

```

# Step 1: Upload CSV in Colab
from google.colab import files
import pandas as pd

uploaded = files.upload() # Select your kidney_disease.csv

# Load dataset
df = pd.read_csv("kidney_disease.csv")
print("Shape:", df.shape)
df.head()

<IPython.core.display.HTML object>

Saving kidney_disease.csv to kidney_disease (1).csv
Shape: (400, 26)

{"type": "dataframe", "variable_name": "df"}

import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, StandardScaler

# Drop ID column if it exists
if 'id' in df.columns:
    df = df.drop('id', axis=1)

# Handle missing values (fill numeric with median, categorical with mode)
for col in df.columns:
    if df[col].dtype == 'object':
        df[col].fillna(df[col].mode()[0], inplace=True)
    else:
        df[col].fillna(df[col].median(), inplace=True)

# Encode categorical columns
le = LabelEncoder()
for col in df.columns:
    if df[col].dtype == 'object':
        df[col] = le.fit_transform(df[col])

# Define X and y
X = df.drop('classification', axis=1) # assuming 'classification' is target
y = df['classification']

# Scale features
scaler = StandardScaler()
X = scaler.fit_transform(X)

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y,

```

```
test_size=0.2, random_state=42, stratify=y)
```

```
print("Training shape:", X_train.shape)
```

```
print("Testing shape:", X_test.shape)
```

```
Training shape: (320, 24)
```

```
Testing shape: (80, 24)
```

/tmp/ipython-input-214255797.py:14: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.  
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
df[col].fillna(df[col].median(), inplace=True)
```

/tmp/ipython-input-214255797.py:12: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.  
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
df[col].fillna(df[col].mode()[0], inplace=True)
```

```
from xgboost import XGBClassifier
```

```
from sklearn.metrics import accuracy_score, classification_report
```

```
# Train model
```

```
model = XGBClassifier(use_label_encoder=False, eval_metric='logloss',  
random_state=42)
```

```
model.fit(X_train, y_train)
```

```
# Predictions
```

```
y_pred = model.predict(X_test)
```

```
# Evaluate
```

```
acc = accuracy_score(y_test, y_pred)
```

```

print("XGBoost Accuracy:", acc)
print(classification_report(y_test, y_pred))

/usr/local/lib/python3.12/dist-packages/xgboost/training.py:183:
UserWarning: [14:05:54] WARNING: /workspace/src/learner.cc:738:
Parameters: { "use_label_encoder" } are not used.

bst.update(dtrain, iteration=i, fobj=obj)

```

	precision	recall	f1-score	support
0	0.96	1.00	0.98	50
2	1.00	0.93	0.97	30
accuracy			0.97	80
macro avg	0.98	0.97	0.97	80
weighted avg	0.98	0.97	0.97	80

```

from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report

# Train
rf = RandomForestClassifier(random_state=42)
rf.fit(X_train, y_train)

# Predict
y_pred_rf = rf.predict(X_test)

# Evaluate
acc_rf = accuracy_score(y_test, y_pred_rf)
print("Random Forest Accuracy:", acc_rf)
print(classification_report(y_test, y_pred_rf))

```

	precision	recall	f1-score	support
0	0.98	1.00	0.99	50
2	1.00	0.97	0.98	30
accuracy			0.99	80
macro avg	0.99	0.98	0.99	80
weighted avg	0.99	0.99	0.99	80

```

from lightgbm import LGBMClassifier
from sklearn.metrics import accuracy_score, classification_report

# Train with regularization + constraints
lgb = LGBMClassifier(

```



[illegible]

[illegible]

[illegible]

[illegible]



[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

□ LightGBM Accuracy: 0.9875

	precision	recall	f1-score	support
0	0.98	1.00	0.99	50
2	1.00	0.97	0.98	30
accuracy			0.99	80



macro avg	0.99	0.98	0.99	80
weighted avg	0.99	0.99	0.99	80

```
/usr/local/lib/python3.12/dist-packages/sklearn/utils/
validation.py:2739: UserWarning: X does not have valid feature names,
but LGBMClassifier was fitted with feature names
warnings.warn(
```

```
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
import numpy as np # Import numpy
```

```
# 1. Class distribution
```

```
sns.countplot(x='classification', data=df)
plt.title("Class Distribution")
plt.show()
```

```
# 2. Correlation heatmap
```

