

# Scikit-learn:

Machine learning in Python land

Gaël Varoquaux

inria



- 
- 1 The vision: enabling machine learning
  - 2 The tool: a Python library
  - 3 The project: a community

# 1 The vision: enabling machine learning

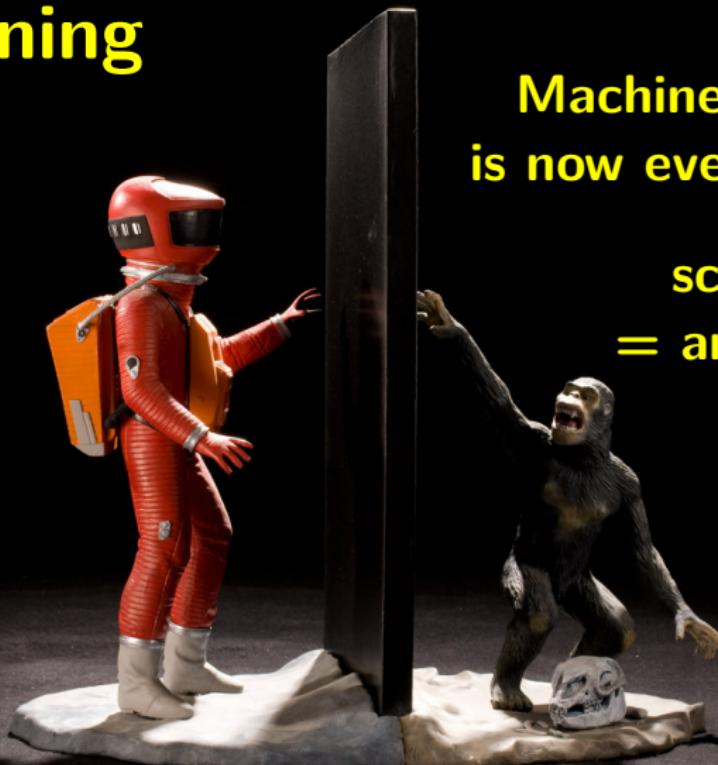


# 1 The vision: enabling machine learning

Machine-learning  
is now everywhere



# 1 The vision: enabling machine learning

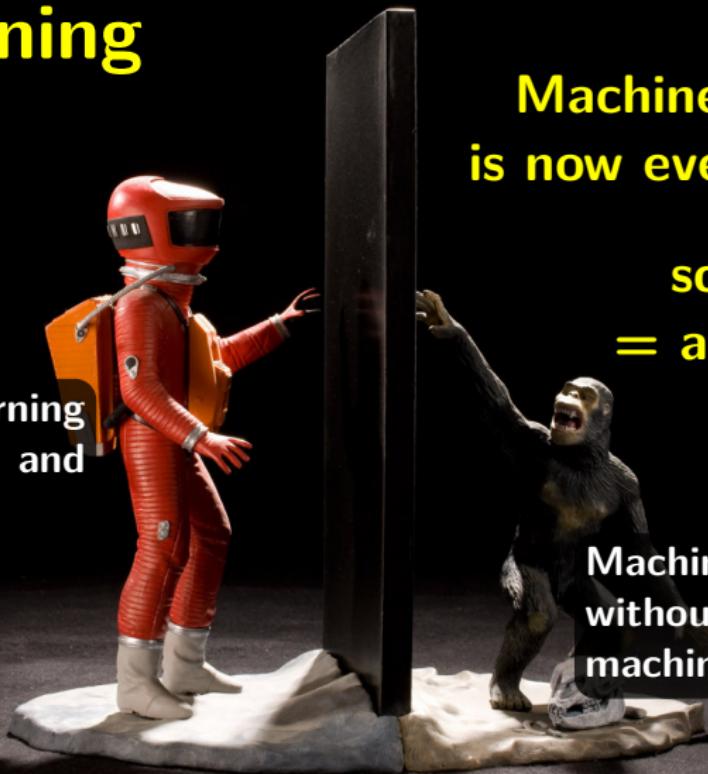


Machine-learning  
is now everywhere

scikit-learn  
= an enabler

# 1 The vision: enabling machine learning

Machine learning  
for everybody and  
for everything



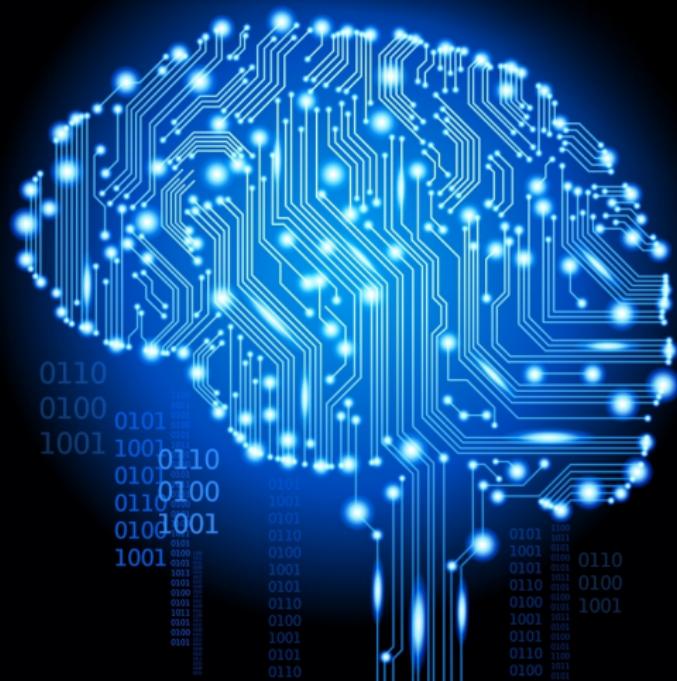
Machine-learning  
is now everywhere

scikit-learn  
= an enabler

Machine learning  
without learning the  
machinery

## Machine learning in a nutshell

Machine learning is about making prediction from data

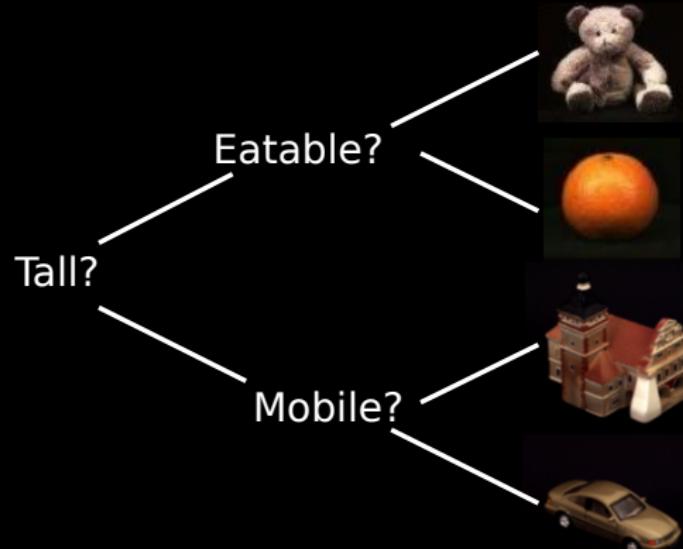


# 1 Machine learning: a historical perspective

## Expert systems

- Building decision rules

The 80s



# 1 Machine learning: a historical perspective

## Expert systems

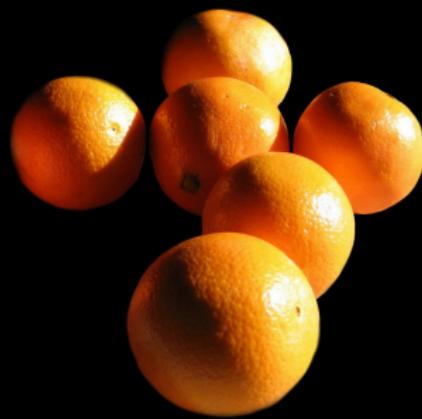
The 80s

- Building decision rules

