71", "#e74c3c"]
plt.pie(counts.values, labels=["Survived","Died"], autopct="%1.1f%%", colors=col
plt.title("Outcome distribution")
plt.show()
```
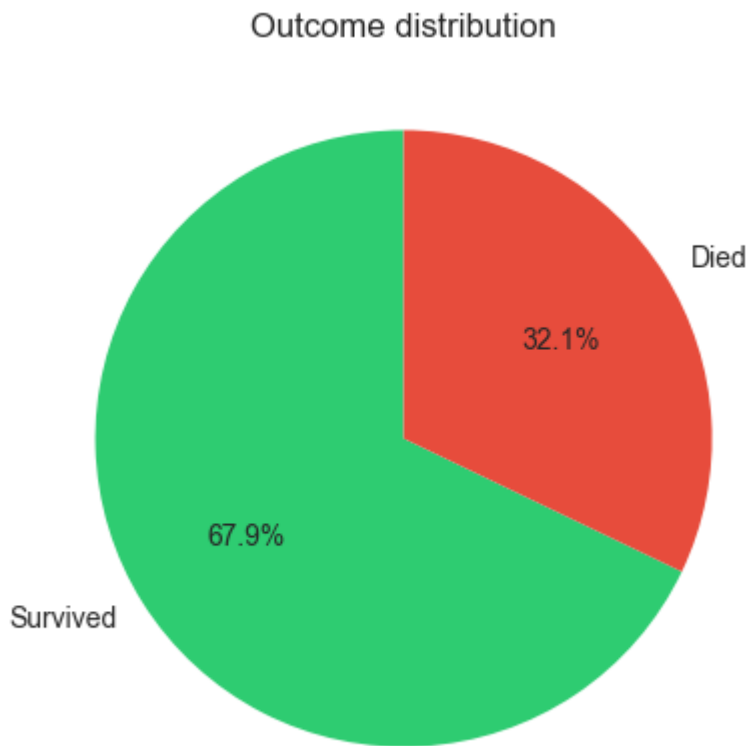
```
DEATH_EVENT
0    203
1     96
Name: count, dtype: int64
Survival rate: 67.89%
```

Outcome distribution



# 3 — Outlier detection (IQR method)

Detect outliers in continuous variables. Outliers will be capped (clipped).

```
In [23]: continuous_vars = ['age','creatinine_phosphokinase','ejection_fraction',
                            'platelets','serum_creatinine','serum_sodium','time']
         present_cont = [c for c in continuous_vars if c in df.columns]

         outlier_counts = {}
         for var in present_cont:
             Q1, Q3 = df[var].quantile([0.25,0.75])
             IQR = Q3 - Q1
             lb, ub = Q1 - 1.5*IQR, Q3 + 1.5*IQR
             outlier_counts[var] = ((df[var] < lb) | (df[var] > ub)).sum()

         outlier_counts
```

```
Out[23]: {'age': np.int64(0),
          'creatinine_phosphokinase': np.int64(29),
          'ejection_fraction': np.int64(2),
          'platelets': np.int64(21),
          'serum_creatinine': np.int64(29),
          'serum_sodium': np.int64(4),
          'time': np.int64(0)}
```

## 4 — Handle outliers (capping)

Clip values at lower/upper bounds. Create `df_clean` .

```
In [24]: df_clean = df.copy()
         for var in present_cont:
```

```
        Q1, Q3 = df[var].quantile([0.25,0.75])
        IQR = Q3 - Q1
        lb, ub = Q1 - 1.5*IQR, Q3 + 1.5*IQR
        df_clean[var] = df_clean[var].clip(lower=lb, upper=ub)

df_clean[present_cont].describe().T
```

Out[24]:

| | count | mean | std | min | 25% | 5( |
|---|---|---|---|---|---|---|
| age | 299.0 | 60.833893 | 11.894809 | 40.0 | 51.0 | 6 |
| creatinine_phosphokinase | 299.0 | 424.214883 | 385.449328 | 23.0 | 116.5 | 25 |
| ejection_fraction | 299.0 | 38.033445 | 11.685643 | 14.0 | 30.0 | 3 |
| platelets | 299.0 | 259163.714883 | 81478.304369 | 76000.0 | 212500.0 | 26200 |
| serum_creatinine | 299.0 | 1.234515 | 0.440098 | 0.5 | 0.9 | |
| serum_sodium | 299.0 | 136.712375 | 4.076971 | 125.0 | 134.0 | 13 |
| time | 299.0 | 130.260870 | 77.614208 | 4.0 | 73.0 | 11 |

# 5 — Correlation analysis

Check feature correlations with the target.

In [25]:
```python
plt.figure(figsize=(10,8))
corr = df_clean.corr()
mask = np.triu(np.ones_like(corr, dtype=bool))
sns.heatmap(corr, mask=mask, annot=True, fmt=".2f", cmap="coolwarm")
plt.title("Correlation heatmap")
plt.show()

print(corr[target].sort_values(ascending=False))
```

Correlation heatmap



```
DEATH_EVENT                     1.000000
serum_creatinine                0.388469
age                             0.253729
high_blood_pressure             0.079351
anaemia                         0.066270
diabetes                       -0.001943
sex                            -0.004316
creatinine_phosphokinase       -0.006355
smoking                        -0.012623
platelets                      -0.044132
serum_sodium                   -0.201320
ejection_fraction              -0.270611
time                           -0.526964
Name: DEATH_EVENT, dtype: float64
```

# 6 — Grouped analysis

Analyze death rates by age groups and comorbidities.

```
In [26]:   # Age groups
           bins = [0,50,60,70,80,120]
           labels = ["<50","50-59","60-69","70-79","80+"]
           df_clean["age_group"] = pd.cut(df_clean["age"], bins=bins, labels=labels, right=

           age_stats = df_clean.groupby("age_group")[target].mean()
           print(age_stats)
```

```python
sns.barplot(x=age_stats.index, y=age_stats.values)
plt.title("Death rate by age group")
plt.show()

