

W3.code

110

In []: `from sklearn.model_selection import train_test_split
from sklearn.datasets import load_digits`

`dataset = load_digits()
X, y = dataset.data, dataset.target`

`for class_name, class_count in zip(dataset.target_names, np.bincount(dataset.target)):
 print(class_name, class_count)`

In []: `# Creating a dataset with imbalanced binary classes:
Negative class (0) is 'not digit'
Positive class (1) is 'digit'
y_binary_imbalanced = y.copy() # dump first few entries from original labels + binary label
y_binary_imbalanced[y_binary_imbalanced != 1] = 0`

`print('Original labels:\t', y[1:30])
print('New binary labels:\t', y_binary_imbalanced[1:30])`

In []: `np.bincount(y_binary_imbalanced)` *# Negative class (0) is the most frequent class*

*array([1615, 182])
- + → 9:1* *dataset = class imbalance*

0.1 × 0.05 = 0.005 = 1/20

<https://github.com/Qian-Han/coursera-Applied-Data-Science...pply-Machine-Learning-In-Python/week3/Module%203.ipynb> 2020/4/12 下午11:38 第1页 (共8页)

Dummy Classifiers

DummyClassifier is a classifier that makes predictions using simple rules, which can be useful as a baseline for comparison against actual classifiers, especially with imbalanced classes.

In []: `from sklearn.dummy import DummyClassifier` *we = regular classifier - x look individual data*

Negative class (0) is most frequent just call - x look real data & thumb rule - make prediction

dummy_majority = DummyClassifier(strategy = 'most_frequent').fit(X_train, y_train) *most frequent class*

Therefore the dummy 'most_frequent' classifier always predicts class 0 which class in training set

y_dummy_predictions = dummy_majority.predict(X_test) *most frequent*

y_dummy_predictions always predict 0 / negative class for every instance in test set.

In []: `dummy_majority.score(X_test, y_test)` *get accuracy = 90%, = gbm classifier, = very slightly better > dummy classifier*

In []: `svm = SVC(kernel='linear', C=1).fit(X_train, y_train)` *change kernel param = Rbf → linear*

accuracy = 98%

Confusion matrices

Binary (two-class) confusion matrix

In []: `from sklearn.metrics import confusion_matrix`

most-frequent

Negative class (0) is most frequent

dummy_majority = DummyClassifier(strategy = 'most_frequent').fit(X_train, y_train) *predict P ✓ - during confusion*

y_majority_predicted = dummy_majority.predict(X_test) *✓ predict V*

confusion = confusion_matrix(y_test, y_majority_predicted) *predict P ✓ - during confusion*

print('Most frequent class (dummy classifier)\n', confusion) *random* *TN FP* *FN TP* *Success prediction*

Stratified

In []: `# produces random predictions w/ same class proportion as training set`

dummy_classprop = DummyClassifier(strategy='stratified').fit(X_train, y_train) *1st col 2nd col*

y_classprop_predicted = dummy_classprop.predict(X_test)

confusion = confusion_matrix(y_test, y_classprop_predicted) *✓ occasionally predict P*

print('Random class-proportional prediction (dummy classifier)\n', confusion) *62 correct*

SVC

In []: `svm = SVC(kernel='linear', C=1).fit(X_train, y_train)` *440 correct*

svm_predicted = svm.predict(X_test)

confusion = confusion_matrix(y_test, svm_predicted)

print('Support vector machine classifier (linear kernel, C=1)\n', confusion)

Logistic Regress

In []: `from sklearn.linear_model import LogisticRegression`

lr = LogisticRegression().fit(X_train, y_train)

lr_predicted = lr.predict(X_test)

confusion = confusion_matrix(y_test, lr_predicted)

print('Logistic regression classifier (default settings)\n', confusion) *418 correct*

Decision Tree

In []: `from sklearn.tree import DecisionTreeClassifier`

dt = DecisionTreeClassifier(max_depth=2).fit(X_train, y_train)

tree_predicted = dt.predict(X_test)

confusion = confusion_matrix(y_test, tree_predicted)

print('Decision tree classifier (max_depth = 2)\n', confusion)

