

Jan 10th 2021

Ref: hands-on ML Chapter 3 classification.

Loading dataset in sklearn:

- sklearn provides helper function for downloading datasets
eg. `>>> from sklearn.datasets import fetch_openml`.
`>>> mnist = fetch_openml('mnist_784', version=1)`

Contains multiple cleaned datasets.

- The dataset returned from the above operation contains many useful fields:

- Among them:
- `DESCR`: Description of the dataset.
 - `data`: data entries (one row per instance, one col per feature).
 - `target`: Array with labels.

Naive Way to split data in to train and test.

`X-train, X-test, y-train, y-test = x[:6000], x[6000:], y[:6000], y[6000:]`.

Training a Binary Classifier.

- modify labels from categories to true and false.

eg. `>>> y-train-5 = (y-train == 5).`
new label array : [T, F, ..]
↑
is-5 not-5

- `sgd-clf = SGDClassifier(random_state=42)`

SGDClassifier relies on randomness during training.

* same results will be reproduced if same random_state is set the same.

- evaluation:

- Accuracy: percentage of correct predictions.

This is normally not a good measurement.

for example: in the case of binary classification where positive cases are a lot less than negatives, even a naive model (such as a model that always predict false) will produce a high accuracy score.

- Confusion matrix

Sklearn = `>>> y-train-pred = cross_val_predict(sgd-clf, x-train, y-train-5, cv=3).`
return prediction results (labels).

`>>> confusion_matrix(y-train-5, y-train-pred).`

matrix :

TN	FP	Actual
A	B	Negative

FN	TP	positive
C	D	
Predicted negative	positive	

- Precision, Recall, F1.

$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$ How many correct predictions among predictions that label positive?

$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$ How many positives are detected.

$F_1 = 2 \times \frac{\text{Precision} \times \text{recall}}{\text{Precision} + \text{recall}} = \frac{\text{TP}}{\text{TP} + \frac{\text{FN} + \text{FP}}{2}}$

F_1 : harmonic mean of precision and recall. Which gives much more weight to low values.

high F_1 if both recall and precision are high

F_1 score favors classifiers that have similar precision and recall.

```
>>> from sklearn.metrics import f1_score
```

```
f1_score(y_train_s, y_train_pred)
```