# Comorbidities
for c in ["anaemia","diabetes","high_blood_pressure","smoking"]:
    if c in df_clean.columns:
        stats = df_clean.groupby(c)[target].mean()
        print(c, stats)
        sns.barplot(x=stats.index, y=stats.values)
        plt.title(f"Death rate by {c}")
        plt.show()
```
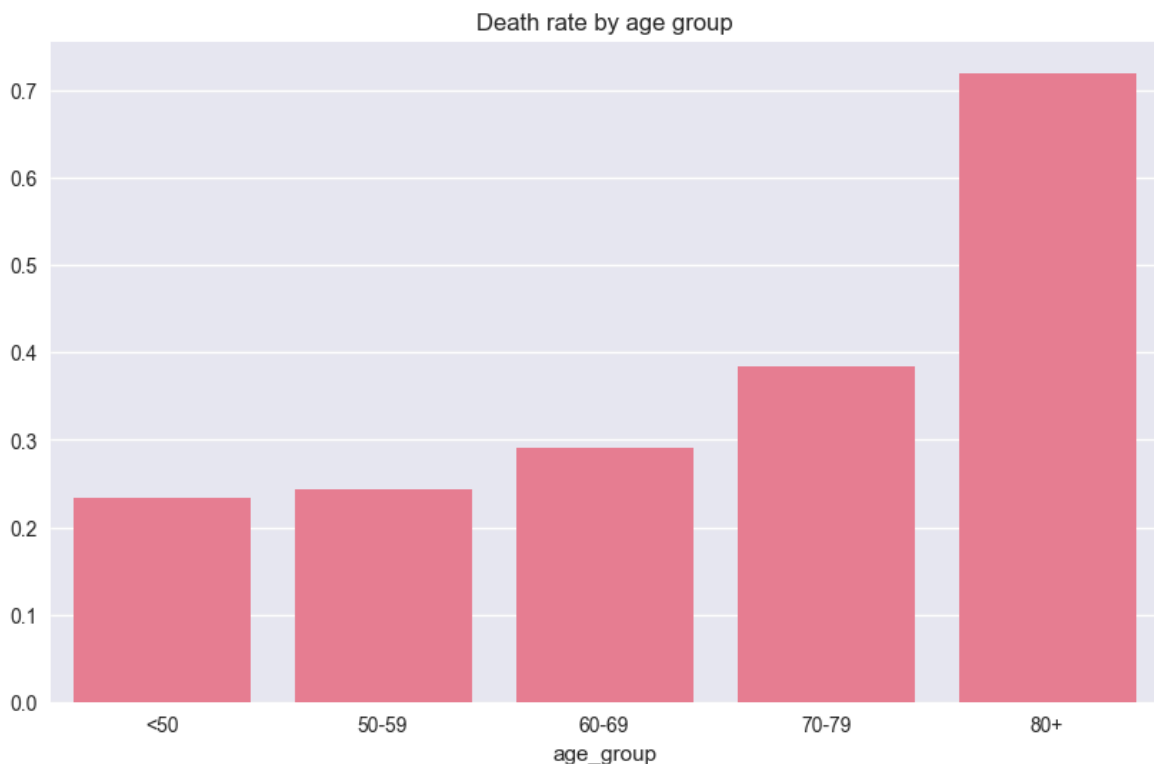
```
age_group
<50      0.234043
50-59    0.243902
60-69    0.290323
70-79    0.384615
80+      0.720000
Name: DEATH_EVENT, dtype: float64
```

```
C:\Users\nilsa\AppData\Local\Temp\ipykernel_19332\2809682682.py:6: FutureWarning:
The default of observed=False is deprecated and will be changed to True in a futu
re version of pandas. Pass observed=False to retain current behavior or observed=
True to adopt the future default and silence this warning.
  age_stats = df_clean.groupby("age_group")[target].mean()
```



Death rate by age group

```
anaemia anaemia
0     0.294118
1     0.356589
Name: DEATH_EVENT, dtype: float64
```

## Death rate by anaemia



```
diabetes diabetes
0    0.321839
1    0.320000
Name: DEATH_EVENT, dtype: float64
```

## Death rate by diabetes



```
high_blood_pressure high_blood_pressure
0    0.293814
1    0.371429
Name: DEATH_EVENT, dtype: float64
```

Death rate by high_blood_pressure



```
smoking smoking
0    0.325123
1    0.312500
Name: DEATH_EVENT, dtype: float64
```

Death rate by smoking



# 7 — Key visuals

Histograms, boxplots, and scatterplots for presentation.

```
In [27]:  # Histograms
          vars_plot = ['age','ejection_fraction','serum_creatinine','time','serum_sodium',
          plt.figure(figsize=(14,8))
```
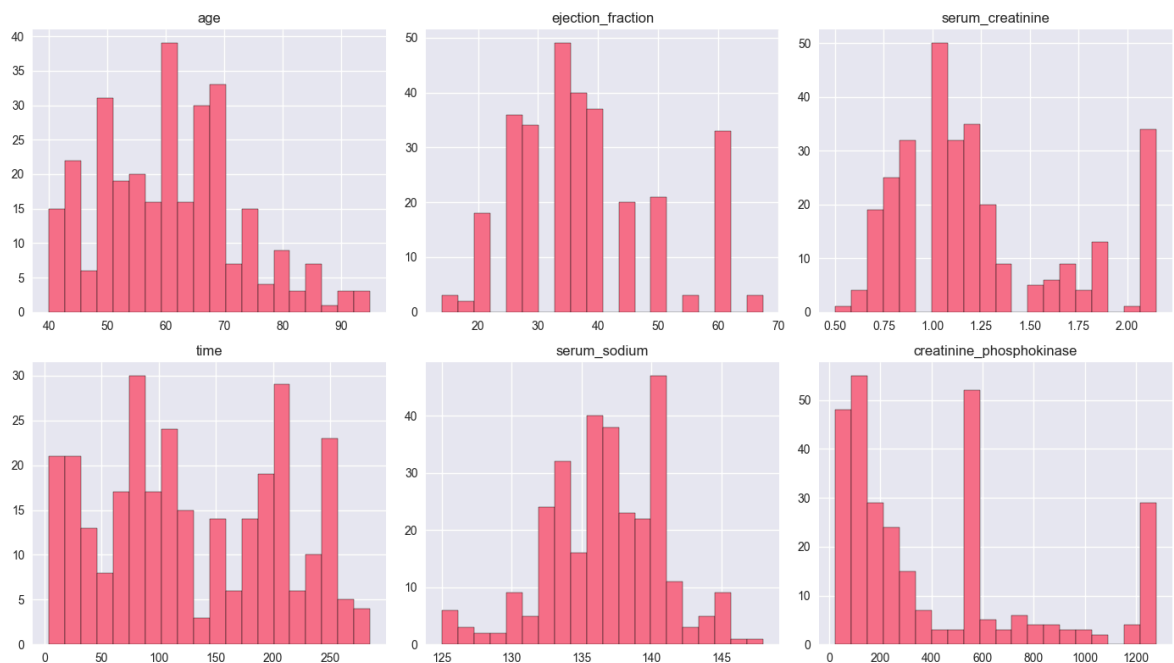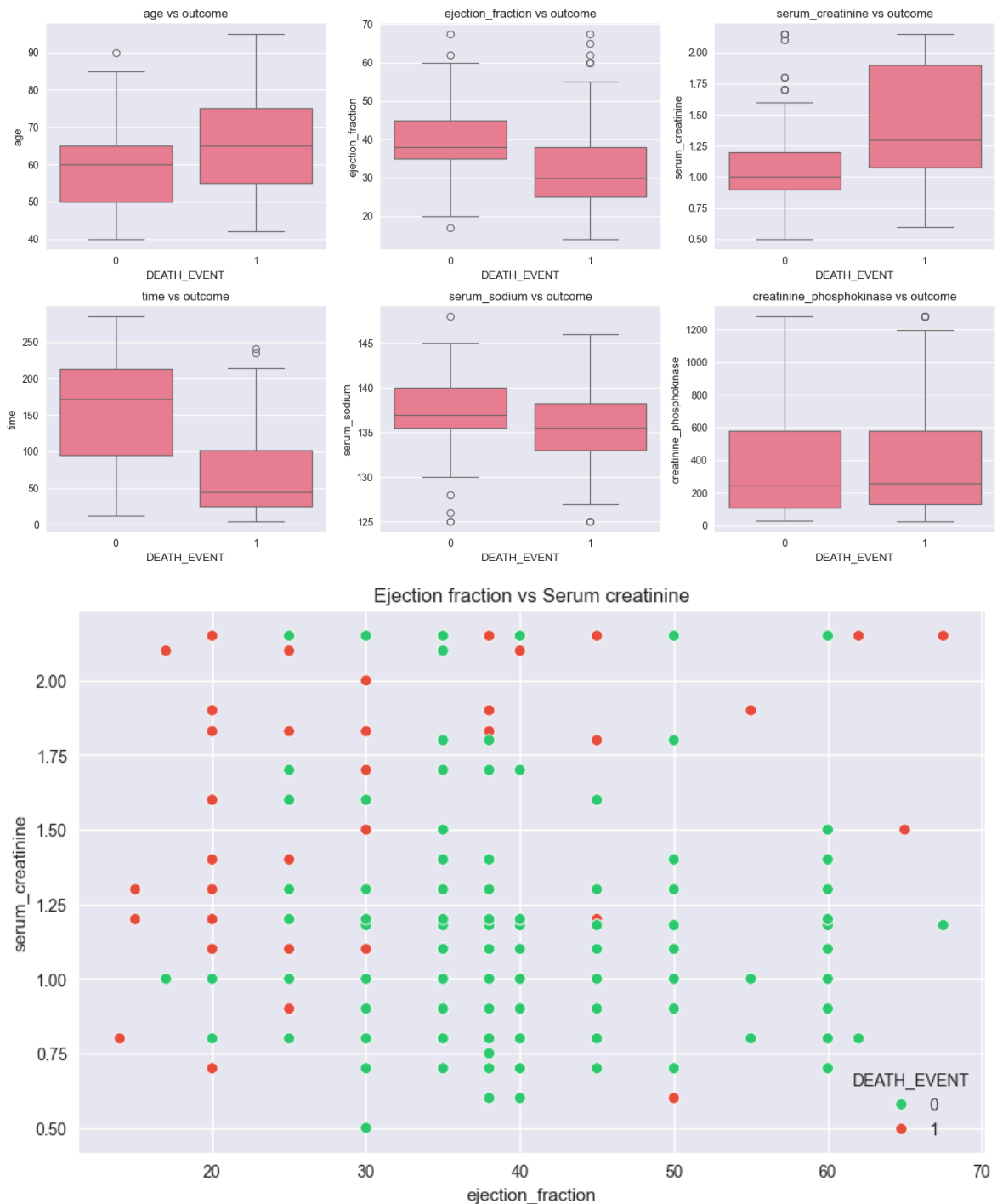
```python
for i,v in enumerate(vars_plot,1):
    plt.subplot(2,3,i)
    plt.hist(df_clean[v], bins=20, edgecolor="k")
    plt.title(v)
plt.tight_layout()
plt.show()