<https://github.com/Qian-Han/coursera-Applied-Data-Science...pply-Machine-Learning-In-Python/week3/Module%203.ipynb> 2020/4/12 下午11:38 第2页 (共8页)

*TN FP
FN TP*

tree_predicted = dt.predict(X_test)

confusion = confusion_matrix(y_test, tree_predicted)

print('Decision tree classifier (max_depth = 2)\n', confusion) *426 correct*

2 Evaluation metrics for binary classification

In []: `from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score`

Accuracy = TP + TN / (TP + TN + FP + FN)

Precision = TP / (TP + FP)

Recall = TP / (TP + FN) Also known as sensitivity, or True Positive Rate

*# F1 = 2 * Precision * Recall / (Precision + Recall)* *predict label of test set*

print('Accuracy: {:.2f}'.format(accuracy_score(y_test, tree_predicted))) *predict label of test set*

print('Precision: {:.2f}'.format(precision_score(y_test, tree_predicted))) *decision tree classifier*

print('Recall: {:.2f}'.format(recall_score(y_test, tree_predicted)))

print('F1: {:.2f}'.format(f1_score(y_test, tree_predicted)))

In []: `# Combined report with all above metrics`

`from sklearn.metrics import classification_report`

the label predict label const formed by output

print(classification_report(y_test, tree_predicted, target_names=['not 1', '1'])) *support: no. of instances*

4 different classifier required total class as in output table

In []: `print('Random class-proportional (dummy)\n',`

`classification_report(y_test, y_classprop_predicted, target_names=['not 1', '1']))` *P.R + → very low*

`print('SVM\n',`

`classification_report(y_test, svm_predicted, target_names = ['not 1', '1']))`

`print('Logistic regression\n',`

`classification_report(y_test, lr_predicted, target_names = ['not 1', '1']))`

`print('Decision tree\n',`

`classification_report(y_test, tree_predicted, target_names = ['not 1', '1']))`

3 Decision functions

Classification score

In []: `X_train, X_test, y_train, y_test = train_test_split(X, y_binary_imbalanced, random_state=0)`

`y_scores_lr = lr.fit(X_train, y_train).decision_function(X_test)`

`y_score_list = list(zip(y_test[0:20], y_scores_lr[0:20]))`

`# show the decision_function scores for first 20 instances`

`y_score_list`

Prediction probability

In []: `X_train, X_test, y_train, y_test = train_test_split(X, y_binary_imbalanced, random_state=0)`

`y_proba_lr = lr.fit(X_train, y_train).predict_proba(X_test)`

`y_proba_list = list(zip(y_test[0:20], y_proba_lr[0:20,1]))`

`# show the probability of positive class for first 20 instances` *(0.995 - 1 = 0.005)*

```
In [ ]: from sklearn.metrics import precision_recall_curve

precision, recall, thresholds = precision_recall_curve(y_test, y_scores_lr)
closest_zero = np.argmin(np.abs(thresholds))
closest_zero_p = precision[closest_zero]
closest_zero_r = recall[closest_zero]

plt.figure()
plt.xlim([0.0, 1.01])
plt.ylim([0.0, 1.01])
plt.plot(precision, recall, label='Precision-Recall Curve')
plt.plot(closest_zero_p, closest_zero_r, 'o', markersize = 12, fillstyle = 'none', c='r', mew=3)
plt.xlabel('Precision', fontsize=16)
plt.ylabel('Recall', fontsize=16)
plt.axes().set_aspect('equal')
plt.show()
```

ROC curves, Area Under Curve (AUC)

<https://github.com/Qian-Han/coursera-Applied-Data-Science---Applied-Machine-Learning-In-Python/week3/Module%203.ipynb>

2020/4/12 下午11:38
第3页 (共8页)