## Machine learning

The 90s

- Learn these from observations



# 1 Machine learning: a historical perspective

## Expert systems

- Building decision rules

The 80s

## Machine learning

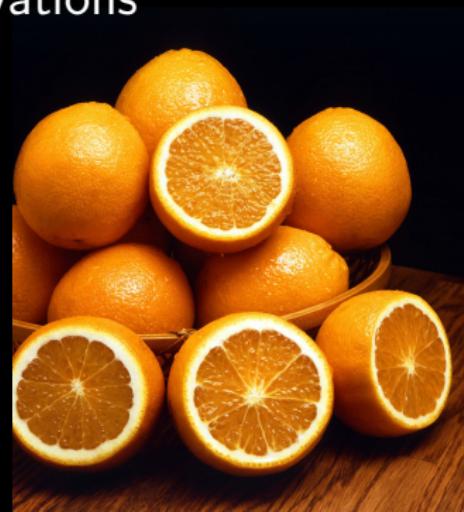
- Learn these from observations

The 90s

## Statistical learning

- Model the noise in the observations

2000s



# 1 Machine learning: a historical perspective

## Expert systems

- Building decision rules

The 80s

## Machine learning

- Learn these from observations

The 90s

## Statistical learning

- Model the noise in the observations

2000s

## Big data

- Many observations

today



# 1 Machine learning: a historical perspective

## Expert systems

- Building decision rules

The 80s

## Machine learning

- Learn these from observations

The 90s

## Statistical learning

- Model the noise in the observations

2000s

## Big data

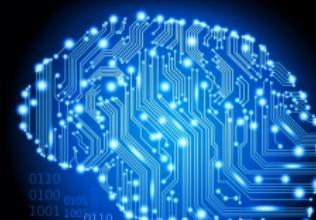
- Many observations

today

## Artificial Intelligence

- Machine learning on hard problems

tomorrow



# 1 Machine learning in a nutshell: an example

## Face recognition



Andrew



Bill



Charles