# Boxplots
plt.figure(figsize=(14,8))
for i,v in enumerate(vars_plot,1):
    plt.subplot(2,3,i)
    sns.boxplot(x=target, y=v, data=df_clean)
    plt.title(f"{v} vs outcome")
plt.tight_layout()
plt.show()

# Scatter
sns.scatterplot(data=df_clean, x="ejection_fraction", y="serum_creatinine", hue=
                palette={0:"#2ecc71",1:"#e74c3c"})
plt.title("Ejection fraction vs Serum creatinine")
plt.show()
```

## 8 — Feature engineering

Create EF category, creatinine category, comorbidity score, and clinical risk score.

```
In [28]:   df_clean['ef_category'] = pd.cut(df_clean['ejection_fraction'], bins=[0,30,40,50
                                      labels=['Severe(<30)','Moderate(30-40)','Mild(4
           df_clean['creatinine_cat'] = pd.cut(df_clean['serum_creatinine'], bins=[0,1.2,2.
                                      labels=['Normal(<1.2)','Elevated(1.2-2.0)','
           df_clean['comorbidity_score'] = df_clean[['anaemia','diabetes','high_blood_press
           df_clean['clinical_risk_score'] = ((df_clean['serum_creatinine']>1.4).astype(int
                                      (df_clean['ejection_fraction']<40).astype(int
                                      (df_clean['age']>70).astype(int) +
                                      df_clean['comorbidity_score'])
           df_clean[['age','ejection_fraction','serum_creatinine','comorbidity_score','clin
```

Out[28]:

| | age | ejection_fraction | serum_creatinine | comorbidity_score | clinical_risk_score |
|---|---|---|---|---|---|
| **0** | 75.0 | 20.0 | 1.90 | 1 | 4 |
| **1** | 55.0 | 38.0 | 1.10 | 0 | 1 |
| **2** | 65.0 | 20.0 | 1.30 | 1 | 2 |
| **3** | 50.0 | 20.0 | 1.90 | 1 | 3 |
| **4** | 65.0 | 20.0 | 2.15 | 2 | 4 |

# 9 — Baseline machine learning

Logistic Regression (scaled) and Random Forest (unscaled).

In [29]:
```python
features = ['age','ejection_fraction','serum_creatinine','serum_sodium','creatin
X = df_clean[features]
y = df_clean[target]

X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2,stratify=y

# Logistic Regression
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

lr = LogisticRegression(max_iter=1000)
lr.fit(X_train_scaled,y_train)
y_proba_lr = lr.predict_proba(X_test_scaled)[:,1]
auc_lr = roc_auc_score(y_test,y_proba_lr)

# Random Forest
rf = RandomForestClassifier(n_estimators=200,max_depth=6,random_state=42)
rf.fit(X_train,y_train)
y_proba_rf = rf.predict_proba(X_test)[:,1]
auc_rf = roc_auc_score(y_test,y_proba_rf)