ROC curves, Area Under Curve (AUC)

```
In [ ]: from sklearn.metrics import roc_curve, auc

X_train, X_test, y_train, y_test = train_test_split(X, y_binary_imbalanced, random_state=0)

y_score_lr = lr.fit(X_train, y_train).decision_function(X_test)
fpr_lr, tpr_lr, _ = roc_curve(y_test, y_score_lr)
roc_auc_lr = auc(fpr_lr, tpr_lr)

plt.figure()
plt.xlim([-0.01, 1.00])
plt.ylim([-0.01, 1.01])
plt.plot(fpr_lr, tpr_lr, lw=3, label='LogRegr ROC curve (area = {:.2f})'.format(roc_auc_lr))
plt.xlabel('False Positive Rate', fontsize=16)
plt.ylabel('True Positive Rate', fontsize=16)
plt.title('ROC curve (1-of-10 digits classifier)', fontsize=16)
plt.legend(loc='lower right', fontsize=13)
plt.plot([0, 1], [0, 1], color='navy', lw=3, linestyle='--')
plt.axes().set_aspect('equal')
plt.show()
```

```
In [ ]: from matplotlib import cm

X_train, X_test, y_train, y_test = train_test_split(X, y_binary_imbalanced, random_state=0)

plt.figure()
plt.xlim([-0.01, 1.00])
plt.ylim([-0.01, 1.01])
for g in [0.01, 0.1, 0.20, 1]:
    svm = SVC(gamma=g).fit(X_train, y_train)
    y_score_svm = svm.decision_function(X_test)
    fpr_svm, tpr_svm, _ = roc_curve(y_test, y_score_svm)
    roc_auc_svm = auc(fpr_svm, tpr_svm)
    accuracy_svm = svm.score(X_test, y_test)
    print("gamma = {:.2f} accuracy = {:.2f} AUC = {:.2f}".format(g, accuracy_svm,
                                                               roc_auc_svm))
    plt.plot(fpr_svm, tpr_svm, lw=3, alpha=0.7,
             label='SVM (gamma = {:.2f}, area = {:.2f})'.format(g, roc_auc_svm))

plt.xlabel('False Positive Rate', fontsize=16)
plt.ylabel('True Positive Rate (Recall)', fontsize=16)
plt.plot([0, 1], [0, 1], color='k', lw=0.5, linestyle='--')
plt.legend(loc="lower right", fontsize=11)
plt.title('ROC curve: (1-of-10 digits classifier)', fontsize=16)
plt.axes().set_aspect('equal')

plt.show()
```

15 Evaluation measures for multi-class classification

Multi-class confusion matrix

```
In [ ]: dataset = load_digits()
X, y = dataset.data, dataset.target
X_train_mc, X_test_mc, y_train_mc, y_test_mc = train_test_split(X, y, random_state=0)

svm = SVC(kernel = 'linear').fit(X_train_mc, y_train_mc)
svm_predicted_mc = svm.predict(X_test_mc)
confusion_mc = confusion_matrix(y_test_mc, svm_predicted_mc)
df_cm = pd.DataFrame(confusion_mc,
                      index = [i for i in range(0,10)], columns = [i for i in range(0,10)])

plt.figure(figsize=(5.5,4))
sns.heatmap(df_cm, annot=True)
plt.title('SVM Linear Kernel \nAccuracy:{0:.3f}'.format(accuracy_score(y_test_mc,
                                                                      svm_predicted_mc)))
plt.ylabel('True label')
plt.xlabel('Predicted label')
```

most predictions correct
 True & → 3 times
 predict!

<https://github.com/Qian-Han/coursera-Applied-Data-Science---Applied-Machine-Learning-In-Python/week3/Module%203.ipynb>

2020/4/12 下午11:38
第4页 (共8页)

2nd = confusion matrix = worse
 64 instances one digit 2 17 classify ✓
 74 classify → digit 4 X

Show problem: By feature preprocessing

```

df_cm = pd.DataFrame(confusion_mc, index = [i for i in range(0,10)],
                     columns = [i for i in range(0,10)])

plt.figure(figsize = (5.5,4))
sns.heatmap(df_cm, annot=True)
plt.title('SVM RBF Kernel \nAccuracy:{0:.3f}'.format(accuracy_score(y_test_mc,
                                                                     svm_predicted_mc)))
plt.ylabel('True label')
plt.xlabel('Predicted label');

```

Multi-class classification report

```
In [ ]: print(classification_report(y_test_mc, svm_predicted_mc)) avg → each class
```