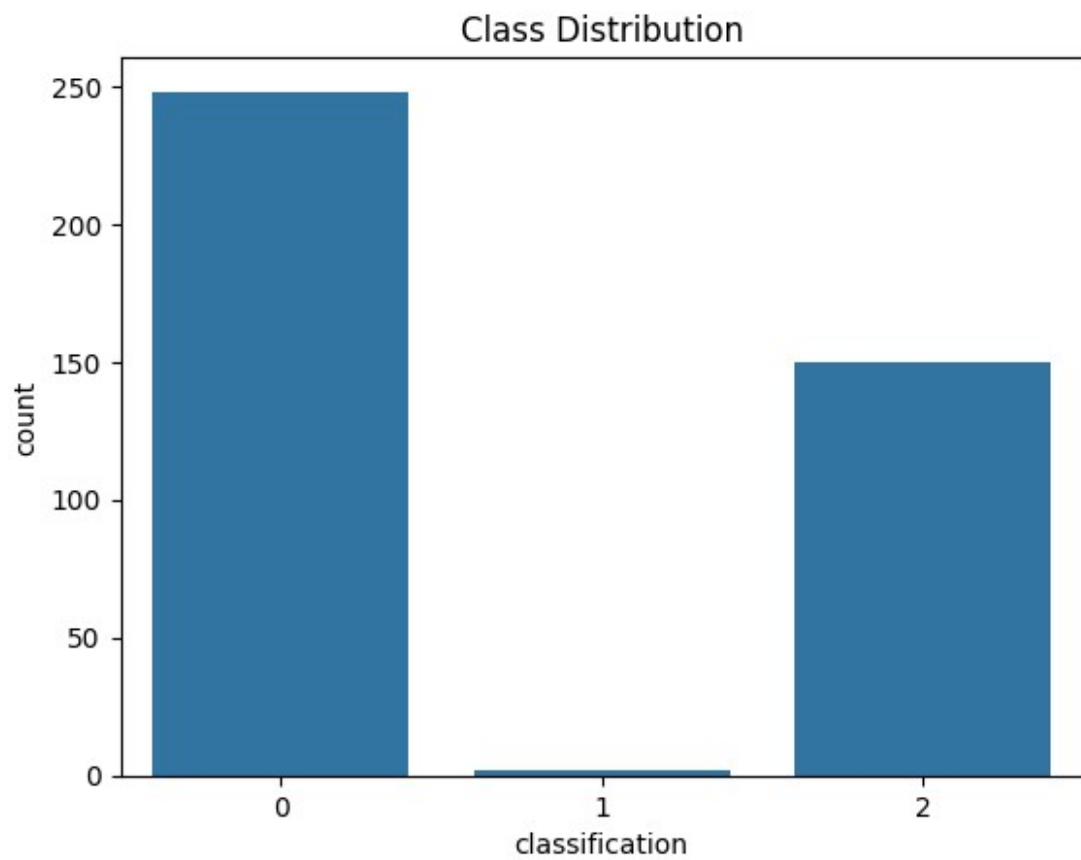
```
plt.figure(figsize=(12,8))
sns.heatmap(df.corr(), cmap="coolwarm", annot=False)
plt.title("Feature Correlation Heatmap")
plt.show()
```

