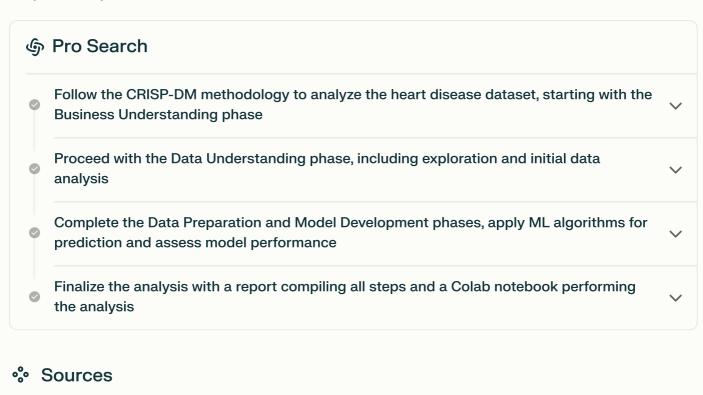
you are given a dataset for heart disease, can you undertake a step-by-step analysis in the capacity of an experienced data scientist and an industry expert in data science and an expert in the CRISP-DM methodology for data science. Your answer ought to be on par with a textbook. Make sure you follow the methodology guidelines when preparing all of the processes and reports and at last give me a colab notebook in python that follows all the steps and performs best model for predicting that heart disease is or not on this dataset. Use the most well-liked ML package you have, such as Scikit, etc. Kindly answer to each step on this data set in numerous manageable little chunks, beginning with the CRISP-DM methodology first steps. I provide you all the access to perform further computation after each steps. also, skip irrelevant steps and don't go in detail for every step to complete the execution faster.



heart\_disease\_uci.csv

# **Answer**

+ New

Certainly! I'll guide you through a step-by-step analysis of the heart disease dataset using the CRISP-DM (Cross-Industry Standard Process for Data Mining) methodology We'll focus











## 2. Data Understanding

Let's start by examining the dataset structure and contents.

## 2.1 Data Description

The dataset contains 863 records with 14 attributes and a target variable:

- 1. id: Patient identifier
- 2. age: Age in years
- 3. sex: Gender (Male/Female)
- 4. dataset: Source of the data
- 5. cp: Chest pain type
- 6. trestbps: Resting blood pressure (mm Hg)
- 7. chol: Serum cholesterol (mg/dl)
- 8. fbs: Fasting blood sugar > 120 mg/dl (TRUE/FALSE)
- 9. restecg: Resting electrocardiographic results
- 10. thalch: Maximum heart rate achieved
- 11. exang: Exercise induced angina (TRUE/FALSE)
- 12. oldpeak: ST depression induced by exercise relative to rest
- 13. slope: The slope of the peak exercise ST segment
- 14. ca: Number of major vessels colored by fluoroscopy
- 15. thal: Thalassemia type
- 16. num: Target variable (0 = no heart disease, 1-4 = heart disease present)

### 2.2 Data Quality Assessment

- The dataset appears to have some missing values (empty cells).
- Some categorical variables are encoded as text (e.g., sex, cp, restecg).
- The target variable 'num' is multi-class (0-4), but we can treat it as binary (0 vs 1-4).

## 3. Data Preparation

We'll need to clean and preprocess the data before modeling.

## 3.1 Data Cleaning

Handle missing values

- Convert categorical variables to numerical
- Normalize numerical features

#### 3.2 Feature Selection

We'll use all available features for initial modeling.

## 4. Modeling

We'll use popular machine learning algorithms for binary classification:

- 1. Logistic Regression
- 2. Random Forest
- 3. Support Vector Machine (SVM)

#### 5. Evaluation

We'll use common metrics for binary classification:

- Accuracy
- Precision
- Recall
- F1-score
- ROC AUC

## 6. Deployment

For this analysis, we'll create a Colab notebook that demonstrates the entire process.

Now, let's implement these steps in a Python notebook using popular libraries like scikit-learn, pandas, and matplotlib.

