# scikit-learn
# From Zero to Hero!

Alexandre Gramfort

http://alexandre.gramfort.net
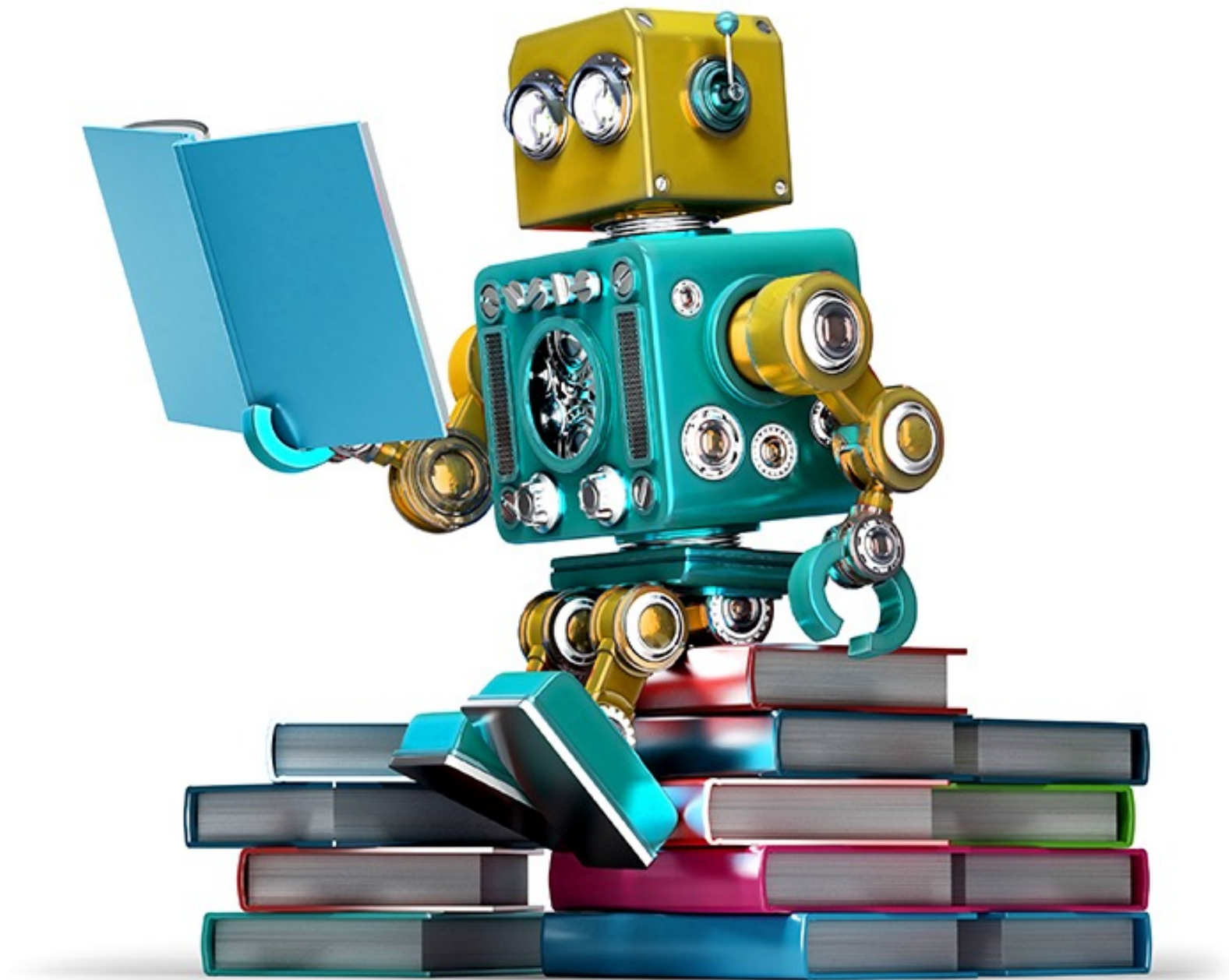
https://github.com/agramfort/blend2017

*GitHub : @agramfort*          *Twitter : @agramfort*

**scikit-learn** is a ML library in **python**

but what is machine learning?

