

## ▼ CSE 572: Lab 11

In this lab, you will practice implementing ensemble models.

To execute and make changes to this notebook, click File > Save a copy to save your own version in your Google Drive or Github. Read the step-by-step instructions below carefully. To execute the code, click on each cell below and press the SHIFT-ENTER keys simultaneously or by clicking the Play button.

When you finish executing all code/exercises, save your notebook then download a copy (.ipynb file). Submit the following **three** things:

1. a link to your Colab notebook,
2. the .ipynb file, and
3. a pdf of the executed notebook on Canvas.

To generate a pdf of the notebook, click File > Print > Save as PDF.

```
# Import libraries
import numpy as np
import pandas as pd

# Set the random seed for reproducibility
seed = 0
np.random.seed(0)
```

## ▼ Ensemble of hybrid models

One straightforward approach for constructing an ensemble classifier is to train  $k$  separate classifiers using different classification methods and then combine their predictions using majority vote. In the first exercise, you will use Scikit-learn to train a  $k$  nearest neighbors, naive Bayes, and logistic regression classifier separately and then combine their predictions using a VotingClassifier.

## ▼ Load the dataset

We will use the Wisconsin breast cancer dataset with class values 'benign' or 'malignant'. Below, we will load the dataset and perform preprocessing as in previous labs.

```
# Load the Wisconsin breast cancer dataset
data = pd.read_csv('https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/brea
data.columns = ['Sample code', 'Clump Thickness', 'Uniformity of Cell Size', 'Uniformity of Cell Shape',
                'Marginal Adhesion', 'Single Epithelial Cell Size', 'Bare Nuclei', 'Bland Chromatin',
                'Normal Nucleoli', 'Mitoses', 'Class']

data = data.drop(['Sample code'],axis=1)

data = data.replace('?',np.NaN)
data['Bare Nuclei'] = pd.to_numeric(data['Bare Nuclei'])

data
```



	Clump Thickness	Uniformity of Cell Size	Uniformity of Cell Shape	Marginal Adhesion	Single Epithelial Cell Size	Bare Nuclei	Bland Chromatin	No Nucl
0	5	1	1	1	2	1.0	3	
1	5	4	4	5	7	10.0	3	
2	3	1	1	1	2	2.0	3	
3	6	8	8	1	3	4.0	3	
4	4	1	1	3	2	1.0	3	
...	...	...	...	...	...	...	...	
694	3	1	1	1	3	2.0	1	
695	2	1	1	1	2	1.0	1	
696	5	10	10	3	7	3.0	8	
697	4	8	6	4	3	4.0	10	
698	4	8	8	5	4	5.0	10	

After loading the dataset, we clean it by removing samples with missing data, duplicates, or outliers using the code from Labs 2-3.



```
def inds_nans(df):
    inds = df.isna().any(axis=1)
    # print('Found {} rows that had NaN values.'.format(inds.sum()))
    return inds

def inds_dups(df):
    inds = df.duplicated()
    # print('Found {} rows that were duplicates.'.format(inds.sum()))
    return inds

def inds_outliers(df):
    # In this example, we defined outliers as values that are +/- 3 standard deviations
    # from the mean value. To identify such values, we need to compute the Z score for
    # every value by subtracting the feature-wise mean and dividing by the feature-wise
    # standard deviation (also known as standardizing the data).
    df = df[df.columns[:-1]]
    Z = (df-df.mean())/df.std()
    # The below code will give a value of True or False for each row. The row will be
    # True if all of the feature values for that row were within 3 standard deviations of
    # the mean. The row will be False if at least one of the feature values for that row
    # was NOT within 3 standard deviations of the mean.
    inlier_inds = ((Z > -3).sum(axis=1)==9) & ((Z <= 3).sum(axis=1)==9)
    # The outliers are the inverse boolean values of the above
    outlier_inds = ~inlier_inds
    # print('Found {} rows that were outliers.'.format(outlier_inds.sum()))
    return outlier_inds

# Select only the rows at index locations that were not nans, duplicates, or outliers
data_clean = data.loc[~((inds_nans(data) | inds_dups(data)) | inds_outliers(data)),:]

data_clean
```

	Clump Thickness	Uniformity of Cell Size	Uniformity of Cell Shape	Marginal Adhesion	Single Epithelial Cell Size	Bare Nuclei	Bland Chromatin	Ni
0	5	1	1	1	2	1.0	3	
1	5	4	4	5	7	10.0	3	
2	3	1	1	1	2	2.0	3	
3	6	8	8	1	3	4.0	3	
4	4	1	1	3	2	1.0	3	
...	...	...	...	...	...	...	...	
693	3	1	1	1	2	1.0	2	
694	3	1	1	1	3	2.0	1	
696	5	10	10	3	7	3.0	8	
697	4	8	6	4	3	4.0	10	
698	4	8	8	5	4	5.0	10	

Next we normalize the data using the code from Lab 3 so the features will have approximately normal distributions.



```
from sklearn import preprocessing
```

```
# Normalize the feature columns
```

```
data_clean[data_clean.columns[:-1]] = preprocessing.normalize(data_clean[data_clean.columns[:-1]], norm='l
```

```
/usr/local/lib/python3.8/dist-packages/pandas/core/frame.py:3678: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexi
self[col] = igetitem(value, i)
```



```
data_clean
```