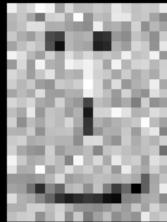


Dave

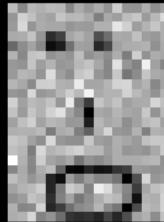


# 1 Machine learning in a nutshell: an example

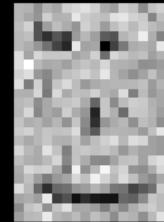
## Face recognition



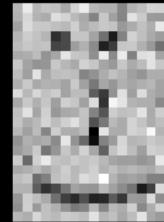
Andrew



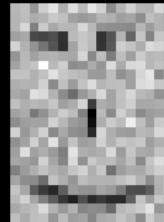
Bill



Charles



Dave



?

# 1 Machine learning in a nutshell

A simple method:

- 1 Store all the known (noisy) images and the names that go with them.
- 2 From a new (noisy) images, find the image that is most similar.

**“Nearest neighbor” method**



# 1 Machine learning in a nutshell

## A simple method:

- 1 Store all the known (noisy) images and the names that go with them.
- 2 From a new (noisy) images, find the image that is most similar.

## “Nearest neighbor” method

How many errors on already-known images?

...

0: no errors

## 1 Machine learning in a nutshell

A simple method:

- 1 Store all the known (noisy) images and the names that go with them.
- 2 From a new (noisy) images, find the image that is most similar.

### “Nearest neighbor” method

How many errors on already-known images?

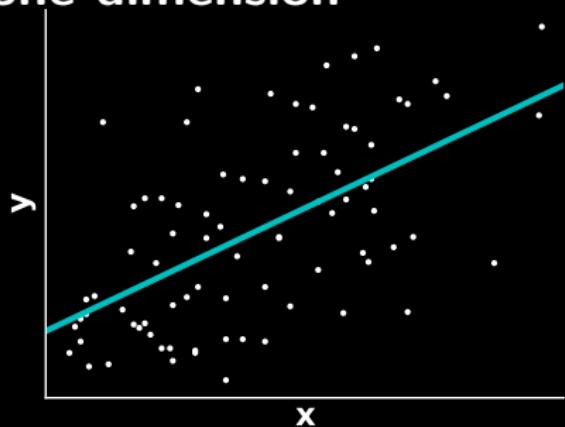
...