**Artificial Intelligence**

**Predictive Modeling (Data Analytics)**

IBM Watson

Movie recommendations

Self-driving cars

Predictive Maintenance

Slide courtesy of @ogrisel

Slide courtesy of @ogrisel

Slide courtesy of @ogrisel

**Definition:** Machine learning consists in teaching a computer to make decisions based on examples

machine learning is one (very effective) way to solve AI problems

Source: https://www.openhub.net/p/scikit-learn

In a Nutshell, scikit-learn...

... has had 22,256 commits made by 1,135 contributors
representing 139,215 lines of code

... took an estimated 36 years of effort (COCOMO model)
starting with its first commit in January, 2010
ending with its most recent commit about 20 hours
ago

- > 21500 followers and 11500 forks on github.com

- 474k entries on google.com

http://scikit-learn.org

Source: https://www.openhub.net/p/scikit-learn
In a Nutshell, scikit-learn...

... has had 22,256 commits made by 1,135 contributors
representing 139,215 lines of code

... took an estimated 36 years of effort (COCOMO model)
starting with its first commit in January, 2010
ending with its most recent commit about 20 hours
ago

• > 21500 followers and 11500 forks on github.com
• 474k entries on google.com

Funding:

http://scikit-learn.org

http://scikit-learn.org/stable/testimonials/testimonials.html

- Installed on 1% of debian systems
- 1200 job offers on Stack-Overflow
- Users: 60% academics & 40% industry
- > 400,000 regular users (web stats)

http://scikit-learn.org/stable/testimonials/testimonials.html

# Machine Learning Taxonomy

Machine Learning

Unsupervised learning — Supervised learning

Clustering — Dimensionality Reduction — Anomaly Detection

Regression — Classification — Ranking

# Machine Learning Taxonomy

Machine Learning

Unsupervised learning     Supervised learning

Clustering    Dimensionality Reduction    Anomaly Detection      Regression    Classification    Ranking

let's start with supervised learning workflow....

| type (category) | # rooms (int) | surface (float m2) | public trans (boolean) |
|---|---|---|---|
| Apartment | 3 | 50 | TRUE |
| House | 5 | 254 | FALSE |
| Duplex | 4 | 68 | TRUE |
| Apartment | 2 | 32 | TRUE |

| type (category) | # rooms (int) | surface (float m2) | public trans (boolean) | sold (float k€) |
|---|---|---|---|---|
| Apartment | 3 | 50 | TRUE | 450 |
| House | 5 | 254 | FALSE | 430 |
| Duplex | 4 | 68 | TRUE | 712 |
| Apartment | 2 | 32 | TRUE | 234 |

|  | features |  |  | target |
| :---: | :---: | :---: | :---: | :---: |
| **type (category)** | **# rooms (int)** | **surface (float m2)** | **public trans (boolean)** | **sold (float k€)** |
| Apartment | 3 | 50 | TRUE | 450 |
| House | 5 | 254 | FALSE | 430 |
| Duplex | 4 | 68 | TRUE | 712 |
| Apartment | 2 | 32 | TRUE | 234 |

samples (train)

|  | features | | | | target |
| --- | --- | --- | --- | --- | --- |
| | **type (category)** | **# rooms (int)** | **surface (float m2)** | **public trans (boolean)** | **sold (float k€)** |
| **samples (train)** | Apartment | 3 | 50 | TRUE | 450 |
| | House | 5 | 254 | FALSE | 430 |
| | Duplex | 4 | 68 | TRUE | 712 |
| | Apartment | 2 | 32 | TRUE | 234 |
| **samples (test)** | Apartment | 2 | 33 | TRUE | ? |
| | House | 4 | 210 | TRUE | ? |

features    target

| type (category) | # rooms (int) | surface (float m2) | public trans (boolean) | sold (float k€) |
|---|---|---|---|---|
| Apartment | 3 | 50 | TRUE | 450 |
| House | | | FALSE | |
| Duplex | 4 | 68 | TRUE | 712 |
| Apartment | 2 | 32 | TRUE | 234 |

