

# Supervised Learning with scikit-learn

## 3). Finetuning your model

### a). Metrics for Classification

# Import necessary modules

```
from sklearn.metrics import classification_report
```

```
from sklearn.metrics import confusion_matrix
```

# Create training and test set

```
X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.4, random_state=42)
```

# Instantiate a k-NN classifier: knn

```
knn = KNeighborsClassifier(n_neighbors=6)
```

# Fit the classifier to the training data

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

# Predict the labels of the test data: y\_pred

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

# Generate the confusion matrix and classification report

```
print(confusion_matrix(y_test, y_pred))
```

```
print(classification_report(y_test, y_pred))
```

<script.py> output:

```
[[176 30]
```

```
[ 52 50]]
```

```
precision recall f1-score support
```

```
0    0.77    0.85    0.81    206
```

```
1    0.62    0.49    0.55    102
```

```
avg / total    0.72    0.73    0.72    308
```

**c). Building a logisticregression Model****# Import the necessary modules****from sklearn.linear\_model import LogisticRegression****from sklearn.metrics import confusion\_matrix, classification\_report****# Create training and test sets****X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size = 0.4, random\_state=42)****# Create the classifier: logreg****logreg = LogisticRegression()****# Fit the classifier to the training data****logreg.fit(X\_train,y\_train)****# Predict the labels of the test set: y\_pred****y\_pred = logreg.predict(X\_test)****# Compute and print the confusion matrix and classification report****print(confusion\_matrix(y\_test, y\_pred))****print(classification\_report(y\_test, y\_pred))**

&lt;script.py&gt; output:

[[176 30]

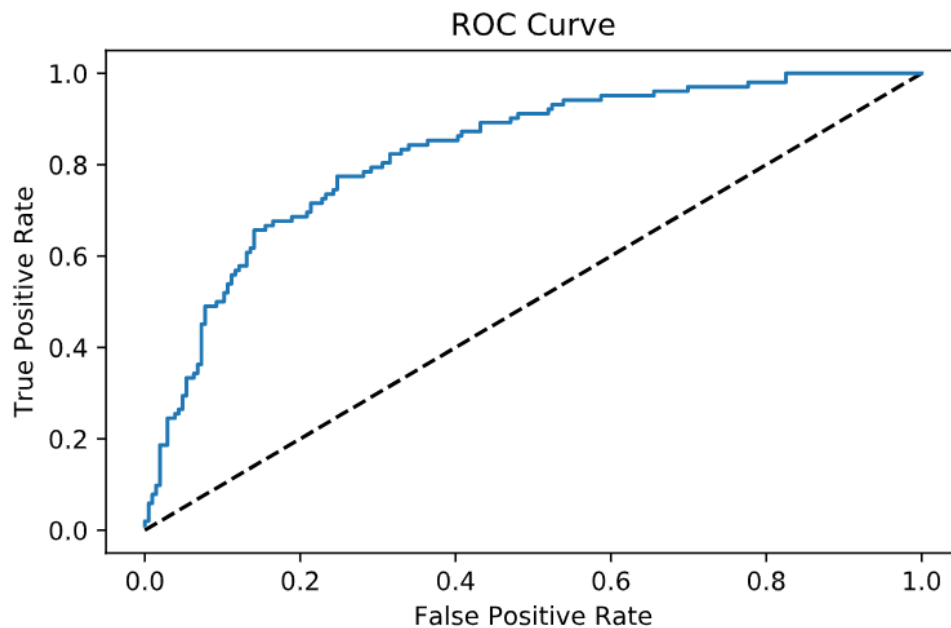
[ 35 67]]