Micro- vs. macro-averaged metrics

```
In [ ]: print('Micro-averaged precision = {:.2f} (treat instances equally)'
           .format(precision_score(y_test_mc, svm_predicted_mc, average = 'micro')))
print('Macro-averaged precision = {:.2f} (treat classes equally)'
      .format(precision_score(y_test_mc, svm_predicted_mc, average = 'macro')))
```

```
In [ ]: print('Micro-averaged f1 = {:.2f} (treat instances equally)'
           .format(f1_score(y_test_mc, svm_predicted_mc, average = 'micro')))
print('Macro-averaged f1 = {:.2f} (treat classes equally)'
      .format(f1_score(y_test_mc, svm_predicted_mc, average = 'macro')))
```

Regression evaluation metrics

```

In [ ]: %matplotlib notebook
import matplotlib.pyplot as plt
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn import datasets
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
from sklearn.dummy import DummyRegressor

diabetes = datasets.load_diabetes()

X = diabetes.data[:, None, 6]
y = diabetes.target

X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)

lm = LinearRegression().fit(X_train, y_train)
lm_dummy_mean = DummyRegressor(strategy = 'mean').fit(X_train, y_train)

y_predict = lm.predict(X_test)
y_predict_dummy_mean = lm_dummy_mean.predict(X_test)

print('Linear model, coefficients: ', lm.coef_)
print("Mean squared error (dummy): {:.2f}".format(mean_squared_error(y_test,
                                                                     y_predict_dummy_mean)))
print("Mean squared error (linear model): {:.2f}".format(mean_squared_error(y_test, y_predict)))
print("r2_score (dummy): {:.2f}".format(r2_score(y_test, y_predict_dummy_mean)))
print("r2_score (linear model): {:.2f}".format(r2_score(y_test, y_predict)))

# Plot outputs
plt.scatter(X_test, y_test, color='black')
plt.plot(X_test, y_predict, color='green', linewidth=2)
plt.plot(X_test, y_predict_dummy_mean, color='red', linestyle = 'dashed',
         linewidth=2, label = 'dummy')

plt.show()

```

dummy = $r^2 = 0$ constant prediction
linear $r^2 > 0$, ob slightly better fit

<https://github.com/Qian-Han/coursera-Applied-Data-Science---Applied-Machine-Learning-In-Python/week3/Module%203.ipynb>

2020/4/12 下午11:38
第5页 (共8页)

7 Model selection using evaluation metrics

Cross-validation example

```

In [ ]: from sklearn.model_selection import cross_val_score
from sklearn.svm import SVC

dataset = load_digits()
# again, making this a binary problem with 'digit 1' as positive class
# and 'not 1' as negative class
X, y = dataset.data, dataset.target == 1 5 folds

clf = SVC(kernel='rbf', C=1)

# accuracy is the default scoring metric
print('Cross-validation (accuracy)', cross_val_score(clf, X, y, cv=5))
# use AUC as scoring metric
print('Cross-validation (AUC)', cross_val_score(clf, X, y, cv=5, scoring = 'roc_auc'))
# use recall as scoring metric
print('Cross-validation (recall)', cross_val_score(clf, X, y, cv=5, scoring = 'recall'))

```

default accuracy - evaluation metric
AUC evaluation metric

resulting test = 5 evaluation values - 1 / fold each metric

Grid search example

```

In [ ]: from sklearn.svm import SVC
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import roc_auc_score

dataset = load_digits()
X, y = dataset.data, dataset.target == 1
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)

clf = SVC(kernel='rbf') set width/influence of kernel
grid_values = {'gamma': [0.001, 0.01, 0.05, 0.1, 1, 10, 100]}

# default metric to optimize over grid parameters: accuracy
grid_clf_acc = GridSearchCV(clf, param_grid = grid_values) and Y-optimize evaluation metrics
grid_clf_acc.fit(X_train, y_train)
y_decision_fn_scores_acc = grid_clf_acc.decision_function(X_test)

print('Grid best parameter (max. accuracy): ', grid_clf_acc.best_params_)
print('Grid best score (accuracy): ', grid_clf_acc.best_score_)

# alternative metric to optimize over grid parameters: AUC
grid_clf_auc = GridSearchCV(clf, param_grid = grid_values, scoring = 'roc_auc')
grid_clf_auc.fit(X_train, y_train)
y_decision_fn_scores_auc = grid_clf_auc.decision_function(X_test)

print('Test set AUC: ', roc_auc_score(y_test, y_decision_fn_scores_auc))
print('Grid best parameter (max. AUC): ', grid_clf_auc.best_params_)
print('Grid best score (AUC): ', grid_clf_auc.best_score_)