samples (train) — X train — y train

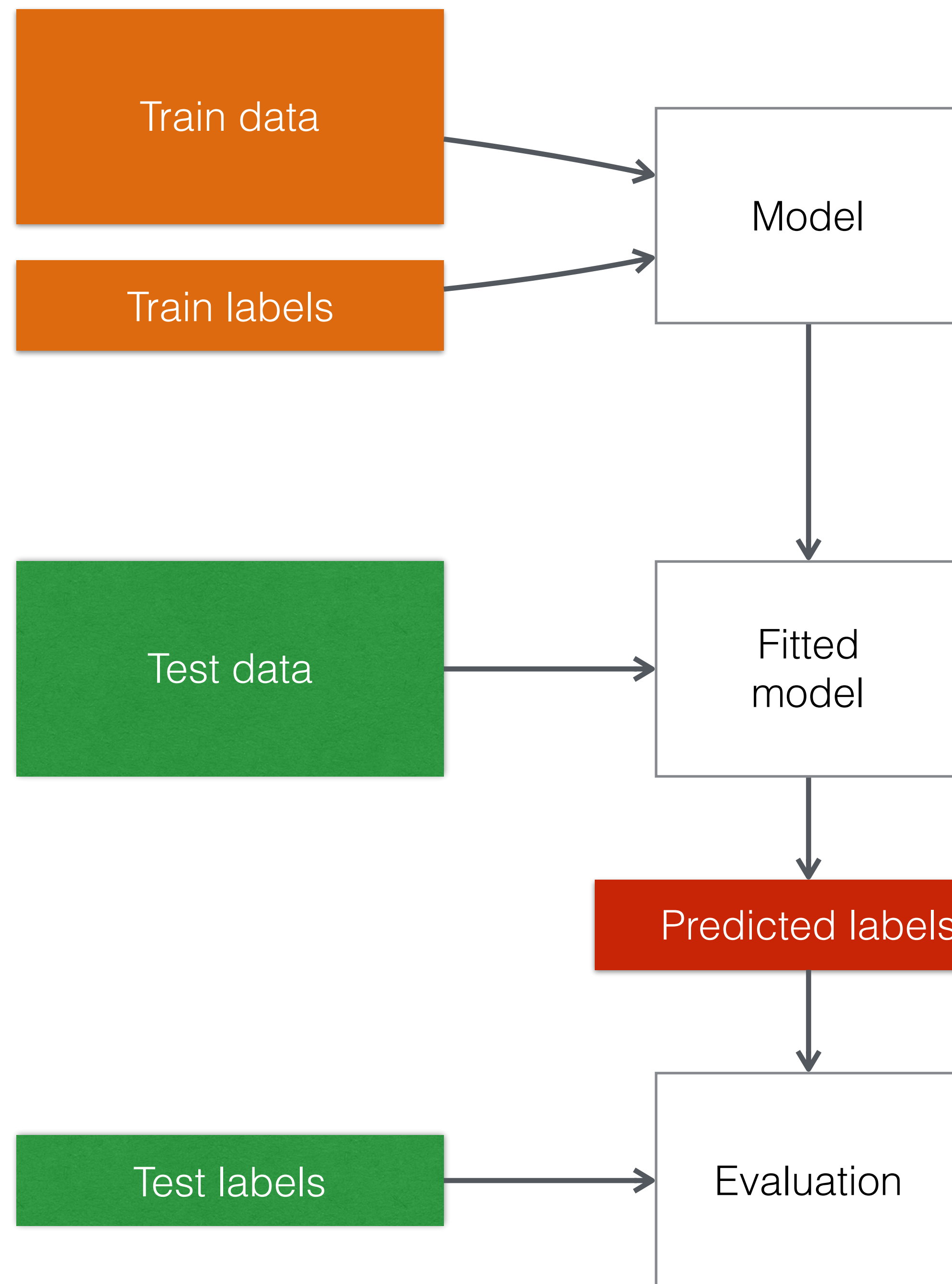| type (category) | # rooms (int) | surface (float m2) | public trans (boolean) | sold (float k€) |
|---|---|---|---|---|
| Apartment | | | TRUE | |
| House | | 210 | TRUE | |

