```
TEST
In [ ]:
from google.colab import drive
drive.mount('/content/drive')
Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount
("/content/drive", force remount=True).
In [ ]:
import pandas as pd
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
# Load the dataset
file path = '/content/drug200.csv'
data = pd.read csv(file path)
# Define categorical and numerical features based on actual data columns
categorical_features = ['Sex', 'BP', 'Cholesterol', 'Drug']
numerical features = ['Age', 'Na to K']
# Create transformers for the pipeline
categorical transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='most frequent')),
    ('onehot', OneHotEncoder(handle unknown='ignore'))
1)
numerical transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='mean')),
    ('scaler', StandardScaler())
])
# Combine transformers into a preprocessor with ColumnTransformer
preprocessor = ColumnTransformer(
    transformers=[
       ('num', numerical_transformer, numerical_features),
        ('cat', categorical transformer, categorical_features)
    ])
# Apply transformations
data preprocessed = preprocessor.fit transform(data)
# Convert preprocessed data back to a DataFrame (optional, for visualization)
columns transformed = preprocessor.named transformers ['cat'].named steps['onehot'].get f
eature names out (categorical features)
new columns = list(numerical features) + list(columns transformed)
data preprocessed df = pd.DataFrame(data preprocessed, columns=new columns)
# Print the first few rows of the preprocessed data
print(data preprocessed df.head())
            Na to K Sex F Sex M BP HIGH BP LOW BP NORMAL
        Age
0 -1.291591 1.286522
                      1.0
                             0.0
                                        1.0
                                                0.0
                                                            0.0
                        0.0
 0.162699 -0.415145
                                        0.0
                                                            0.0
                               1.0
                                                1.0
 0.162699 -0.828558
                        0.0
                               1.0
                                        0.0
                                                1.0
                                                            0.0
3 -0.988614 -1.149963
                        1.0
                              0.0
                                        0.0
                                                0.0
                                                            1.0
4 1.011034 0.271794
                       1.0
                             0.0
                                        0.0
                                                1.0
                                                            0.0
   Cholesterol HIGH Cholesterol NORMAL Drug drugA Drug drugB Drug drugC
0
                1.0
                                                0.0
                                    0.0
                                                            0.0
                                                                        0.0
1
                1.0
                                    0.0
                                                0.0
                                                            0.0
                                                                        1.0
```

0.0

0.0

0.0

1.0

2

1.0

```
3
                 ⊥.∪
                                        \cup . \cup
                                                     \cup . \cup
                                                                  \cup . \cup
                                                                                \cup . \cup
4
                 1.0
                                        0.0
                                                     0.0
                                                                  0.0
                                                                                0.0
   Drug_drugX Drug_drugY
0
          0.0
           0.0
                        0.0
2
                        0.0
           0.0
3
           1.0
                        0.0
4
           0.0
                        1.0
In [ ]:
# Load the dataset
data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 6 columns):
                  Non-Null Count Dtype
 # Column
___
 0
                   200 non-null int64
   Age
 1 Sex
                   200 non-null object
                   200 non-null
                                    object
 3
    Cholesterol 200 non-null
                                    object
   Na to K
                   200 non-null
                                    float64
 5
    Drug
                   200 non-null
                                    object
dtypes: float64(1), int64(1), object(4)
memory usage: 9.5+ KB
In [ ]:
data preprocessed df.head()
Out[]:
          Na_to_K Sex_F Sex_M BP_HIGH BP_LOW BP_NORMAL Cholesterol_HIGH Cholesterol_NORMAL Drug_drugA
           1.286522
                     1.0
                           0.0
                                            0.0
                                                       0.0
                                                                      1.0
                                                                                        0.0
                                                                                                  0.0
                                    1.0
   1.291591
1 0.162699
                     0.0
                           1.0
                                    0.0
                                            1.0
                                                       0.0
                                                                      1.0
                                                                                        0.0
                                                                                                  0.0
           0.415145
2 0.162699
                     0.0
                           1.0
                                    0.0
                                            1.0
                                                       0.0
                                                                      1.0
                                                                                        0.0
                                                                                                  0.0
           0.828558
                     1.0
                           0.0
                                    0.0
                                            0.0
                                                       1.0
                                                                      1.0
                                                                                        0.0
                                                                                                  0.0
  0.988614 1.149963
4 1.011034 0.271794
                     1.0
                           0.0
                                    0.0
                                            1.0
                                                       0.0
                                                                      1.0
                                                                                        0.0
                                                                                                  0.0
In [ ]:
from sklearn.model selection import train test split
from sklearn.preprocessing import LabelEncoder
# Define features and target
X = data[['Age', 'Sex', 'BP', 'Cholesterol', 'Na to K']]
y = data['Drug']
# Encode categorical features
X = pd.get dummies(X, columns=['Sex', 'BP', 'Cholesterol'])
# Encode the labels
label encoder = LabelEncoder()
y_encoded = label_encoder.fit_transform(y)
```