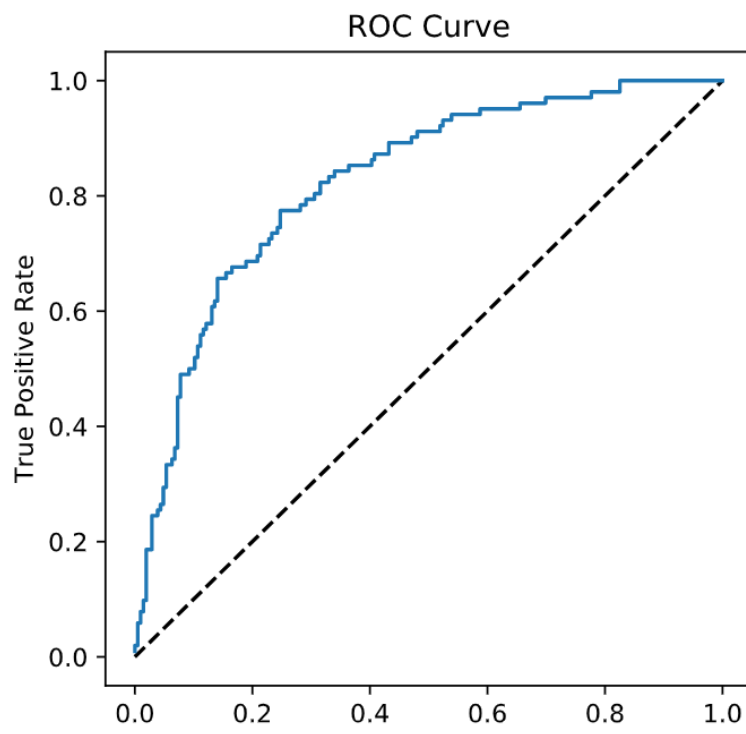
	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.83	0.85	0.84	206
---	------	------	------	-----

1	0.69	0.66	0.67	102
---	------	------	------	-----

avg / total	0.79	0.79	0.79	308
-------------	------	------	------	-----

**d). Plotting an ROC curve****# Import necessary modules****from sklearn.metrics import roc\_curve****# Compute predicted probabilities: y\_pred\_prob****y\_pred\_prob = logreg.predict\_proba(X\_test)[: ,1]****# Generate ROC curve values: fpr, tpr, thresholds****fpr, tpr, threshold = roc\_curve(y\_test, y\_pred\_prob)****# Plot ROC curve****plt.plot([0, 1], [0, 1], 'k--')****plt.plot(fpr, tpr)****plt.xlabel('False Positive Rate')****plt.ylabel('True Positive Rate')****plt.title('ROC Curve')****plt.show()**

**e). AUC Computation****# Import necessary modules****from sklearn.metrics import roc\_auc\_score****from sklearn.model\_selection import cross\_val\_score****# Compute predicted probabilities: y\_pred\_prob****y\_pred\_prob = logreg.predict\_proba(X\_test)[:,1]****# Compute and print AUC score****print("AUC: {}".format(roc\_auc\_score(y\_test, y\_pred\_prob)))****# Compute cross-validated AUC scores: cv\_auc****cv\_auc = cross\_val\_score(logreg, X, y, cv=5, scoring='roc\_auc')****# Print list of AUC scores****print("AUC scores computed using 5-fold cross-validation: {}".format(cv\_auc))**

**f). Hyperparameter tuning with gridsearchCV()****# Import necessary modules****from sklearn.linear\_model import LogisticRegression****from sklearn.model\_selection import GridSearchCV****# Setup the hyperparameter grid****c\_space = np.logspace(-5, 8, 15)****param\_grid = {'C': c\_space}****# Instantiate a logistic regression classifier: logreg****logreg = LogisticRegression()****# Instantiate the GridSearchCV object: logreg\_cv****logreg\_cv = GridSearchCV(logreg, param\_grid, cv=5)****# Fit it to the data****logreg\_cv.fit(X,y)****# Print the tuned parameters and score****print("Tuned Logistic Regression Parameters: {}".format(logreg\_cv.best\_params\_))****print("Best score is {}".format(logreg\_cv.best\_score\_))**

&lt;script.py&gt; output:

