

```
from sklearn.datasets import fetch_openml
from sklearn import metrics
from sklearn.model_selection import train_test_split
import pandas as pd
```

```
import datetime
x1 = datetime.datetime.now()
print(x1)
```

```
↳ 2019-12-06 21:54:12.049991
```

```
mnist = fetch_openml('mnist_784', version=1, cache=True )
```

```
#mnist
```

```
mnist.data.shape
```

```
↳ (70000, 784)
```

```
mnist.target.shape
```

```
↳ (70000,)
```

```
# test_size: what proportion of original data is used for test set
train_img, test_img, train_lbl, test_lbl = train_test_split(
    mnist.data, mnist.target, test_size=1/7.0, random_state=0)
```

```
print(train_img.shape)
```

```
↳ (60000, 784)
```

```
print(train_lbl.shape)
```

```
↳ (60000,)
```

```
print(test_img.shape)
```

```
↳ (10000, 784)
```

```
print(test_lbl.shape)
```

```
↳ (10000,)
```

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
```

```
# Fit on training set only.
scaler.fit(train_img)

# Apply transform to both the training set and the test set.
train_img = scaler.transform(train_img)
test_img = scaler.transform(test_img)
```

▼ Step 1: Import the model you want to use

In sklearn, all machine learning models are implemented as Python classes

```
from sklearn import svm
clf = svm.SVC(gamma=0.001, C=100) # các tham số cho mô hình hoạt động tốt hơn
```

▼ Step 2: Training the model on the data, storing the information learned from the data

Model is learning the relationship between x (digits) and y (labels)

```
clf.fit(train_img, train_lbl)

☞ SVC(C=100, cache_size=200, class_weight=None, coef0=0.0,
      decision_function_shape='ovr', degree=3, gamma=0.001, kernel='rbf',
      max_iter=-1, probability=False, random_state=None, shrinking=True,
      tol=0.001, verbose=False)
```

```
y_pred = clf.predict(test_img)
y_pred
```

```
☞ array(['0', '4', '1', ..., '1', '3', '0'], dtype=object)
```

▼ Measuring Model Performance

accuracy (fraction of correct predictions): correct predictions / total number of data points

Basically, how the model performs on new data (test set)

```
from sklearn.metrics import accuracy_score
print("Accuracy is ", accuracy_score(test_lbl, y_pred)*100, "%")
```

```
☞ Accuracy is 97.31 %
```

```
score = clf.score(test_img, test_lbl)
print(score)
```

```
☞ 0.9731
```

▼ Step 3: Predict the labels of new data (new images)

Uses the information the model learned during the model training process

```
new = clf.predict(test_img[0].reshape(1,-1))  
new
```

```
↳ array(['0'], dtype=object)
```

```
x2 = datetime.datetime.now()  
print(x2)
```

```
↳ 2019-12-06 22:13:04.694182
```

```
# distance  
d = x2 - x1  
print(d)
```

```
↳ 0:18:52.644191
```