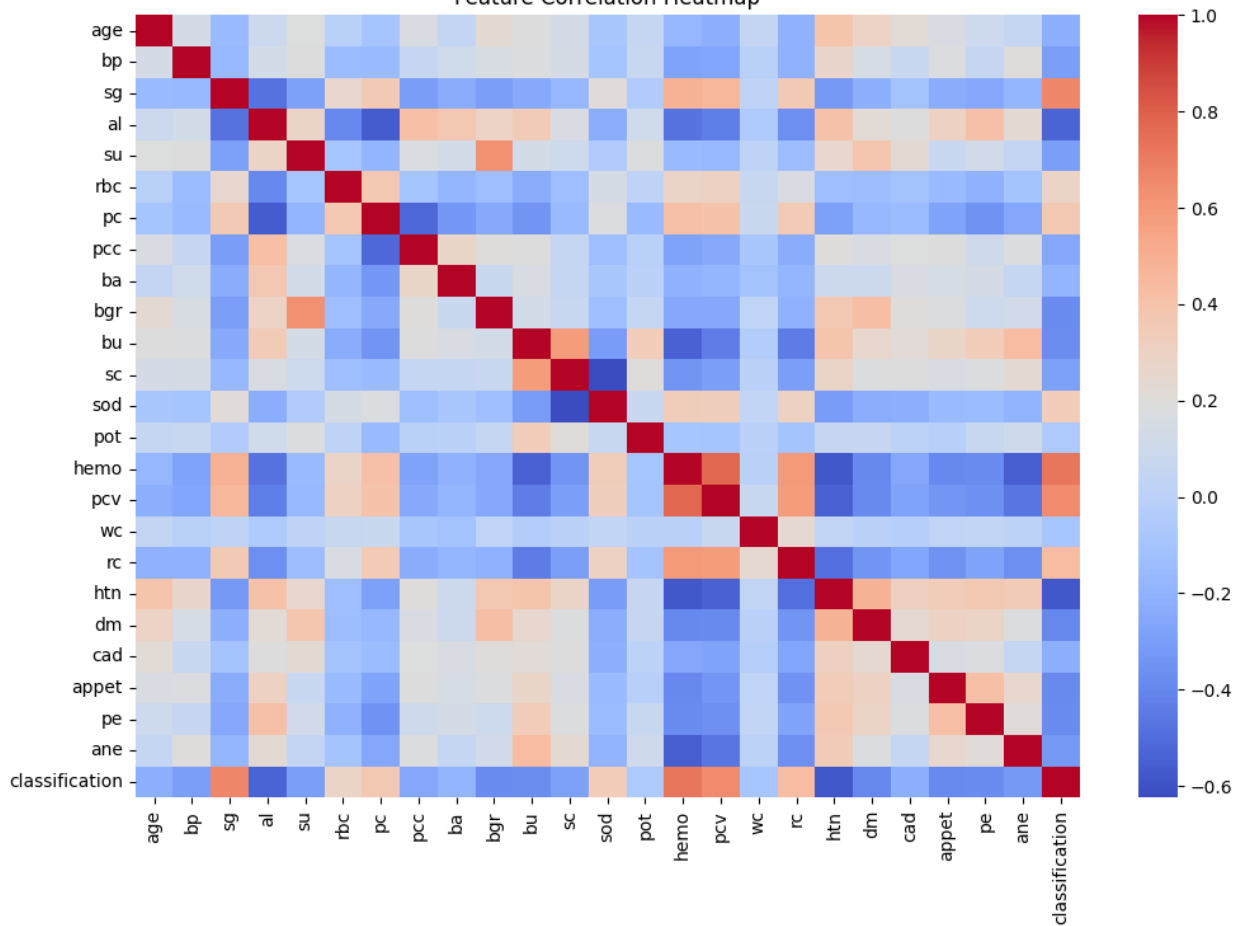
```
# 3. Confusion Matrix
```

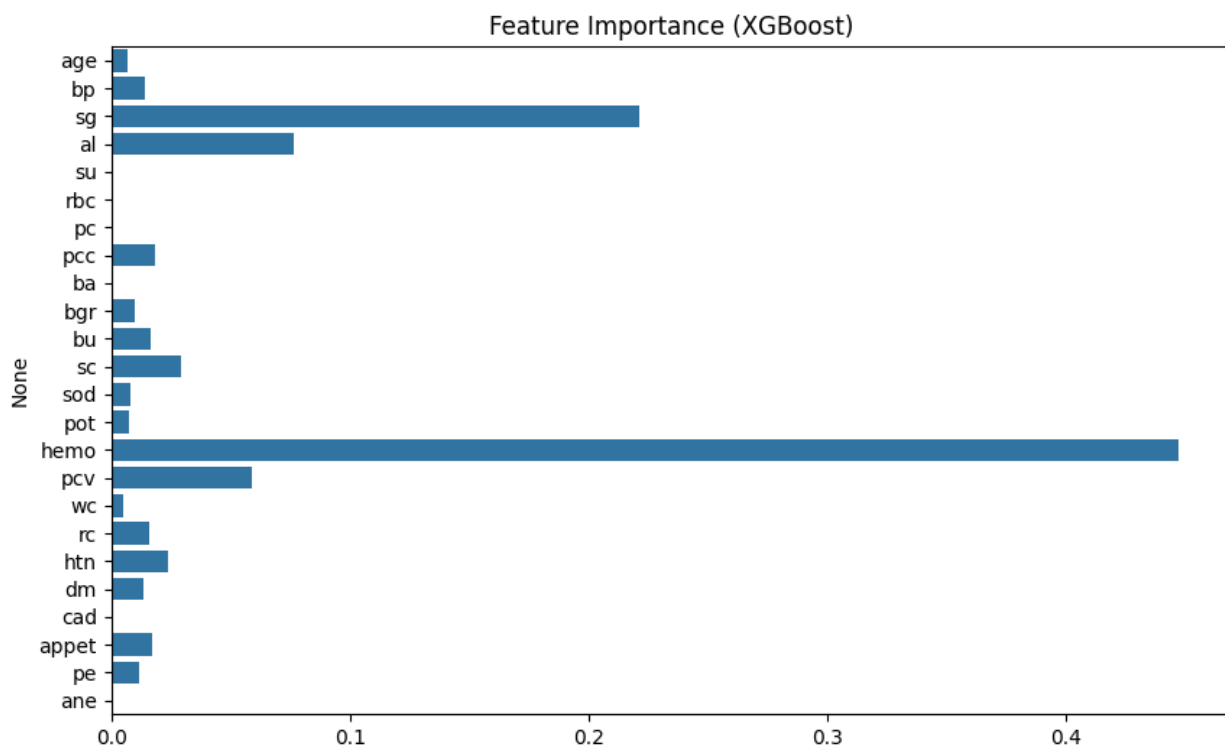
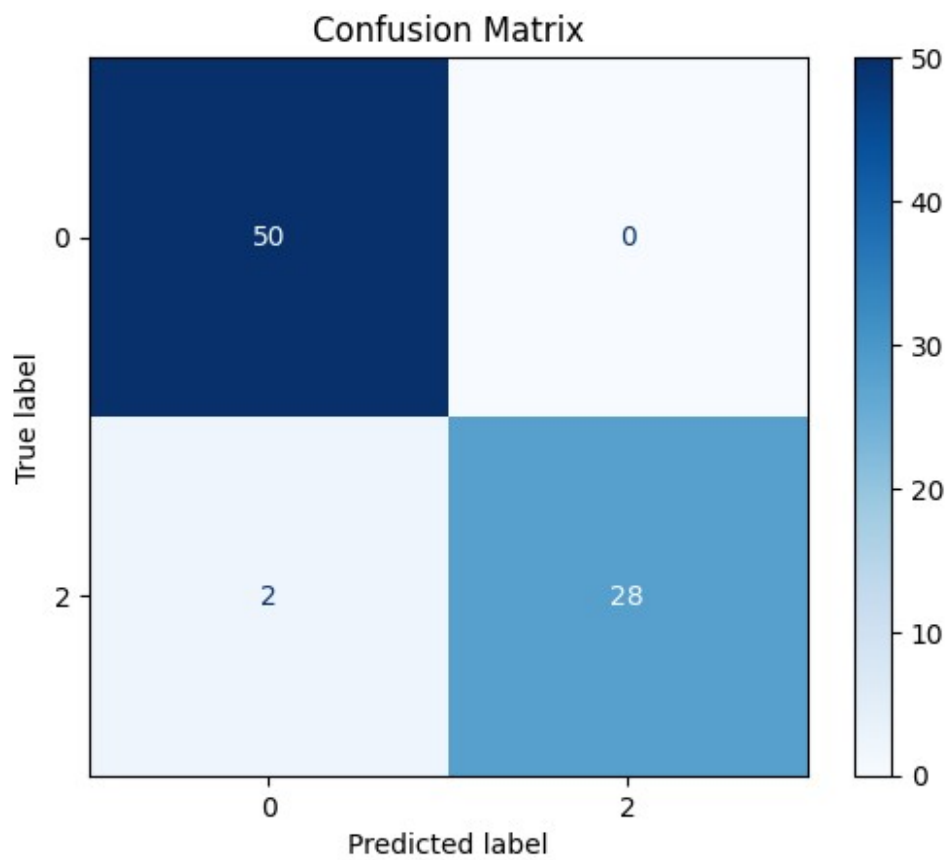
```
cm = confusion_matrix(y_test, y_pred)
# Use unique values from y_test for display labels
disp = ConfusionMatrixDisplay(confusion_matrix=cm,
display_labels=np.unique(y_test))
disp.plot(cmap="Blues")
plt.title("Confusion Matrix")
plt.show()
```

```
# 4. Feature Importance
```

```
plt.figure(figsize=(10,6))
sns.barplot(x=model.feature_importances_, y=df.drop('classification',
axis=1).columns)
plt.title("Feature Importance (XGBoost)")
plt.show()
```







```

from google.colab import files
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, StandardScaler
from xgboost import XGBClassifier

try:
    df = pd.read_csv("kidney_disease.csv")
    print("Shape:", df.shape)
except FileNotFoundError:
    print("kidney_disease.csv not found. Please upload the file using  

the 'files.upload()' cell.")
    # Provide a placeholder DataFrame or exit if file not found
    df = pd.DataFrame() # Placeholder

if not df.empty:
    # Drop ID column if it exists
    if 'id' in df.columns:
        df = df.drop('id', axis=1)

        if df[col].dtype == 'object':
            df[col] = df[col].fillna(df[col].mode()[0])
        else:
            df[col] = df[col].fillna(df[col].median())

    le = LabelEncoder()
    for col in df.columns:
        if df[col].dtype == 'object':
            df[col] = le.fit_transform(df[col])

    X = df.drop('classification', axis=1) # assuming 'classification'
is target
    y = df['classification']

    # Scale features
    scaler = StandardScaler()
    X = scaler.fit_transform(X)

    # Train-test split
    X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42, stratify=y)

    print("Training shape:", X_train.shape)
    print("Testing shape:", X_test.shape)

    model = XGBClassifier(use_label_encoder=False,