samples (test) — X test — y pred

```
>>> model = LogisticRegression(C=1)
>>> model.fit(X_train, y_train)
```

```
>>> model = LogisticRegression(C=1)
>>> model.fit(X_train, y_train)
```

```
>>> y_pred = model.predict(X_test)
```
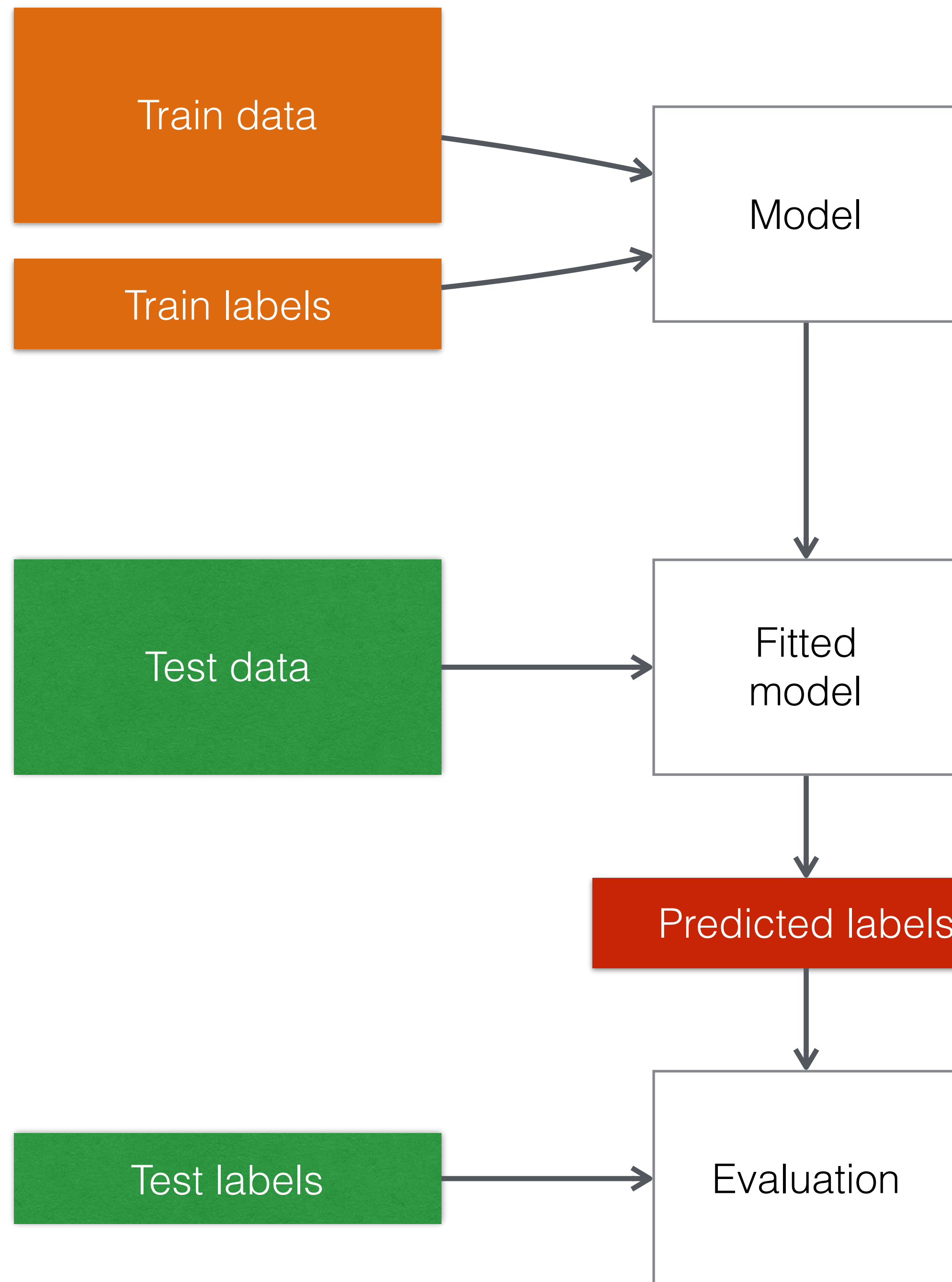
| Train data | | Model |
| Train labels | | |

```
>>> model = LogisticRegression(C=1)
>>> model.fit(X_train, y_train)
```

| Test data | | Fitted model |

```
>>> y_pred = model.predict(X_test)
```

Predicted labels

| Test labels | | Evaluation |

```
>>> accuracy_score(y_test, y_pred)
```

```
>>> model = LogisticRegression(C=1)
>>> model.fit(X_train, y_train)
```