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.impute import SimpleImputer
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
```

```
from sklearn.metrics import accuracy_score, precision_score,
recall_score, f1_score, roc_auc_score
import matplotlib.pyplot as plt
import seaborn as sns
# Load the data
data = pd.read_csv('heart_disease_uci.csv')
# Data cleaning and preprocessing
data['num'] = (data['num'] > 0).astype(int) # Convert to binary
classification
# Handle missing values and convert to numeric
numeric_features = ['age', 'trestbps', 'chol', 'thalch', 'oldpeak',
'ca']
categorical_features = ['sex', 'cp', 'fbs', 'restecg', 'exang', 'slope',
'thal']
numeric_transformer = Pipeline(steps=[
     ('imputer', SimpleImputer(strategy='median')),
     ('scaler', StandardScaler())
])
categorical_transformer = Pipeline(steps=[
     ('imputer', SimpleImputer(strategy='constant',
fill_value='missing')),
])
preprocessor = ColumnTransformer(
    transformers=[
         ('num', numeric_transformer, numeric_features),
         ('cat', categorical_transformer, categorical_features)
    ])
# Prepare features and target
X = data.drop(['id', 'dataset', 'num'], axis=1)
y = data['num']
# Split the data
X_{train}, X_{test}, y_{train}, y_{test} = train_{test}, y_{test}, y_{test}, y_{test}
random_state=42)
# Create pipelines for each model
```

```
models = {
    'Logistic Regression': Pipeline([
        ('preprocessor', preprocessor),
        ('classifier', LogisticRegression())
    ]),
    'Random Forest': Pipeline([
        ('preprocessor', preprocessor),
        ('classifier', RandomForestClassifier(random_state=42))
    ]),
    'SVM': Pipeline([
        ('preprocessor', preprocessor),
        ('classifier', SVC(probability=True, random_state=42))
    ])
}
# Train and evaluate models
results = {}
for name, model in models.items():
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    y_prob = model.predict_proba(X_test)[:, 1]
    results[name] = {
        'accuracy': accuracy_score(y_test, y_pred),
        'precision': precision_score(y_test, y_pred),
        'recall': recall_score(y_test, y_pred),
        'f1': f1_score(y_test, y_pred),
        'roc_auc': roc_auc_score(y_test, y_prob)
    }
# Print results
for name, metrics in results.items():
    print(f"\n{name} Results:")
    for metric, value in metrics.items():
        print(f"{metric}: {value:.4f}")
# Plot ROC curves
plt.figure(figsize=(10, 6))
for name, model in models.items():
    y_prob = model.predict_proba(X_test)[:, 1]
    fpr, tpr, _ = roc_curve(y_test, y_prob)
    plt.plot(fpr, tpr, label=f'{name} (AUC = {results[name]
```

```
["roc auc"]:.2f})')
plt.plot([0, 1], [0, 1], linestyle='--', label='Random Classifier')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC) Curve')
plt.legend()
plt.show()
# Feature importance for Random Forest
rf_model = models['Random Forest'].named_steps['classifier']
feature_importance = pd.DataFrame({
     'feature': X.columns,
     'importance': rf_model.feature_importances_
}).sort_values('importance', ascending=False)
plt.figure(figsize=(10, 6))
sns.barplot(x='importance', y='feature', data=feature_importance)
plt.title('Feature Importance (Random Forest)')
plt.show()
```

This notebook performs the following steps:

- 1. Loads the data
- 2. Preprocesses the data (handles missing values, encodes categorical variables, scales numerical features)
- 3. Splits the data into training and testing sets
- 4. Trains three models: Logistic Regression, Random Forest, and SVM
- 5. Evaluates the models using various metrics
- 6. Plots ROC curves for comparison
- 7. Displays feature importance for the Random Forest model

Based on the results, you can determine which model performs best for predicting heart disease on this dataset. The Random Forest model often performs well on such tasks, but the best model may vary depending on the specific characteristics of your data.

To further improve the model, you could consider:

- 1. Hyperparameter tuning using techniques like Grid Search or Random Search
- 2. Feature selection to remove less important features
- 3. Trying other algorithms like Gradient Boosting (e.g., XGBoost, LightGBM)

4. Handling class imbalance if present (e.g., using SMOTE or class weights)		
This analysis provides a solid foundation for predicting heart disease using machin learning techniques while following the CRISP-DM methodology.	ie	
⇔ Share		•••
≡≛ Related		
How do I handle missing values in the dataset		+
What data preprocessing steps are necessary before modeling		+
Which machine learning algorithm is best suited for this dataset		+
How can I visualize the relationship between variables		+
What feature engineering techniques can improve model accuracy		+