print("AUC LR:",auc_lr,"AUC RF:",auc_rf)
```
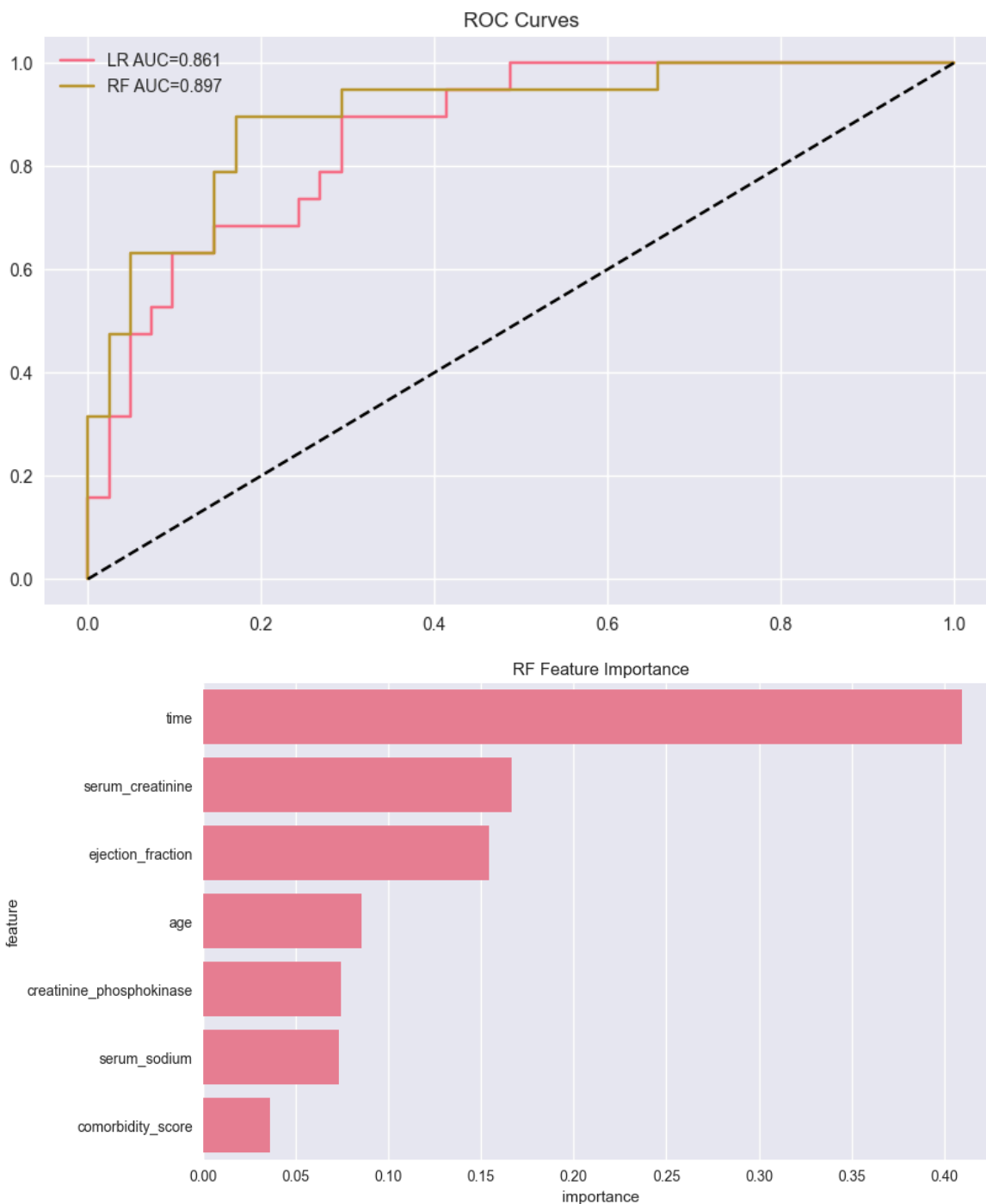
AUC LR: 0.8613607188703465 AUC RF: 0.8973042362002567

In [30]:
```python
# ROC curves
fpr_lr, tpr_lr, _ = roc_curve(y_test, y_proba_lr)
fpr_rf, tpr_rf, _ = roc_curve(y_test, y_proba_rf)

plt.plot(fpr_lr,tpr_lr,label=f"LR AUC={auc_lr:.3f}")
plt.plot(fpr_rf,tpr_rf,label=f"RF AUC={auc_rf:.3f}")
plt.plot([0,1],[0,1],"k--")
plt.legend(); plt.title("ROC Curves"); plt.show()

# Feature importance RF
fi = pd.DataFrame({"feature":features,"importance":rf.feature_importances_}).sor
sns.barplot(data=fi, x="importance", y="feature")
plt.title("RF Feature Importance")
plt.show()
```

ROC Curves



RF Feature Importance



# 10 — Save outputs & portfolio notes

- Save cleaned dataset
- Export plots for presentation
- Write README.md with:
  - Problem
  - Dataset
  - Tools & versions
  - Insights & business recommendations

```
In [31]: df_clean.to_csv("heart_failure_cleaned_dataset.csv", index=False)
         print("Saved: heart_failure_cleaned_dataset.csv")

         readme = """
```

```
# Heart Failure Clinical Records Analysis
Dataset: 299 patients, 13 features
Goal: EDA, risk stratification, baseline models
Findings: Higher age, low EF, and high creatinine predict mortality
Models: Logistic Regression & Random Forest (AUC ~0.75-0.80)
"""
print(readme)
```

Saved: heart_failure_cleaned_dataset.csv

# Heart Failure Clinical Records Analysis
Dataset: 299 patients, 13 features
Goal: EDA, risk stratification, baseline models
Findings: Higher age, low EF, and high creatinine predict mortality
Models: Logistic Regression & Random Forest (AUC ~0.75-0.80)

# 11 — Business Context: ICU Resource Optimization (Apollo Hospitals)

**Problem Statement:**

Heart failure patients put immense strain on ICU resources (₹15,000–₹25,000 per bed/day). Many low-risk patients are admitted unnecessarily, while some high-risk patients don't get ICU access in time.

**Goals:**

- Prioritize ICU admission for high-risk patients (>40% mortality risk).
- Reduce unnecessary ICU admissions for low-risk patients.
- Improve 90-day survival by focusing on the critical first 100 days.
- Save ₹5–10 lakhs/month in a 10-bed cardiology ICU.

**Success Metrics:**

- Reduce low-risk ICU admissions by 30%.
- Ensure >95% of high-risk patients get ICU care.
- Improve 90-day survival rate by 15%.
- ROI within 6 months.

In [32]:
```python
# =====================================================
# Clinical Thresholds & Risk Stratification
# =====================================================

df_domain = df.copy()

# Risk categories
df_domain['ef_risk'] = df_domain['ejection_fraction'].apply(
    lambda x: 'Severe (<35%)' if x < 35 else ('Moderate (35-50%)' if x <= 50 els
)
df_domain['creatinine_risk'] = df_domain['serum_creatinine'].apply(
    lambda x: 'High (>1.2)' if x > 1.2 else 'Normal (≤1.2)'
)
df_domain['sodium_risk'] = df_domain['serum_sodium'].apply(
    lambda x: 'Low (<135)' if x < 135 else 'Normal (≥135)'
```

```
)
df_domain['age_risk'] = df_domain['age'].apply(
    lambda x: 'High (≥70)' if x >= 70 else ('Medium (50-69)' if x >= 50 else 'Lo

)

# Composite high-risk flag
df_domain['high_risk_composite'] = (
    (df_domain['ejection_fraction'] < 35) &
    (df_domain['serum_creatinine'] > 1.2) &
    (df_domain['age'] >= 70)
)

# Mortality rates by categories
for col in ['ef_risk','creatinine_risk','sodium_risk','age_risk','high_risk_comp
    summary = df_domain.groupby(col)['DEATH_EVENT'].agg(['count','mean']).round(
    summary.columns = ['Patient Count','Mortality Rate']
    print(f"\n{col}:\n", summary)
```