```

Evaluation metrics supported for model selection

```
In [ ]: from sklearn.metrics.scorer import SCORERS
print(sorted(list(SCORERS.keys())))

```

Two-feature classification example using the digits dataset

Optimizing a classifier using different evaluation metrics

```

In [ ]: from sklearn.datasets import load_digits
from sklearn.model_selection import train_test_split
from adspy_shared_utilities import plot_class_regions_for_classifier_subplot
from sklearn.svm import SVC
from sklearn.model_selection import GridSearchCV

```

*optional class weight ↑ → ↑ weight → positive class
 ↗ 2 classes - train*

<https://github.com/Qian-Han/coursera-Applied-Data-Science---Applied-Machine-Learning-In-Python/week3/Module%203.ipynb>

2020/4/12 下午11:38
第6页 (共8页)

```

dataset = load_digits()
X, y = dataset.data, dataset.target == 1
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)

# Create a two-feature input vector matching the example plot above
# We jitter the points (add a small amount of random noise) in case there are areas
# in feature space where many instances have the same features.
jitter_delta = 0.25
X_twovar_train = X_train[:,[20,59]]+ np.random.rand(X_train.shape[0], 2) - jitter_delta
X_twovar_test = X_test[:,[20,59]] + np.random.rand(X_test.shape[0], 2) - jitter_delta

clf = SVC(kernel = 'linear').fit(X_twovar_train, y_train)
grid_values = {'class_weight':[{'balanced', {1:2},{1:3},{1:4},{1:5},{1:10},{1:20},{1:50}}]}

```

```

print('Grid best score ({0}): {1}'.format(eval_metric, grid_clf_custom.best_score_))
plt.subplots_adjust(wspace=0.3, hspace=0.3)
plot_class_regions_for_classifier_subplot(grid_clf_custom, X_twovar_test, y_test, None,
                                           None, plt.subplot(2, 2, i+1))

plt.title(eval_metric+'-oriented SVC')
plt.tight_layout()
plt.show()

find most positive class points -> positive class predictors, recall & F1 - boy
Precision-recall curve for the default SVC classifier (with balanced class weights) precision & F1 - boy
F1 - boy - integrated Charnier
AUC {125}

```

In []:

```

from sklearn.model_selection import train_test_split
from sklearn.metrics import precision_recall_curve
from adspy_shared_utilities import plot_class_regions_for_classifier
from sklearn.svm import SVC

dataset = load_digits()
X, y = dataset.data, dataset.target == 1
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)

# create a two-feature input vector matching the example plot above
jitter_delta = 0.25
X_twovar_train = X_train[:,[20,59]] + np.random.rand(X_train.shape[0], 2) - jitter_delta
X_twovar_test = X_test[:,[20,59]] + np.random.rand(X_test.shape[0], 2) - jitter_delta

clf = SVC(kernel='linear', class_weight='balanced').fit(X_twovar_train, y_train)

y_scores = clf.decision_function(X_twovar_test)

precision, recall, thresholds = precision_recall_curve(y_test, y_scores)
closest_zero = np.argmin(np.abs(thresholds))
closest_zero_p = precision[closest_zero]
closest_zero_r = recall[closest_zero]

plot_class_regions_for_classifier(clf, X_twovar_test, y_test)
plt.title("SVC, class_weight = 'balanced', optimized for accuracy")
plt.show()

plt.figure()
plt.xlim([0.0, 1.01])
plt.ylim([0.0, 1.01])
plt.title ("Precision-recall curve: SVC, class_weight = 'balanced'")
plt.plot(precision, recall, label = 'Precision-Recall Curve')
plt.plot(closest_zero_p, closest_zero_r, 'o', markersize=12, fillstyle='none', c='r', mew=3)
plt.xlabel('Precision', fontsize=16)
plt.ylabel('Recall', fontsize=16)
plt.axes().set_aspect('equal')
plt.show()
print('At zero threshold, precision: {:.2f}, recall: {:.2f}'
      .format(closest_zero_p, closest_zero_r))

```