X_train, X_test, y_train, y_test = train_test_split(X_scaled, y encoded, test size=0.2,

Scale features

Split data

scaler = StandardScaler()

X scaled = scaler.fit transform(X)

```
# Check the shapes of the splits
print("Training features shape:", X_train.shape)
print("Testing features shape:", X_test.shape)
print("Training labels shape:", y_train.shape)
print("Testing labels shape:", y_test.shape)
```

```
Training features shape: (160, 9)
Testing features shape: (40, 9)
Training labels shape: (160,)
Testing labels shape: (40,)
```

Model Implementation

svm_with_smo

```
In [ ]:
```

```
import numpy as np
# this library is used to show the progress of the training only
from tqdm import tqdm
def svm with smo(X, y, C, eps, max iter, kernel):
   n samples, n features = X.shape
    # Initialize alpha to 0
   alpha = np.zeros(n samples)
   # Initialize b to 0
   b = 0
    # Pre-compute the kernel matrix
   K = np.array([[kernel(X[i], X[j]) for j in range(n_samples)] for i in range(n sample
s)])
   # Start iterations
   k = 0
   while k < max iter:</pre>
       num alpha changed = 0
       for i in range(n samples):
            Ei = b + np.sum(alpha * y * K[i]) - y[i]
            if (y[i] * Ei < -eps and alpha[i] < C) or (y[i] * Ei > eps and alpha[i] > 0
):
                j = np.random.randint(0, n_samples)
                while j == i:
                    j = np.random.randint(0, n_samples)
                Ej = b + np.sum(alpha * y * K[j]) - y[j]
                alpha old i = alpha[i]
                alpha old j = alpha[j]
                if y[i] != y[j]:
                    L = max(0, alpha[j] - alpha[i])
                    H = min(C, C + alpha[j] - alpha[i])
                else:
                   L = max(0, alpha[i] + alpha[j] - C)
                    H = min(C, alpha[i] + alpha[j])
                if L == H:
                    continue
                eta = 2 * K[i, j] - K[i, i] - K[j, j]
                if eta >= 0:
                    continue
                alpha[j] -= y[j] * (Ei - Ej) / eta
```

```
alpha[j] = np.clip(alpha[j], L, H)
            if abs(alpha[j] - alpha_old j) < 0.00001:
                continue
            alpha[i] += y[i] * y[j] * (alpha old j - alpha[j])
            b1 new term = y[i] * (alpha[i] - alpha old i) * K[i, i]
            b2 new term = y[j] * (alpha[j] - alpha old j) * K[j, j]
            b1 = b - Ei - b1_new_term - b2_new_term
            b2 = b - Ej - b1 new term - b2 new term
            if 0 < alpha[i] < C:</pre>
               b = b1
            elif 0 < alpha[j] < C:</pre>
               b = b2
            else:
                b = (b1 + b2) / 2
            num alpha changed += 1
    if num alpha changed == 0:
        k += 1
    else:
       k = 0
# Compute the weights
w = np.sum(alpha * y * X.T, axis=1)
return w, b
```

One vs all classification

return W, b

In []:

```
def one versus all(X, y, C, eps, kernel, max iter):
   Implementation of the one-versus-all strategy for multi-class classification.
   Args:
       train X: The training data, shape (n samples, n features).
        train_y: The training labels, shape (n_samples,).
       C: The regularization strength.
        eps: The tolerance for stopping criterion.
       max_iter: The maximum number of iterations.
   Returns:
       W: The weights, shape (n classes, n features).
       b: The bias terms, shape (n classes,).
   n samples, n features = X.shape
   n classes = len(np.unique(y))
    # Initialize W and b
   W = np.zeros((n classes, n features))
   b = np.zeros(n classes)
    # Train a binary classifier for each class
   for i in tqdm(range(n_classes)):
        # Create a copy of the labels
       y_{train} = np.copy(y)
        # Set all the labels to -1
       y_{train}[y_{train} != i] = -1
        # Set the labels of the current class to 1
       y_train[y_train == i] = 1
        # Train the binary classifier
        W[i], b[i] = svm_with_smo(X, y_train, C, eps, max_iter, kernel)
```