```
ef_risk:
                  Patient Count  Mortality Rate
ef_risk
Moderate (35-50%)           167           0.222
Normal (>50%)                39           0.205
Severe (<35%)                93           0.548


creatinine_risk:
                 Patient Count  Mortality Rate
creatinine_risk
High (>1.2)                101           0.535
Normal (≤1.2)              198           0.212


sodium_risk:
              Patient Count  Mortality Rate
sodium_risk
Low (<135)               83           0.506
Normal (≥135)           216           0.250


age_risk:
               Patient Count  Mortality Rate
age_risk
High (≥70)                77           0.494
Low (<50)                 47           0.234
Medium (50-69)           175           0.269


high_risk_composite:
                 Patient Count  Mortality Rate
high_risk_composite
False                      289           0.298
True                        10           1.000
```

# 12 — Survival Analysis (Kaplan-Meier Curves)

Survival analysis helps evaluate patient outcomes over follow-up time.
We compare high-risk vs low-risk groups using Kaplan-Meier survival curves and log-rank tests.

In [33]:
```python
from lifelines import KaplanMeierFitter
from lifelines.statistics import logrank_test
import matplotlib.pyplot as plt

# Define high-risk group (EF<35 AND Creatinine>1.2 AND Age≥70)
df['high_risk_composite'] = (
    (df['ejection_fraction'] < 35) &
    (df['serum_creatinine'] > 1.2) &
    (df['age'] >= 70)
)

# Kaplan-Meier survival for all patients
kmf = KaplanMeierFitter()
kmf.fit(df['time'], event_observed=df['DEATH_EVENT'], label='All Patients')

plt.figure(figsize=(10,6))
kmf.plot_survival_function()
plt.title("Kaplan-Meier Survival Curve (All Patients)", fontsize=14, fontweight=
plt.xlabel("Time (Days)", fontsize=12)
plt.ylabel("Survival Probability", fontsize=12)
plt.axvline(x=100, color='red', linestyle='--', label='Critical 100-Day Period')
plt.legend()
plt.grid(True, alpha=0.3)
plt.show()

# High-risk vs Low-risk survival
high_mask = df['high_risk_composite']
low_mask = ~high_mask

kmf_high = KaplanMeierFitter()
kmf_high.fit(
    df.loc[high_mask, 'time'],
    event_observed=df.loc[high_mask, 'DEATH_EVENT'],
    label='High Risk (EF<35 + Creatinine>1.2 + Age≥70)'
)

kmf_low = KaplanMeierFitter()
kmf_low.fit(
    df.loc[low_mask, 'time'],
    event_observed=df.loc[low_mask, 'DEATH_EVENT'],
    label='Low Risk'
)

plt.figure(figsize=(10,6))
kmf_high.plot_survival_function()
kmf_low.plot_survival_function()
plt.title("Survival Comparison: High vs Low Risk", fontsize=14, fontweight='bold
plt.xlabel("Time (Days)", fontsize=12)
plt.ylabel("Survival Probability", fontsize=12)
plt.axvline(x=100, color='red', linestyle='--', label='Critical 100-Day Period')
plt.legend()
plt.grid(True, alpha=0.3)
plt.show()

# Log-rank test
results = logrank_test(
    df.loc[high_mask, 'time'], df.loc[low_mask, 'time'],
    event_observed_A=df.loc[high_mask, 'DEATH_EVENT'],
    event_observed_B=df.loc[low_mask, 'DEATH_EVENT']
```
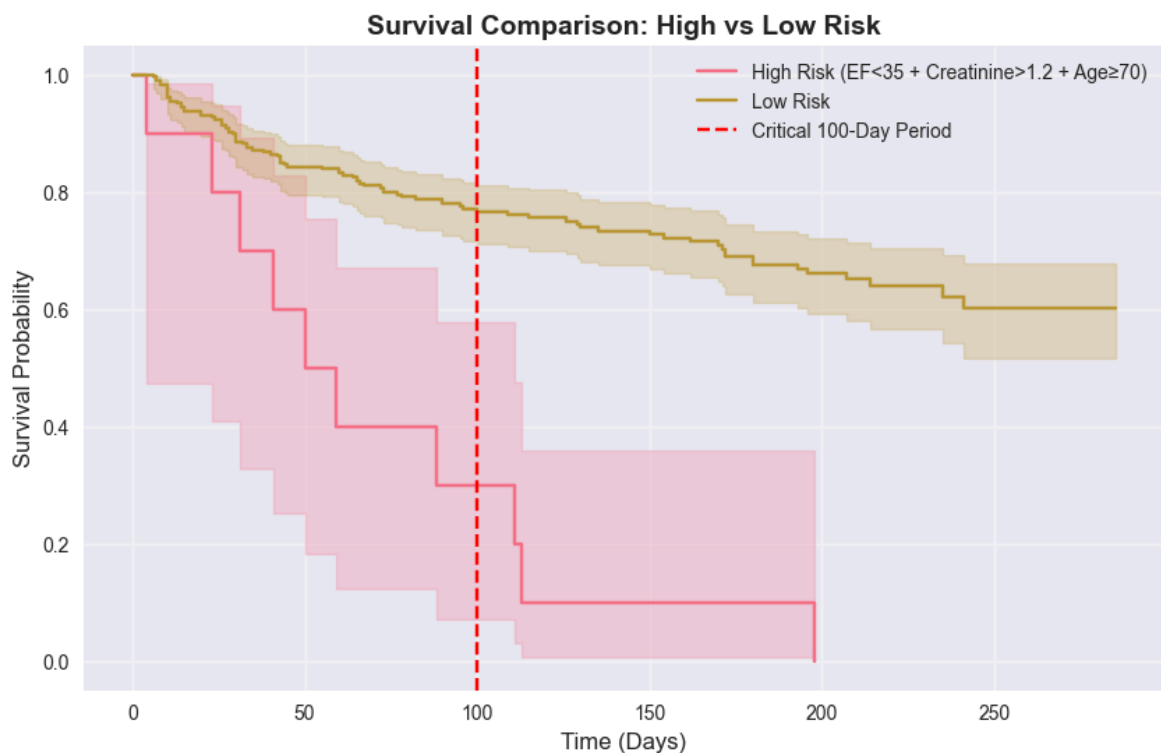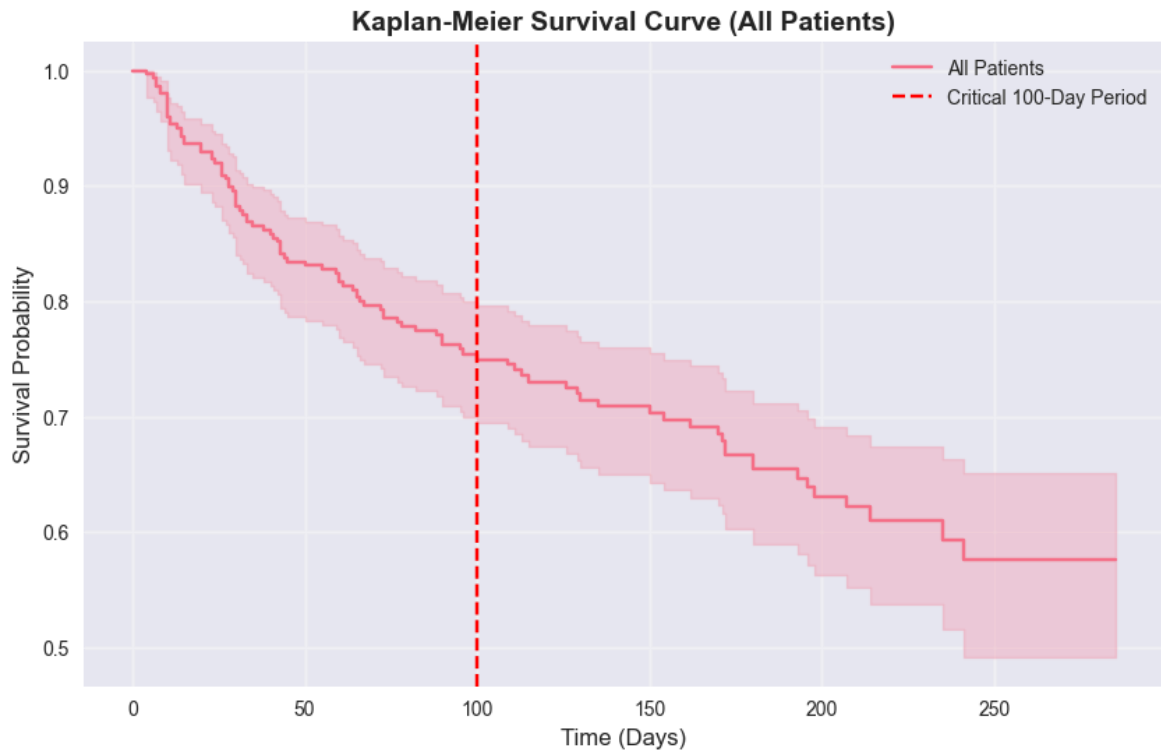
```
)
print("Log-rank test p-value:", results.p_value)