Split the data into a training and test set with 70% train and 30% test. Use the `seed` variable to set the random state.

```
from sklearn.model_selection import train_test_split
```

```
data_new = data_clean[data_clean.columns[:-1]]
```

```
target = data_clean[data_clean.columns[-1]]
```

```
# YOUR CODE HERE
```

```
X_train, X_test, y_train, y_test = train_test_split(data_new, target, test_size=0.3, shuffle=True, random_
```

## ▼ Train a Gaussian Naive Bayes classifier

Use the `GaussianNB` object in `sklearn` to fit a Gaussian Naive Bayes classifier and predict the class labels for the test set based on probabilities estimated from the training set.

```
from sklearn.naive_bayes import GaussianNB
```

```
gnb = GaussianNB()

# Fit the model parameters using the training data
gnb = gnb.fit(X_train, y_train)

# Predict the test set classes using the trained model
y_pred_gnb = gnb.predict(X_test)
```

Compute the accuracy of this model on the test set.

```
from sklearn.metrics import accuracy_score
# YOUR CODE HERE

print('Test data accuracy: {}'.format((accuracy_score(y_test, y_pred_gnb))))

Test data accuracy: 0.8666666666666667
```

## ▼ Train a Logistic Regression classifier

Use the LogisticRegression class in sklearn to fit a Logistic Regression classifier and predict the class labels for the test set.

```
from sklearn.linear_model import LogisticRegression

# Instantiate a logistic regression classifier and fit it to the training data
lr = LogisticRegression(random_state=seed)
lr = lr.fit(X_train, y_train)

# Predict the test set classes using the trained model

y_pred_lr = lr.predict(X_test)
```

Compute the accuracy of this model on the test set.

```
# YOUR CODE HERE

print('Test data accuracy: {}'.format((accuracy_score(y_test, y_pred_lr))))

Test data accuracy: 0.8166666666666667
```

## ▼ Train a k Nearest Neighbors classifier

Use the KNeighborsClassifier class in sklearn to train a kNN classifier and predict the class labels for the test set. We will use Euclidean distance with  $k = 5$ .

```
from sklearn.neighbors import KNeighborsClassifier

knn = KNeighborsClassifier(n_neighbors=5, metric='euclidean')
knn.fit(X_train, y_train)

KNeighborsClassifier(metric='euclidean')

# Predict the test set classes using the trained model

y_pred_knn = knn.predict(X_test) # YOUR CODE HERE
```

Compute the accuracy of this model on the test set.

```
# YOUR CODE HERE
```

```
print('Test data accuracy: {}'.format((accuracy_score(y_test, y_pred_knn))))
```

```
Test data accuracy: 0.8333333333333334
```

**Question 1: What is the test accuracy for each of the 3 models (rounded to 2 decimal places)?**

**Answer:**

YOUR ANSWER HERE

Gaussian NB : 0.87

Logistic: 0.82

KNN: 0.83

### ▼ Creating an ensemble VotingClassifier

We will use the [VotingClassifier](#) class in sklearn to combine the predictions of these 3 models using majority vote.

```
from sklearn.ensemble import VotingClassifier
```

```
ensemble = VotingClassifier(estimators=[('gnb', gnb), ('lr', lr), ('knn', knn)], voting='soft')
```

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

```
VotingClassifier(estimators=[('gnb', GaussianNB()),
                             ('lr', LogisticRegression(random_state=0)),
                             ('knn', KNeighborsClassifier(metric='euclidean'))],
                  voting='soft')
```

**Question 2: What is the effect of setting the voting=soft vs. voting=hard parameter in the VotingClassifier?**

**You should consult the documentation to answer this question.**

**Answer:**

YOUR ANSWER HERE

Voting = hard, the predicted class label for a particular sample is the class label that represents the majority (mode) of the class labels predicted by each individual classifier. It uses the binary class labels for voting.

Voting = Soft, returns the class label as argmax of the sum of predicted probabilities. It uses the class probabilities output by each model for voting.

```
for clf, label in zip([gnb, lr, knn, ensemble], ['Naive Bayes', 'Logistic Regression', 'k Nearest Neighbor',
                                                    'Ensemble']):
    score = accuracy_score(clf.predict(X_test), y_test)
    print("Accuracy: %0.2f [%s]" % (score, label))
```

```
Accuracy: 0.87 [Naive Bayes]
Accuracy: 0.82 [Logistic Regression]
Accuracy: 0.83 [k Nearest Neighbors]
Accuracy: 0.88 [Ensemble]
```

## ▼ Ensemble using boosting

Another method for ensembling classifiers to improve the performance over any single classifier is using a technique called boosting. Boosting uses an iterative procedure to adaptively change the distribution of the training data by focusing more on previously misclassified samples each time a new classifier is trained. In the second exercise, you will implement the [AdaBoost](#) algorithm using a decision tree as the base classifier (the sklearn implementation uses a decision tree by default).

You will use the same Wisconsin breast cancer dataset as in the previous exercise.

## ▼ Train AdaBoost with decision tree

Use the AdaBoost algorithm to train an ensemble of decision trees. Use 50 trees for the ensemble.

```
from sklearn.ensemble import AdaBoostClassifier

ada = AdaBoostClassifier(n_estimators=50, random_state=seed)
# Train the model
ada.fit(X_train, y_train)

AdaBoostClassifier(random_state=0)

# Predict the test set classes using the trained model

y_pred_ada = ada.predict(X_test) # YOUR CODE HERE

Compute the accuracy of this model on the test set.

# YOUR CODE HERE

print('Test data accuracy: {}'.format((accuracy_score(y_test, y_pred_ada))))

Test data accuracy: 0.8916666666666667
```

✓ 0s completed at 6:44 AM