One versus one classification

```
In [ ]:
```

```
def one versus one (X, y, C, eps, kernel, max iter):
   Implementation of the one-versus-one strategy for multi-class classification.
   Args:
       X: The training data, shape (n_samples, n_features).
       y: The training labels, shape (n_samples,).
       C: The regularization strength.
       eps: The tolerance for stopping criterion.
       max iter: The maximum number of iterations.
       kernel: The kernel function.
   Returns:
       W: The weights, shape (n classes * (n classes - 1) / 2, n features).
       b: The bias terms, shape (n classes * (n classes - 1) / 2,).
   n samples, n features = X.shape
   classes = np.unique(y)
   n classes = len(classes)
    # Initialize W and b
   W = []
   b = []
    # Train a binary classifier for each pair of classes
   for i in tqdm(range(n_classes)):
        for j in range(i + 1, n classes):
            # Create a binary label array
           binary y = np.where((y == classes[i]) | (y == classes[j]), y, 0)
            binary y = np.where(binary y == classes[i], -1, binary y)
            binary y = np.where(binary y == classes[j], 1, binary y)
            # Exclude samples that don't belong to the i-th or j-th class
            binary y nonzero = binary y[binary y != 0]
            X nonzero = X[binary y != 0]
            # Train the binary classifier
            w, b value = svm with smo(X nonzero, binary y nonzero, C, eps, max iter, ker
nel)
            W.append(w)
            b.append(b_value)
   return np.array(W), np.array(b)
```

Predict functions

```
In [ ]:
```

```
def predict_one_versus_all(X, W, b):
    # Calculate the scores for each classifier
    scores = np.dot(X, W.T) + b
    # The predicted class is the one with the highest score
    return np.argmax(scores, axis=1)
def predict_one_versus_one(X, W, b, classes):
    n_samples = X.shape[0]
    votes = np.zeros((n_samples, len(classes)))

# Iterate over each classifier
    k = 0
    for i in range(len(classes)):
        for j in range(i + 1, len(classes)):
            scores = np.dot(X, W[k]) + b[k]
            predictions = np.where(scores > 0, classes[j], classes[i])
```

```
for p in range(n_samples):
     votes[p, predictions[p]] += 1
     k += 1

# The predicted class is the one with the highest number of votes
return np.argmax(votes, axis=1)
```

Training in drug200 dataset

```
In [ ]:
```

```
import pandas as pd
import numpy as np
from sklearn.model selection import train test split
from sklearn.preprocessing import StandardScaler, LabelEncoder
# Load the dataset
file path = '/content/drug200.csv'
data = pd.read csv(file path)
# Define features and target
X = data[['Age', 'Sex', 'BP', 'Cholesterol', 'Na to K']]
y = data['Drug']
# Encode categorical features
X = pd.get dummies(X, columns=['Sex', 'BP', 'Cholesterol'])
# Encode the labels
label encoder = LabelEncoder()
y encoded = label encoder.fit transform(y)
# Scale features
scaler = StandardScaler()
X scaled = scaler.fit_transform(X)
# Split data
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y_encoded, test size=0.2,
random state=42)
# Define the kernel function
def kernel linear(x1, x2):
   return np.dot(x1, x2)
```