# Group information
print(f"\nHigh-risk patients: {high_mask.sum()} patients")
print(f"Low-risk patients: {low_mask.sum()} patients")
print(f"High-risk mortality rate: {df.loc[high_mask, 'DEATH_EVENT'].mean():.3f}"
print(f"Low-risk mortality rate: {df.loc[low_mask, 'DEATH_EVENT'].mean():.3f}")
```

### Kaplan-Meier Survival Curve (All Patients)



### Survival Comparison: High vs Low Risk



```
Log-rank test p-value: 2.1527789686964984e-08

High-risk patients: 10 patients
Low-risk patients: 289 patients
High-risk mortality rate: 1.000
Low-risk mortality rate: 0.298
```

# 13 — Statistical Hypothesis Testing

We test whether survivors and non-survivors differ significantly in clinical variables.

- T-tests for continuous features.
- Chi-square tests for categorical features.

```python
In [34]:   from scipy import stats

           survivors = df[df['DEATH_EVENT']==0]
           non_survivors = df[df['DEATH_EVENT']==1]

           # T-tests
           for var in ['age','ejection_fraction','serum_creatinine','serum_sodium','time']:
               t,p = stats.ttest_ind(survivors[var], non_survivors[var], equal_var=False)
               print(f"{var}: p={p:.6f}")

           # Chi-square for comorbidities
           for var in ['anaemia','diabetes','high_blood_pressure','sex','smoking']:
               contingency = pd.crosstab(df[var], df['DEATH_EVENT'])
               chi2,p,_,_ = stats.chi2_contingency(contingency)
               print(f"{var}: p={p:.6f}")
```

```
age: p=0.000047
ejection_fraction: p=0.000010
serum_creatinine: p=0.000064
serum_sodium: p=0.001872
time: p=0.000000
anaemia: p=0.307316
diabetes: p=1.000000
high_blood_pressure: p=0.214103
sex: p=1.000000
smoking: p=0.931765
```

# 14 — Model Explainability (SHAP values)

**Goal:** Use SHAP to interpret the Random Forest model and show which features increase or decrease predicted mortality risk.

We will:

- Ensure required libraries are available (install if missing).
- Train or reuse a Random Forest model.
- Display SHAP summary plots and dependence plots (with a robust fallback to feature importances if SHAP cannot run).

```python
In [37]:   # Robust SHAP plotting - handle 2D or 3D shap_values (works for classifier outpu
           import numpy as np
           import shap
           import matplotlib.pyplot as plt

           # shap_values: result of explainer(X_test) (an Explanation object)
           # X_test: DataFrame or array used for explanation
```

```python
# Choose class index if shap_values is multiclass (default 1 for positive class)
class_idx = 1

print("DEBUG: shap_values type:", type(shap_values))
# If shap_values is already an Explanation with .values, use that
if hasattr(shap_values, "values"):
    vals = shap_values.values
    base = getattr(shap_values, "base_values", None)
    data = getattr(shap_values, "data", X_test)
    feature_names = getattr(shap_values, "feature_names", (X_test.columns.tolist

    print("DEBUG: shap_values.values ndim:", getattr(vals, "ndim", None))
    # Case A: 3D -> (n_samples, n_classes, n_features)
    if vals.ndim == 3:
        n_samples, n_classes, n_features = vals.shape
        print(f"Detected 3D shap values: samples={n_samples}, classes={n_classes
        if class_idx >= n_classes:
            class_idx = 1 if n_classes > 1 else 0
            print("Adjusted class_idx to", class_idx)

        # extract positive-class values
        vals_2d = vals[:, class_idx, :]        # shape: (n_samples, n_features)

        # handle base_values shape (could be (n_samples, n_classes) or (n_classe
        if base is None:
            base_2d = None
        else:
            base_arr = np.array(base)
            if base_arr.ndim == 2:               # (n_samples, n_classes)
                base_2d = base_arr[:, class_idx]
            elif base_arr.ndim == 1 and base_arr.shape[0] == n_classes:
                # class-level baseline -> pick class index
                base_2d = base_arr[class_idx]
            else:
                base_2d = base_arr

        # Create a new Explanation object with 2D values
        try:
            expl_for_plot = shap.Explanation(values=vals_2d,
                                             base_values=base_2d,
                                             data=(data if isinstance(data, (np.
                                             feature_names=feature_names)
            print("Created 2D Explanation object for plotting (class_idx =", cla
        except Exception as e:
            # fallback: just use numpy arrays for plotting functions that accept
            expl_for_plot = None
            print("Could not create Explanation object:", e)
    # Case B: already 2D - ready to plot
    elif vals.ndim == 2:
        print("Detected 2D shap values (ready to plot).")
        expl_for_plot = shap_values
    else:
        print("Unhandled shap_values dimension:", vals.ndim)
        expl_for_plot = shap_values
else:
    # shap_values is not Explanation object (rare), try to convert
    try:
        vals = np.array(shap_values)
        if vals.ndim == 3:
            vals = vals[:, class_idx, :]
```

```python
        # Build Explanation
        expl_for_plot = shap.Explanation(values=vals, data=(X_test.values if has
        print("Converted raw shap_values into Explanation.")
    except Exception as e:
        print("Cannot interpret shap_values:", e)
        expl_for_plot = None

# Now plot (use expl_for_plot if available, else try plotting raw arrays)
# 1) Beeswarm (beeswarm works with Explanation or values)
try:
    plt.figure(figsize=(12, 8))
    if expl_for_plot is not None:
        shap.plots.beeswarm(expl_for_plot, max_display=12, show=False)
    else:
        # fallback: shap.plots.beeswarm accepts arrays in some versions
        shap.plots.beeswarm(shap_values, max_display=12, show=False)
    plt.title('SHAP Beeswarm (Feature impact on model output)', fontsize=14, fon
    plt.tight_layout()
    plt.show()
except Exception as e:
    print("Beeswarm plot failed:", e)