**fit**

```
>>> y_pred = model.predict(X_test)
```

**predict**

```
>>> accuracy_score(y_test, y_pred)
```
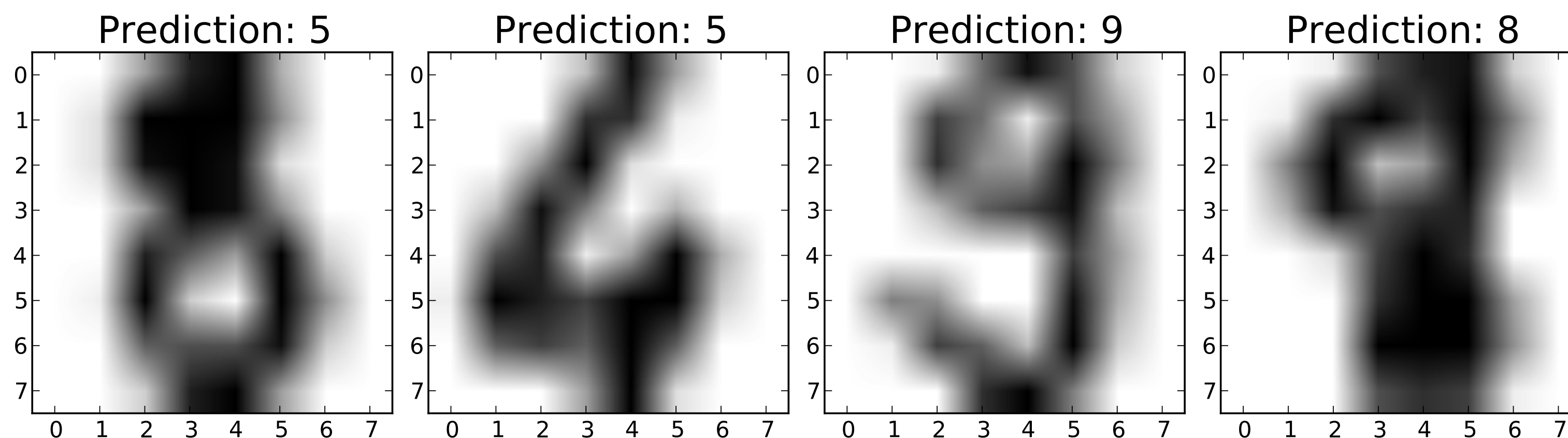
**score**

## Classification of images of digits in a few lines of code

```python
import matplotlib.pyplot as plt
from sklearn import datasets, svm
# Load data
digits = datasets.load_digits()
n_samples = len(digits.images)
data = digits.images.reshape((n_samples, -1))
# Learn ie. fit
classifier = svm.SVC()
classifier.fit(data[:n_samples // 2], digits.target[:n_samples // 2])
# Predict and plot
for index, image in enumerate(digits.images[n_samples // 2:n_samples // 2 + 4]):
    plt.subplot(1, 4, index)
    plt.imshow(image, cmap=plt.cm.gray_r)
    plt.title('Prediction: %i' % classifier.predict(image.ravel()), fontsize=20)
```

Prediction: 5    Prediction: 5    Prediction: 9    Prediction: 8

powerful yet easy!

```
>>> from sklearn import Model
>>> model = Model(param1=1e-8, param2="auto")
>>> print(model.param2)
"auto"
>>> model.fit(X_train, y_train)   # learn from training data
>>> y_pred = model.predict(X_test)   # predict from new data
>>> model.score(X_test, y_test)   # evaluate performance new data
0.96
```

## API: fit, predict, score

features

target

samples

| type (category) | # rooms (int) | surface (float m2) | public trans (boolean) | sold (float k€) |
|---|---|---|---|---|
| Apartment | 3 | 50 | TRUE | 150 |
| House | 5 | 254 | FALSE | 430 |
| Duplex | 4 | 68 | TRUE | 712 |
| Apartment | 2 | 32 | TRUE | 234 |

We have just X (no y)

```
>>> from sklearn import Model
>>> model = Model(param1=1e-8, param2="auto")
>>> print(model.param2)
"auto"
>>> model.fit(X)  # learn from training data (no y)
>>> Xt = model.transform(X)  # transform new or same data
```

API: fit, transform

**Contact:**

Alexandre Gramfort
http://alexandre.gramfort.net

*GitHub : @agramfort*

*Twitter : @agramfort*

informatics / mathematics
Inria