```

```
eval_metric='logloss', random_state=42)
model.fit(X_train, y_train)
```

```
Shape: (400, 26)
```

```
Training shape: (320, 24)
```

```
Testing shape: (80, 24)
```

```
/usr/local/lib/python3.12/dist-packages/xgboost/training.py:183:
UserWarning: [14:18:01] WARNING: /workspace/src/learner.cc:738:
Parameters: { "use_label_encoder" } are not used.
```

```
bst.update(dtrain, iteration=i, fobj=obj)
```

```
# □ Histogram of Age
```

```
plt.figure(figsize=(8,5))
sns.histplot(data=df, x='age', bins=20, kde=True,
hue='classification')
plt.title("Age Distribution by CKD Status")
plt.show()
```

```
# □ Histogram of Blood Pressure
```

```
plt.figure(figsize=(8,5))
sns.histplot(data=df, x='bp', bins=20, kde=True, hue='classification')
plt.title("Blood Pressure Distribution by CKD Status")
plt.show()
```

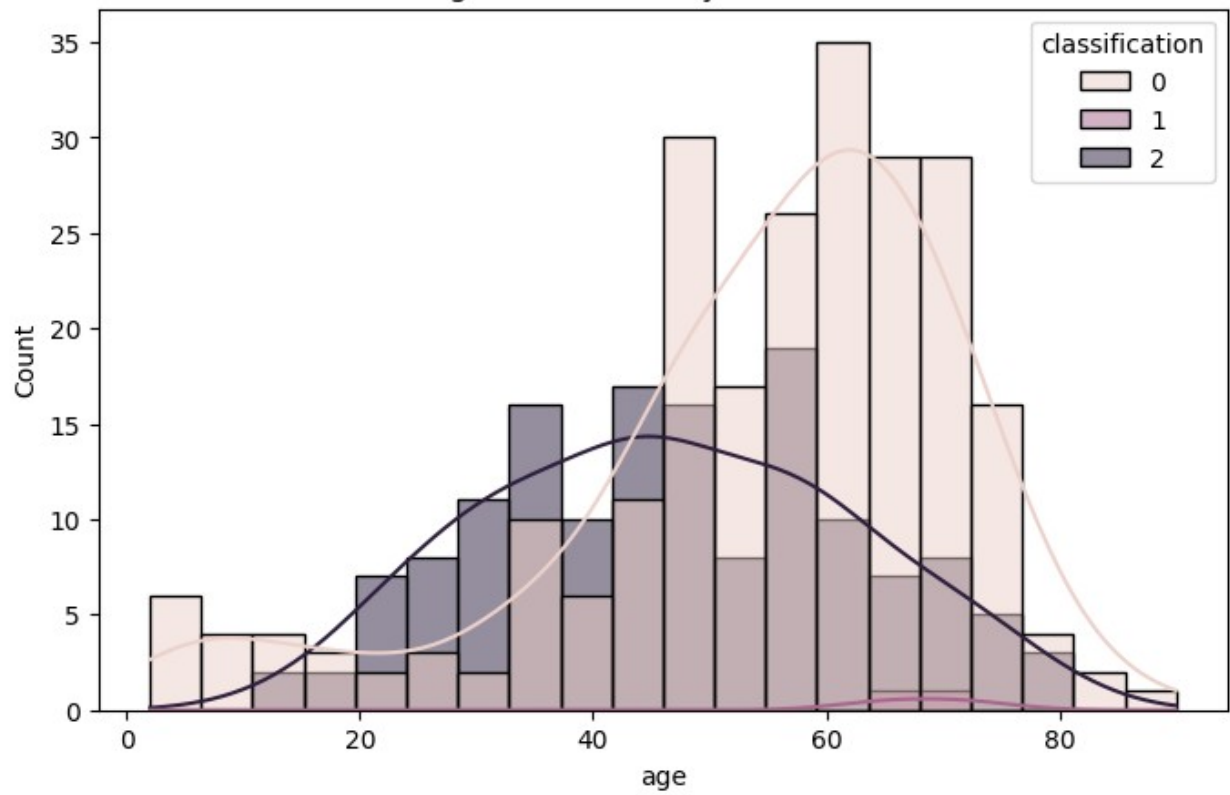
```
# □ Scatter Plot: Age vs Blood Pressure
```

```
plt.figure(figsize=(8,5))
sns.scatterplot(x=df['age'], y=df['bp'], hue=df['classification'],
alpha=0.7)
plt.title("Age vs Blood Pressure (Colored by CKD Status)")
plt.show()
```

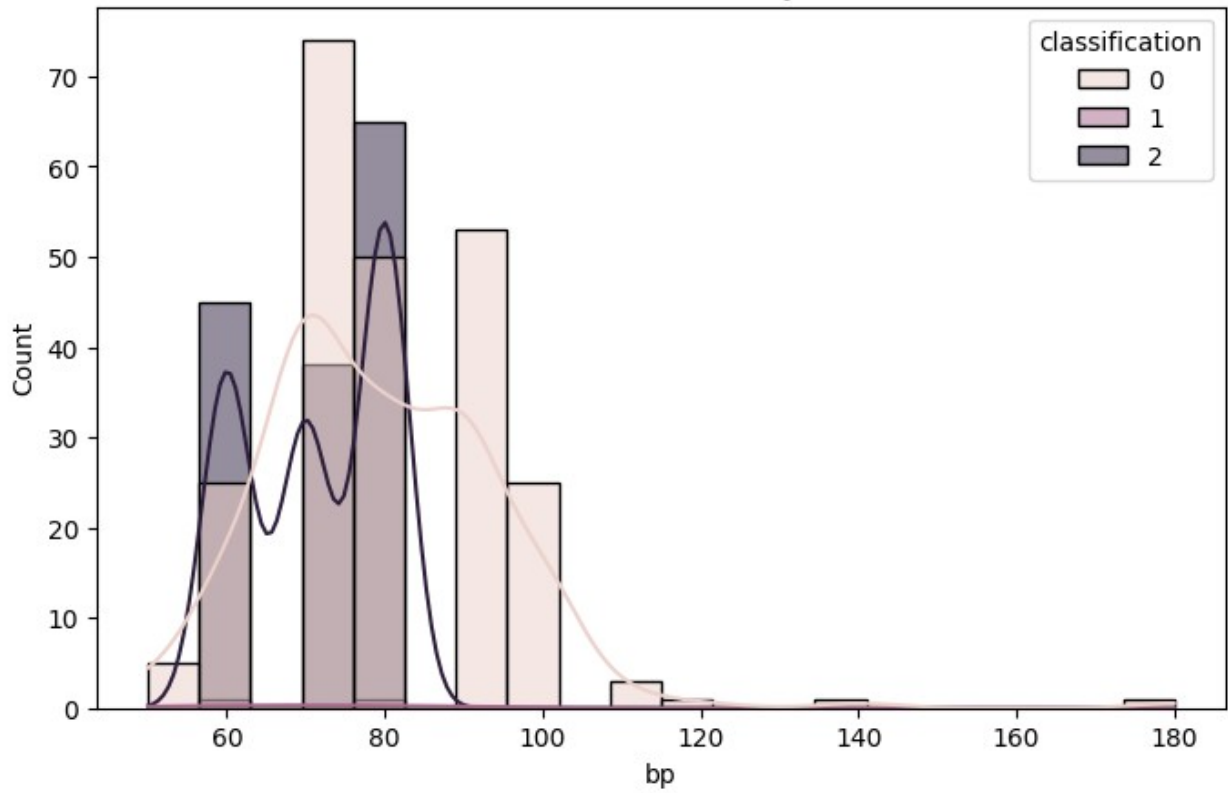
```
# □ Pairplot for selected features
```

```
selected_features = ['age', 'bp', 'sg', 'al', 'su', 'classification']
sns.pairplot(df[selected_features], hue='classification')
plt.suptitle("Pairplot of Selected Features", y=1.02)
plt.show()
```