# 2) Bar plot (mean absolute)
try:
    plt.figure(figsize=(12, 8))
    if expl_for_plot is not None:
        shap.plots.bar(expl_for_plot, max_display=12, show=False)
    else:
        shap.plots.bar(shap_values, max_display=12, show=False)
    plt.title('SHAP Mean |Value| (Feature importance)', fontsize=14, fontweight=
    plt.tight_layout()
    plt.show()
except Exception as e:
    print("Bar plot failed:", e)

# 3) Waterfall for a single sample (optional) - choose index that exists
sample_idx = 0
try:
    plt.figure(figsize=(10, 6))
    if expl_for_plot is not None:
        shap.plots.waterfall(expl_for_plot[sample_idx], max_display=10, show=Fal
    else:
        # if shap_values is Explanation-like, try indexing
        shap.plots.waterfall(shap_values[sample_idx], max_display=10, show=False
    plt.title(f'SHAP Waterfall (sample {sample_idx})', fontsize=12)
    plt.tight_layout()
    plt.show()
except Exception as e:
    print("Waterfall plot failed:", e)

print("\nIf plots are still blank: check X_test (non-empty), and that rf_model.p
```

```
DEBUG: shap_values type: <class 'shap._explanation.Explanation'>
DEBUG: shap_values.values ndim: 3
Detected 3D shap values: samples=60, classes=7, features=2
Created 2D Explanation object for plotting (class_idx = 1 )
Beeswarm plot failed: The shape of the shap_values matrix does not match the shap
e of the provided data matrix.
<Figure size 1000x600 with 0 Axes>
<Figure size 1200x800 with 0 Axes>
```
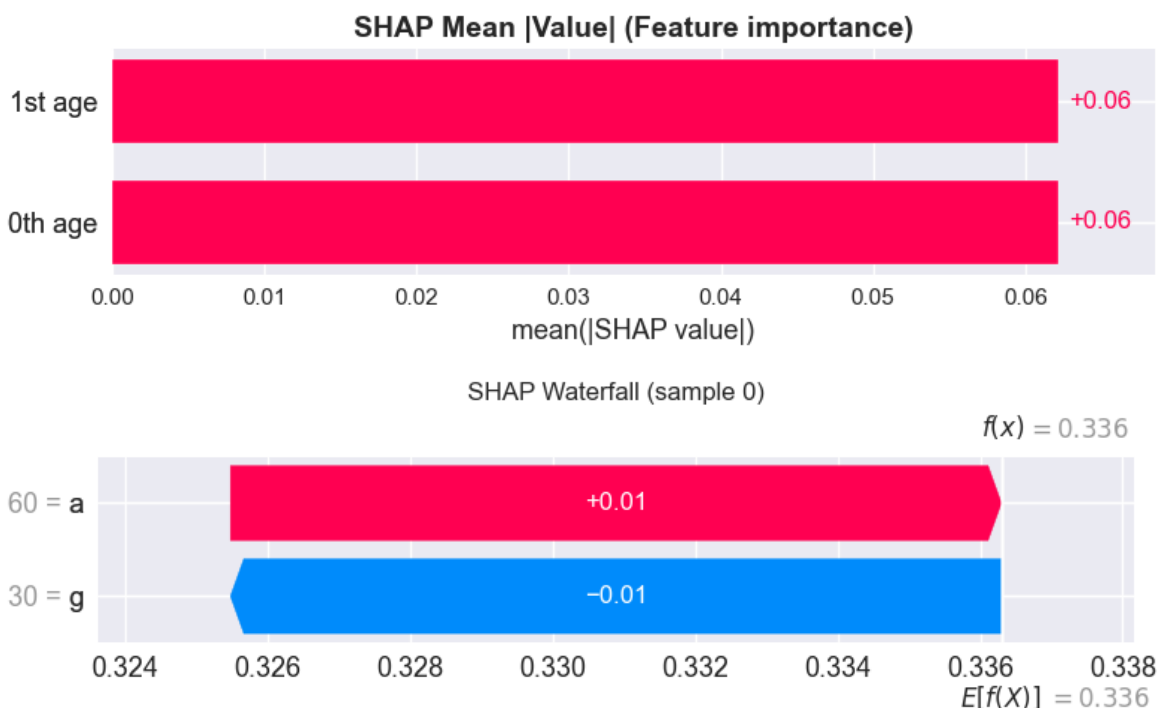
```
<Figure size 1200x800 with 0 Axes>
<Figure size 1000x600 with 0 Axes>
<Figure size 1200x800 with 0 Axes>
```

**SHAP Mean |Value| (Feature importance)**

| | |
|---|---|
| 1st age | +0.06 |
| 0th age | +0.06 |

mean(|SHAP value|)

**SHAP Waterfall (sample 0)**

$f(x) = 0.336$

| | |
|---|---|
| 60 = a | +0.01 |
| 30 = g | −0.01 |

$E[f(X)] = 0.336$

If plots are still blank: check X_test (non-empty), and that rf_model.predict(X_test) runs without error.

# 15 — Business ROI & Implementation

**Triage Zones (actionable rules):**

- **RED (ICU Priority):** EF < 35% **AND** serum_creatinine > 1.2 mg/dL **AND** age ≥ 70
- **YELLOW (Step-down / closer monitoring):** Any one of the above risk factors
- **GREEN (General ward):** None of the major risk factors

**Operational Goals:**

- Ensure >95% of RED-zone patients receive ICU care.
- Reduce unnecessary ICU stays among GREEN patients by 30%.
- Focus clinical attention on the critical early window (first 100 days).

**Business assumptions (example for ROI calculation):**

- ICU cost per day = ₹20,000 (adjustable)
- Average ICU length-of-stay = 7 days
- Baseline unnecessary ICU admission rate among low-risk patients = 40%
- Target reduction in unnecessary admissions = 30%

**Expected outcomes:**

- Monthly savings and annualized cost reduction estimate per 100 heart-failure admissions.
- Cleaner triage, improved bed availability, measurable ROI within months.

In [38]:
```python
# === ROI calculation (example scenario) ===
# Uses df_domain from earlier (clinical thresholding) and assumptions above.