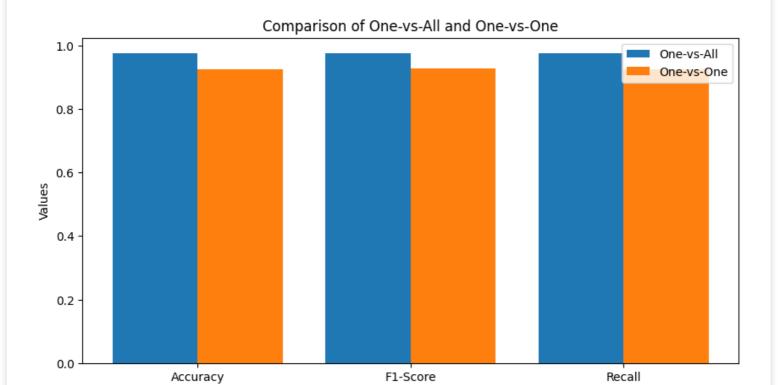
In []:

```
from imblearn.over sampling import SMOTE
from sklearn.metrics import confusion_matrix, classification_report
from sklearn.metrics import f1_score, recall_score
import matplotlib.pyplot as plt
# Apply SMOTE to balance the dataset
smote = SMOTE(random state=42)
X resampled, y resampled = smote.fit resample(X train, y train)
# Train using one-versus-all on resampled data
W one all, b one all = one versus all(X resampled, Y resampled, C=1.0, eps=0.001, kernel
=kernel_linear, max_iter=1000)
# Train using one-versus-one on resampled data
W one one, b one one = one versus one(X resampled, y resampled, C=1.0, eps=0.001, kernel
=kernel linear, max iter=1000)
# Predict using one-versus-all
y_pred_one_all = predict_one_versus_all(X_test, W_one_all, b_one_all)
accuracy_one_all = np.mean(y_pred_one_all == y_test)
f1_one_all = f1_score(y_test, y_pred_one_all, average='weighted')
recall one all = recall score(y test, y pred one all, average='weighted')
print("One-vs-All Prediction Accuracy:", accuracy one all)
```

```
print("One-vs-All Prediction F1-Score:", f1 one all)
print("One-vs-All Prediction Recall:", recall_one_all)
print("One-vs-All Confusion Matrix:\n", confusion matrix(y test, y pred one all))
print("One-vs-All Classification Report:\n", classification report(y test, y pred one all
# Predict using one-versus-one
y pred one one = predict one versus one(X test, W one one, b one one, np.unique(y train)
accuracy one one = np.mean(y pred one one == y test)
f1 one one = f1 score(y test, y pred one one, average='weighted')
recall one one = recall score(y test, y pred one one, average='weighted')
print("One-vs-One Prediction Accuracy:", accuracy one one)
print("One-vs-One Prediction F1-Score:", f1 one one)
print("One-vs-One Prediction Recall:", recall one one)
print("One-vs-One Confusion Matrix:\n", confusion matrix(y test, y pred one one))
print("One-vs-One Classification Report:\n", classification_report(y_test, y_pred_one_one
) )
# Plotting metrics
metrics = ['Accuracy', 'F1-Score', 'Recall']
one_all_values = [accuracy_one_all, f1_one_all, recall_one_all]
one one values = [accuracy one one, f1 one one, recall one one]
x = np.arange(len(metrics))
plt.figure(figsize=(10, 5))
plt.bar(x - 0.2, one all values, 0.4, label='One-vs-All')
plt.bar(x + 0.2, one one values, 0.4, label='One-vs-One')
plt.xlabel('Metrics')
plt.ylabel('Values')
plt.title('Comparison of One-vs-All and One-vs-One')
plt.xticks(x, metrics)
plt.legend()
plt.show()
          | 5/5 [11:47<00:00, 141.47s/it]
100%|
              | 5/5 [06:42<00:00, 80.55s/it]
One-vs-All Prediction Accuracy: 0.975
One-vs-All Prediction F1-Score: 0.9755305039787799
One-vs-All Prediction Recall: 0.975
One-vs-All Confusion Matrix:
 [[ 0 0 0 0 0]]
 [ 0 3 0 0 0]
 [00500]
 [ 0 0 0 11 0]
 [ 1 0 0 0 14]]
One-vs-All Classification Report:
             precision recall f1-score support
                                    0.92
          0
                 0.86 1.00
                                                  6
                                                 3
          1
                  1.00
                          1.00
                                    1.00
          2
                  1.00
                          1.00
                                    1.00
                                                 5
                          1.00
          3
                 1.00
                                    1.00
                                                 11
                 1.00
                          0.93
                                    0.97
                                                 15
                                     0.97
                                                 40
   accuracy
                 0.97
                          0.99
                                    0.98
  macro avq
                                                 40
                 0.98
                            0.97
                                     0.98
                                                 40
weighted avg
One-vs-One Prediction Accuracy: 0.925
One-vs-One Prediction F1-Score: 0.9286440570923331
One-vs-One Prediction Recall: 0.925
One-vs-One Confusion Matrix:
 [[5 1 0 0 0]
 [ 0 3 0 0 0]
     0 5 0
 [ 0
              0]
 [ 1 0 0 10 0]
 [ 1 0 0 0 14]]
```

--- One Oleveitication Demant.

one-vs-one	е ста	SSIIICALION	керогі:		
		precision	recall	f1-score	support
	0	0.71	0.83	0.77	6
	1	0.75	1.00	0.86	3
	2	1.00	1.00	1.00	5
	3	1.00	0.91	0.95	11
	4	1.00	0.93	0.97	15
200117	0.017			0.93	40
accura	_	0 00	0 04		
macro a	avg	0.89	0.94	0.91	40
weighted a	avg	0.94	0.93	0.93	40



Metrics