0: no errors

Test data  $\neq$  Train data

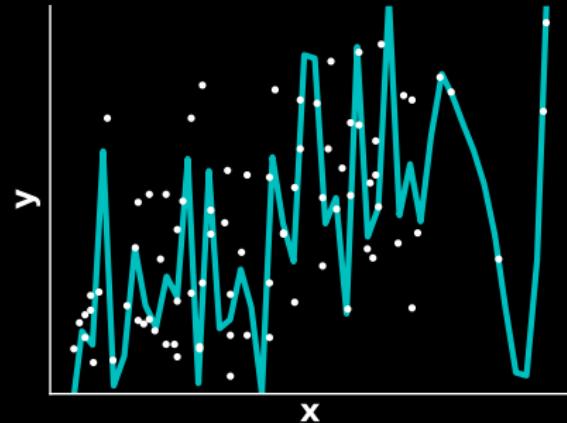
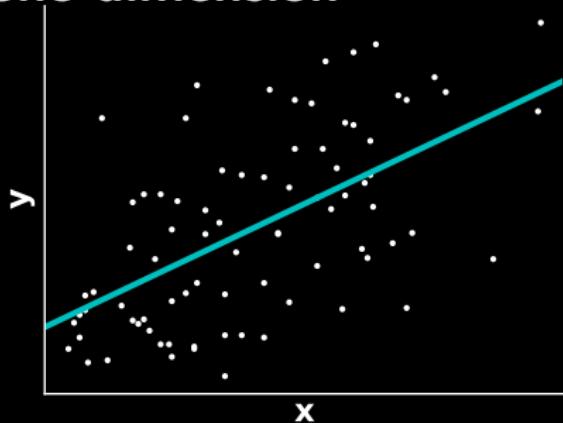
# 1 Machine learning in a nutshell: regression

A single descriptor:  
one dimension



# 1 Machine learning in a nutshell: regression

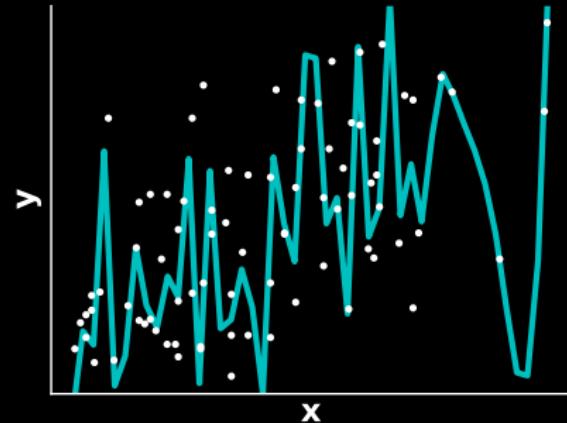
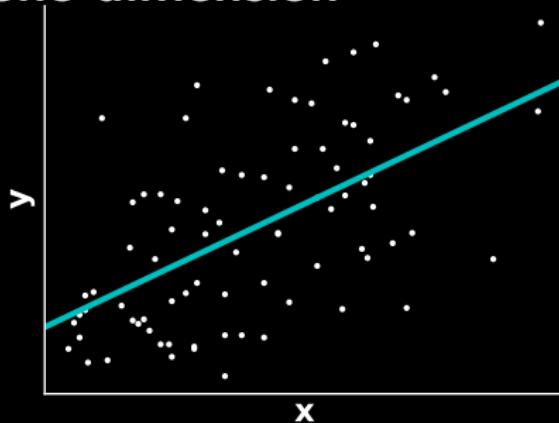
A single descriptor:  
one dimension



Which model to prefer?

# 1 Machine learning in a nutshell: regression

A single descriptor:  
one dimension

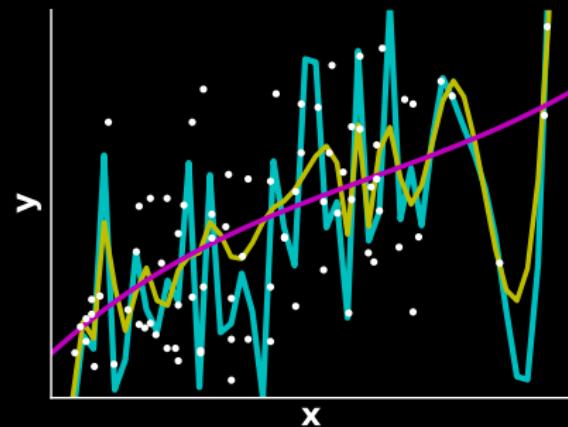
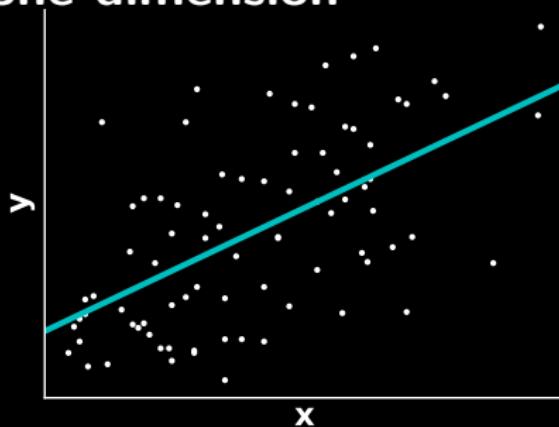