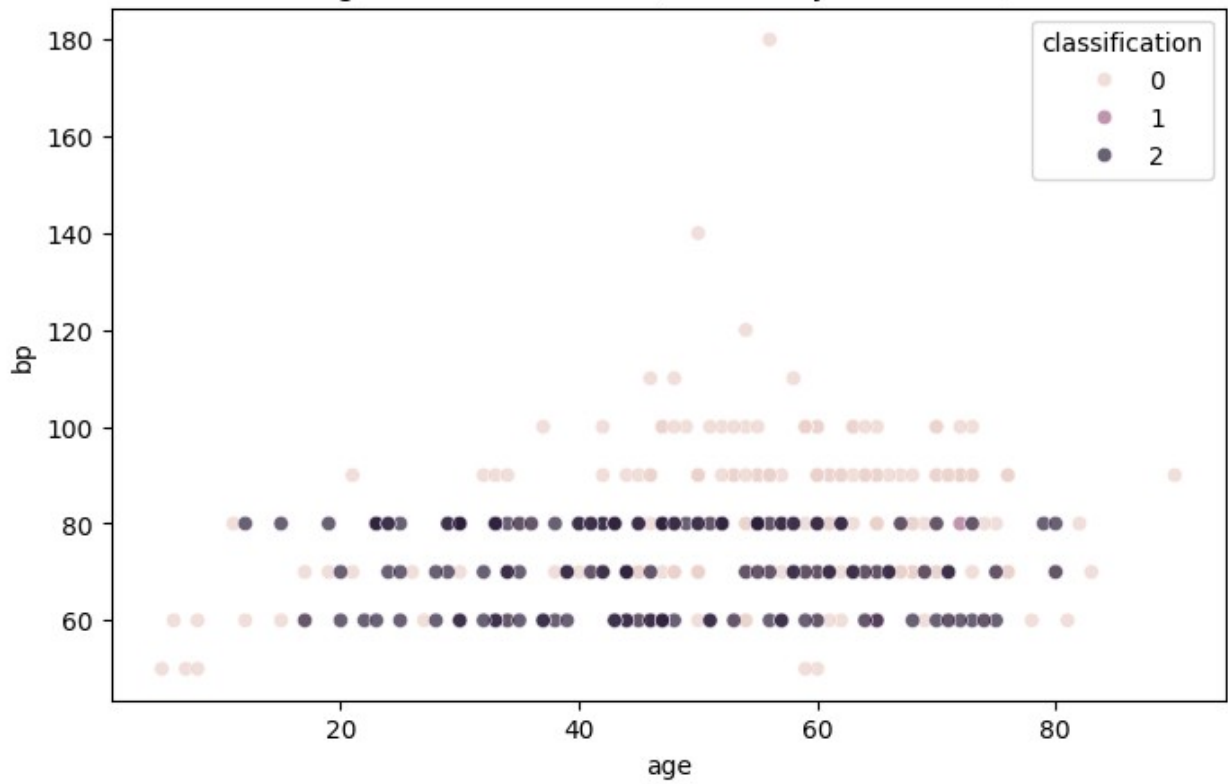
Age Distribution by CKD Status



Blood Pressure Distribution by CKD Status

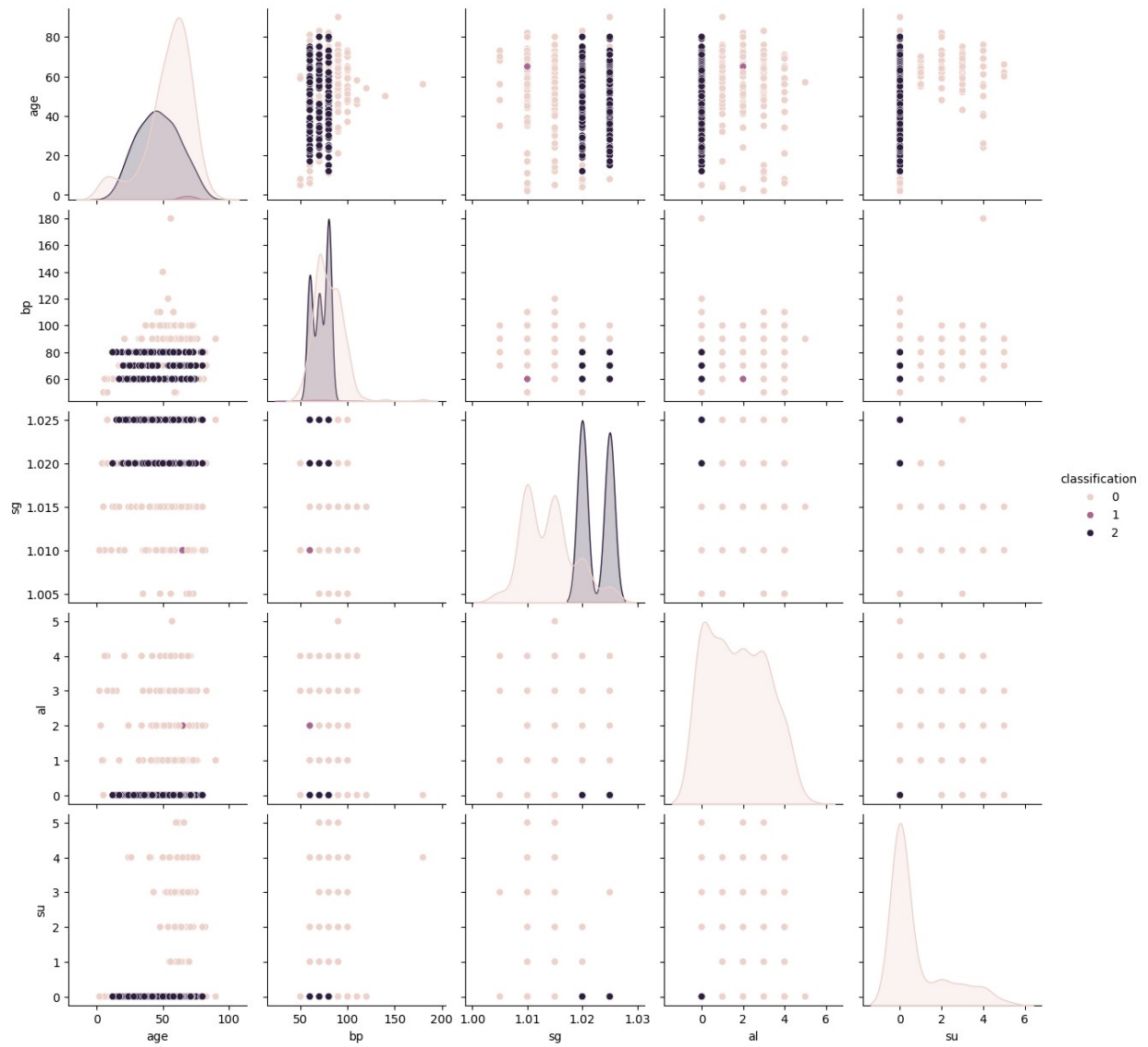


Age vs Blood Pressure (Colored by CKD Status)





Pairplot of Selected Features



```
import shap
explainer = shap.TreeExplainer(model)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values, X_test)
```

