

# Deep learning

---

Romain Tavenard (Université de Rennes 2)

# Administrative details

---

- 20 hours
- Instructor: Romain Tavenard  
[romain.tavenard@univ-rennes2.fr](mailto:romain.tavenard@univ-rennes2.fr) (A201)
- Webpage:  
[rtavenar.github.io/teaching/ml/](https://rtavenar.github.io/teaching/ml/)
- Evaluation
  - Individual project (scikit-learn)
  - End of semester:
    - Multiple Choice Questionnaire (0.5h)
    - Lab session (1.5h)

# Some slides are more important than others...

---

- Slides marked with this symbol:



Are considered basic knowledge to pass the exams

# Contents

---

- Machine Learning in Python with scikit-learn
- Intro to deep learning
- Fully-connected models
- Images & ConvNets
- Generative models

# Machine learning in Python: scikit-learn

---

Romain Tavenard (Université de Rennes 2)

# Installing scikit-learn

---

- On Google Colab:
  - Already installed
- On your machine:
  - `pip install scikit-learn`  
or  
`conda install scikit-learn`
- Loading the module in Python:
  - `import sklearn`
- Most importantly
  - The doc is awesome, use it!  
<https://scikit-learn.org/stable/index.html>



# Preparing your data

---

- In scikit-learn:
  - $X$  is a  $(n, p)$  numpy array  
n input observations in dimension  $p$
  - $y$  is a  $(n, )$  or  $(n, p_{\text{out}})$  numpy array  
expected outputs
- Pre-processing
  - <https://scikit-learn.org/stable/modules/preprocessing.html>
  - Pre-processing methods are Transformers (cf. next slide)



# Transformers

---

- Objects with following methods
  - `fit(X)`
  - `transform(X)`
  - `fit_transform(X)`
- Example (from scikit-learn docs)

```
>>> from sklearn.preprocessing import StandardScaler
>>> data = [[0, 0], [0, 0], [1, 1], [1, 1]]
>>> scaler = StandardScaler()
>>> scaler.fit(data)
>>> print(scaler.mean_)
[0.5 0.5]
>>> print(scaler.transform(data))
[[-1. -1.]
 [-1. -1.]
 [ 1.  1.]
 [ 1.  1.]]
>>> print(scaler.transform([[2, 2]]))
[[3. 3.]]
```





# Estimators

---

- Objects with following methods
  - `fit(X[, y])`
  - `predict(X)`
  - `predict_proba(X)` (but not always...)
- Example (from scikit-learn docs)

```
>>> X = [[0], [1], [2], [3]]
>>> y = [0, 0, 1, 1]
>>> from sklearn.neighbors import KNeighborsClassifier
>>> neigh = KNeighborsClassifier(n_neighbors=3)
>>> neigh.fit(X, y)
KNeighborsClassifier(...)
>>> print(neigh.predict([[1.1]]))
[0]
>>> print(neigh.predict_proba([[0.9]]))
[[0.66666667 0.33333333]]
```

# Pipeline

---

- Objects that represent successive Transformers / Estimators
- Example (adapted from scikit-learn docs)

```
from sklearn import datasets
from sklearn.decomposition import PCA
from sklearn.linear_model import SGDClassifier
from sklearn.pipeline import Pipeline

logistic = SGDClassifier(loss='log', penalty='l2', early_stopping=True,
                        max_iter=10000, tol=1e-5, random_state=0)
pca = PCA(n_components=30)
pipe = Pipeline(steps=[('pca', pca), ('logistic', logistic)])

digits = datasets.load_digits()
X_digits = digits.data
y_digits = digits.target

pipe.fit(X_digits, y_digits)
```



# Model selection

---

- Goal: compare several models / hyper-parameter sets
- Base object in scikit-learn: **GridSearchCV**
- Example (from scikit-learn docs)

```
>>> from sklearn import svm, datasets
>>> from sklearn.model_selection import GridSearchCV
>>> iris = datasets.load_iris()
>>> parameters = {'kernel':('linear', 'rbf'), 'C':[1, 10]}
>>> svc = svm.SVC(gamma="scale")
>>> clf = GridSearchCV(svc, parameters, cv=5)
>>> clf.fit(iris.data, iris.target)
...
GridSearchCV(cv=5, error_score=...,
             estimator=SVC(C=1.0, cache_size=..., class_weight=..., coef0=...,
                           decision_function_shape='ovr', degree=..., gamma=...,
                           kernel='rbf', max_iter=-1, probability=False,
                           random_state=None, shrinking=True, tol=...,
                           verbose=False),
             iid=..., n_jobs=None,
             param_grid=..., pre_dispatch=..., refit=..., return_train_score=...,
             scoring=..., verbose=...)
```

# Conclusion

---

- What you should know about scikit-learn
  - Data formats
  - Basic objects
  - **How to read the docs**
- On the course page
  - Tutorials for basic scikit-learn usage

# Project

---

- Data: Horse Colic Dataset
  - Download: UCI Machine Learning Repository  
<http://archive.ics.uci.edu/ml/datasets/Horse+Colic>
  - Problem type: Classification
  - Variable to be predicted: "Surgical lesion?"
- Goal
  - What is the best model (from sklearn-implemented ones) for this data, and what is its performance?
- Format
  - Upload your notebook on CURSUS
- Deadline
  - December 20th, 23:59