## Problem of “*over-fitting*”

- Minimizing error is not always the best strategy  
(learning noise)
- Test data  $\neq$  train data

# 1 Machine learning in a nutshell: regression

A single descriptor:  
one dimension



Prefer simple models

= concept of “regularization”

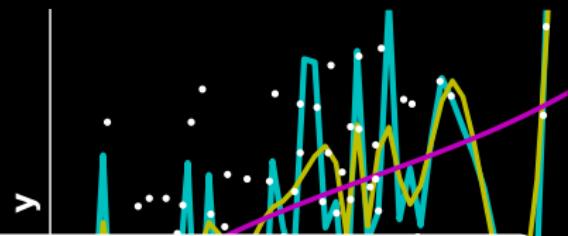
Balance the number of parameters to learn  
with the amount of data

# 1 Machine learning in a nutshell: regression

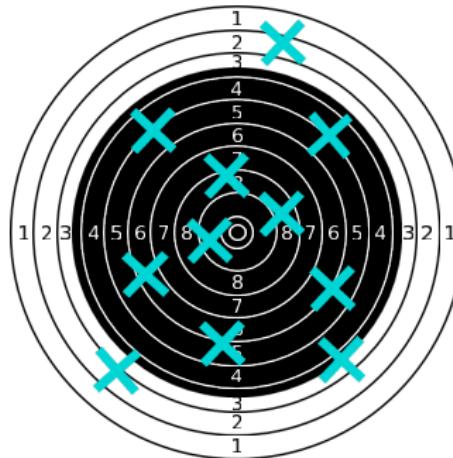
A single descriptor:  
one dimension



Bias