Tuned Logistic Regression Parameters: {'C': 3.7275937203149381}

Best score is 0.7708333333333334

**g). Hyperparameter tuning with RandomizedSearchCV****# Import necessary modules****from scipy.stats import randint****from sklearn.tree import DecisionTreeClassifier****from sklearn.model\_selection import RandomizedSearchCV****# Setup the parameters and distributions to sample from: param\_dist****param\_dist = {"max\_depth": [3, None],****"max\_features": randint(1, 9),****"min\_samples\_leaf": randint(1, 9),****"criterion": ["gini", "entropy"]}****# Instantiate a Decision Tree classifier: tree****tree = DecisionTreeClassifier()****# Instantiate the RandomizedSearchCV object: tree\_cv****tree\_cv = RandomizedSearchCV(tree, param\_dist, cv=5)****# Fit it to the data****tree\_cv.fit(X,y)****# Print the tuned parameters and score****print("Tuned Decision Tree Parameters: {}".format(tree\_cv.best\_params\_))****print("Best score is {}".format(tree\_cv.best\_score\_))**

&lt;script.py&gt; output:

Tuned Decision Tree Parameters: {'criterion': 'entropy', 'max\_features': 7, 'min\_samples\_leaf': 8, 'max\_depth': None}

Best score is 0.7200520833333334

**h). Hold-out set in practice I: Classification**

**# Import necessary modules**

**from sklearn.model\_selection import train\_test\_split**

**from sklearn.linear\_model import LogisticRegression**

**from sklearn.model\_selection import GridSearchCV**

**# Create the hyperparameter grid**

**c\_space = np.logspace(-5, 8, 15)**

**param\_grid = {'C': c\_space, 'penalty': ['l1', 'l2']}**

**# Instantiate the logistic regression classifier: logreg**

**logreg = LogisticRegression()**

**# Create train and test sets**

**X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.4, random\_state=42)**

**# Instantiate the GridSearchCV object: logreg\_cv**

**logreg\_cv = GridSearchCV(logreg, param\_grid, cv=5)**

**# Fit it to the training data**

**logreg\_cv.fit(X\_train, y\_train)**

**# Print the optimal parameters and best score**

**print("Tuned Logistic Regression Parameter: {}".format(logreg\_cv.best\_params\_))**

**print("Tuned Logistic Regression Accuracy: {}".format(logreg\_cv.best\_score\_))**

<script.py> output:

Tuned Logistic Regression Parameter: {'C': 0.43939705607607948, 'penalty': 'l1'}

Tuned Logistic Regression Accuracy: 0.7652173913043478



**i). Hold out set in Practice II: Regression**

```
# Import necessary modules
```

```
from sklearn.linear_model import ElasticNet
```

```
from sklearn.metrics import mean_squared_error
```

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

```
from sklearn.model_selection import GridSearchCV
```

```
# Create train and test sets
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4, random_state=42 )
```

```
# Create the hyperparameter grid
```

```
l1_space = np.linspace(0, 1, 30)
```

```
param_grid = {'l1_ratio': l1_space}
```

```
# Instantiate the ElasticNet regressor: elastic_net
```

```
elastic_net = ElasticNet()
```

```
# Setup the GridSearchCV object: gm_cv
```

```
gm_cv = GridSearchCV(elastic_net, param_grid, cv=5)
```

```
# Fit it to the training data
```

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

```
# Predict on the test set and compute metrics
```

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

```
r2 = gm_cv.score(X_test, y_test)
```

```
mse = mean_squared_error(y_test, y_pred)
```

```
print("Tuned ElasticNet l1 ratio: {}".format(gm_cv.best_params_))
```

```
print("Tuned ElasticNet R squared: {}".format(r2))
```

```
print("Tuned ElasticNet MSE: {}".format(mse))
```

<script.py> output:

Tuned ElasticNet l1 ratio: {'l1\_ratio': 0.20689655172413793}

Tuned ElasticNet R squared: 0.8668305372460283

Tuned ElasticNet MSE: 10.05791413339844