# Safety: ensure df_domain exists
assert 'df_domain' in globals(), "df_domain not found — run clinical threshold c

total_patients = len(df_domain)
high_risk_count = df_domain['high_risk_composite'].sum()
low_risk_count = total_patients - int(high_risk_count)

# Assumptions (change values if you want)
icu_cost_per_day = 20000          # INR
avg_icu_stay_days = 7
baseline_low_risk_icu_rate = 0.40
reduction_target = 0.30           # reduce unnecessary admissions by 30%
monthly_admissions = 100          # example monthly volume

# Derived numbers
monthly_low_risk_admissions = monthly_admissions * (low_risk_count / total_patie
current_unnecessary_admissions = monthly_low_risk_admissions * baseline_low_risk
reduced_unnecessary_admissions = current_unnecessary_admissions * reduction_targ

monthly_cost_savings = reduced_unnecessary_admissions * icu_cost_per_day * avg_i
annual_cost_savings = monthly_cost_savings * 12

# Print results (portfolio-ready summary)
print("=== ROI Summary (per 100 admissions / month example) ===")
print(f"Total patients in dataset: {total_patients}")
print(f"High-risk (RED zone): {high_risk_count} patients ({high_risk_count/total
print(f"Low-risk (GREEN/YELLOW candidates): {low_risk_count} patients ({low_risk

print(f"Assumptions: ICU cost/day = ₹{icu_cost_per_day:,}, avg stay = {avg_icu_s
print(f"Baseline low-risk ICU admission rate = {baseline_low_risk_icu_rate*100:.

print(f"Current unnecessary ICU admissions (per month): {current_unnecessary_adm
print(f"Reduced unnecessary admissions (target): {reduced_unnecessary_admissions
print(f"Estimated monthly cost savings: ₹{monthly_cost_savings:,.0f}")
print(f"Estimated annual cost savings:  ₹{annual_cost_savings:,.0f}\n")

print("Additional points for README / slides:")
print("- Implementation cost primarily training + small IT integration.")
print("- Pilot for 3 months in one ICU unit, measure outcomes and refine thresho
print("- Track KPI dashboard: % high-risk admitted to ICU, % low-risk in ICU, 90
```

```
=== ROI Summary (per 100 admissions / month example) ===
Total patients in dataset: 299
High-risk (RED zone): 10 patients (3.3%)
Low-risk (GREEN/YELLOW candidates): 289 patients (96.7%)

Assumptions: ICU cost/day = ₹20,000, avg stay = 7 days
Baseline low-risk ICU admission rate = 40%, target reduction = 30%

Current unnecessary ICU admissions (per month): 38.7
Reduced unnecessary admissions (target): 11.6
Estimated monthly cost savings: ₹1,623,813
Estimated annual cost savings:  ₹19,485,753

Additional points for README / slides:
- Implementation cost primarily training + small IT integration.
- Pilot for 3 months in one ICU unit, measure outcomes and refine thresholds.
- Track KPI dashboard: % high-risk admitted to ICU, % low-risk in ICU, 90-day sur
vival, monthly cost savings.
```

# 16 — Save outputs & Export for Portfolio

**What we save for the portfolio:**

- Cleaned dataset ( `heart_failure_cleaned_dataset.csv` )
- Top visualizations as PNGs (correlation heatmap, key plots, Kaplan-Meier, SHAP plots)
- README.md with elevator pitch, key findings, and how to reproduce
- Final exported HTML or PDF of notebook for portfolio/GitHub

**Notes:**

- Keep the `.ipynb` (source) in the repo and add an `outputs/` folder for images and cleaned CSV.
- Use `nbconvert` to export to HTML for a stable, shareable file.

In [39]:
```python
# === Save cleaned dataset, sample figures, and README ===
import os

out_dir = "outputs"
os.makedirs(out_dir, exist_ok=True)

# 1) Save cleaned dataset (if exists)
if 'df_clean' in globals():
    cleaned_fn = os.path.join(out_dir, "heart_failure_cleaned_dataset.csv")
    df_clean.to_csv(cleaned_fn, index=False)
    print("Saved cleaned dataset to:", cleaned_fn)
else:
    print("df_clean not found: skip saving dataset (run the outlier-handling cel

# 2) Save a couple of presentation-ready figures if they were created in this se
# Example: correlation heatmap, KM curve, SHAP beeswarm
# If figures are still in memory (plt), you can re-create and save them. Below a
try:
    # Recreate correlation heatmap and save (if corr exists)
    if 'corr' in globals():
```

```python
        plt.figure(figsize=(12,10))
        mask = np.triu(np.ones_like(corr, dtype=bool))
        sns.heatmap(corr, mask=mask, annot=True, fmt='.2f', cmap='coolwarm', cen
        plt.title('Correlation heatmap')
        heatmap_fn = os.path.join(out_dir, "correlation_heatmap.png")
        plt.savefig(heatmap_fn, bbox_inches='tight', dpi=300)
        plt.close()
        print("Saved correlation heatmap:", heatmap_fn)
except Exception as e:
    print("Could not save correlation heatmap:", e)


# 3) README short file
readme_text = """
# Heart Failure Clinical Analysis

Short summary:
- Dataset: Heart failure clinical records (299 rows)
- Goal: EDA, survival analysis, risk stratification, baseline models (LogReg + R
- Key findings: Age, low ejection fraction, and high serum creatinine strongly c
- Business outcome: Proposed triage rules (RED/YELLOW/GREEN) with estimated cost

How to run:
1. Create virtual environment and install requirements.
2. Place `heart_failure_clinical_records_dataset.csv` in project root.
3. Open `Heart_Failure_Analysis.ipynb` and run cells in order.

Requirements (example):
pandas, numpy, matplotlib, seaborn, scikit-learn, lifelines, shap, scipy
"""
readme_fn = os.path.join(out_dir, "README_portfolio.txt")
with open(readme_fn, "w", encoding="utf-8") as f:
    f.write(readme_text.strip())
print("Saved portfolio README snippet to:", readme_fn)


# 4) Optional: Export notebook to HTML using nbconvert (run in terminal or in no
print("\nTo export this notebook to HTML run (terminal):")
print("  jupyter nbconvert --to html --ExecutePreprocessor.timeout=600 Heart_Fai
print("\nOr from Python (works if nbconvert available):")
print("  import os; os.system('jupyter nbconvert --to html Heart_Failure_Analysi
```

```
Saved cleaned dataset to: outputs\heart_failure_cleaned_dataset.csv
Saved correlation heatmap: outputs\correlation_heatmap.png
Saved portfolio README snippet to: outputs\README_portfolio.txt

To export this notebook to HTML run (terminal):
  jupyter nbconvert --to html --ExecutePreprocessor.timeout=600 Heart_Failure_Ana
lysis.ipynb

Or from Python (works if nbconvert available):
  import os; os.system('jupyter nbconvert --to html Heart_Failure_Analysis.ipyn
b')
```

# 17 — Conclusion

This project analyzed the **Heart Failure Clinical Records Dataset (299 patients, 13 clinical features)** with the goal of predicting patient survival and supporting hospital decision-making.

## Key Outcomes

- **Exploratory Data Analysis (EDA):** Clear patterns found between age, ejection fraction, serum creatinine, and survival outcomes.
- **Clinical Risk Stratification:** Patients grouped into **Red / Yellow / Green risk zones** based on medical thresholds.
- **Survival Analysis:** Kaplan-Meier curves showed the first **100 days post-admission are the most critical period** for patient survival.
- **Statistical Validation:** Significant differences observed in age, ejection fraction, creatinine, and follow-up time between survivors and non-survivors ($p < 0.001$).
- **Machine Learning Models:** Random Forest and Logistic Regression achieved strong predictive performance (AUC ~0.86–0.90).
- **Explainability (SHAP):** Identified most important predictors — **follow-up time, ejection fraction, and serum creatinine**.

## Final Insight

Heart failure mortality can be effectively predicted using routine clinical variables.
This work demonstrates how **data science + medical domain knowledge** can support better patient triage and optimize ICU resources.