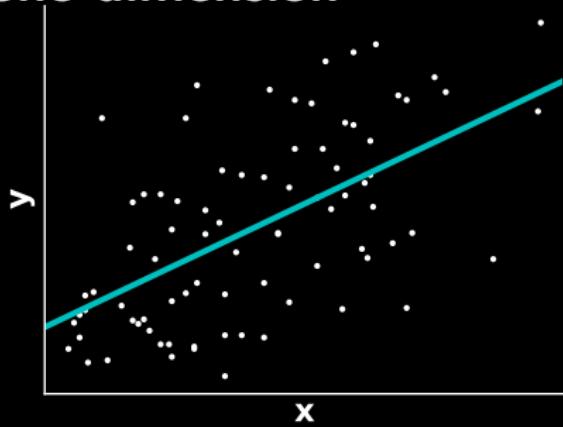


variance tradeoff

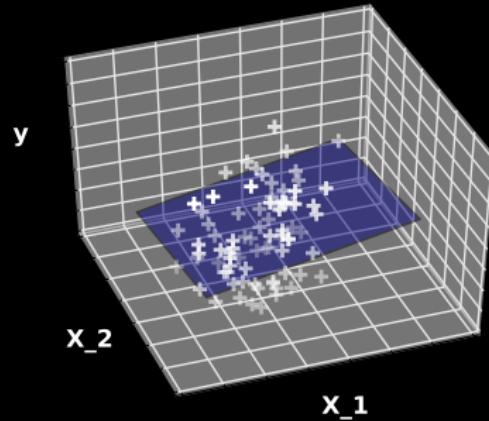


# 1 Machine learning in a nutshell: regression

A single descriptor:  
one dimension



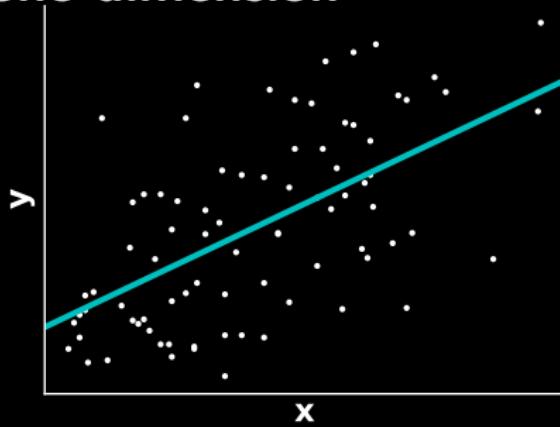
Two descriptors:  
2 dimensions



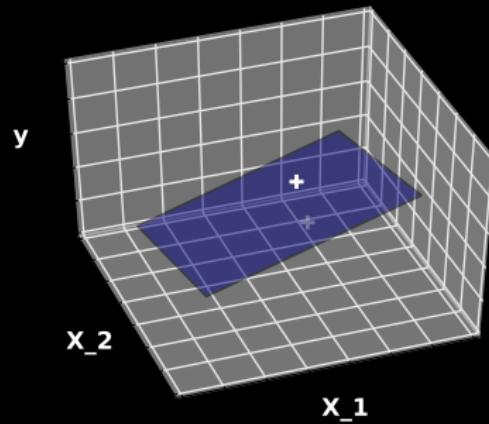
More parameters

# 1 Machine learning in a nutshell: regression

A single descriptor:  
one dimension



Two descriptors:  
2 dimensions

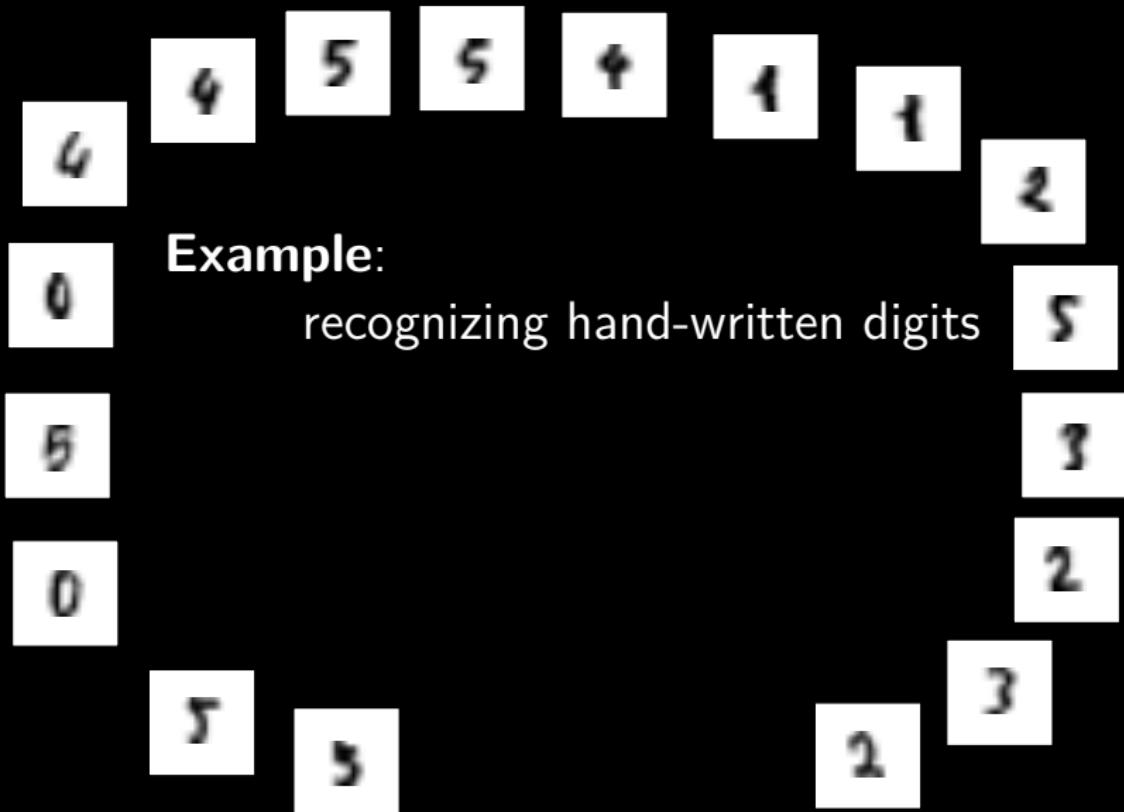


More parameters

⇒ need more data

***“curse of dimensionality”***

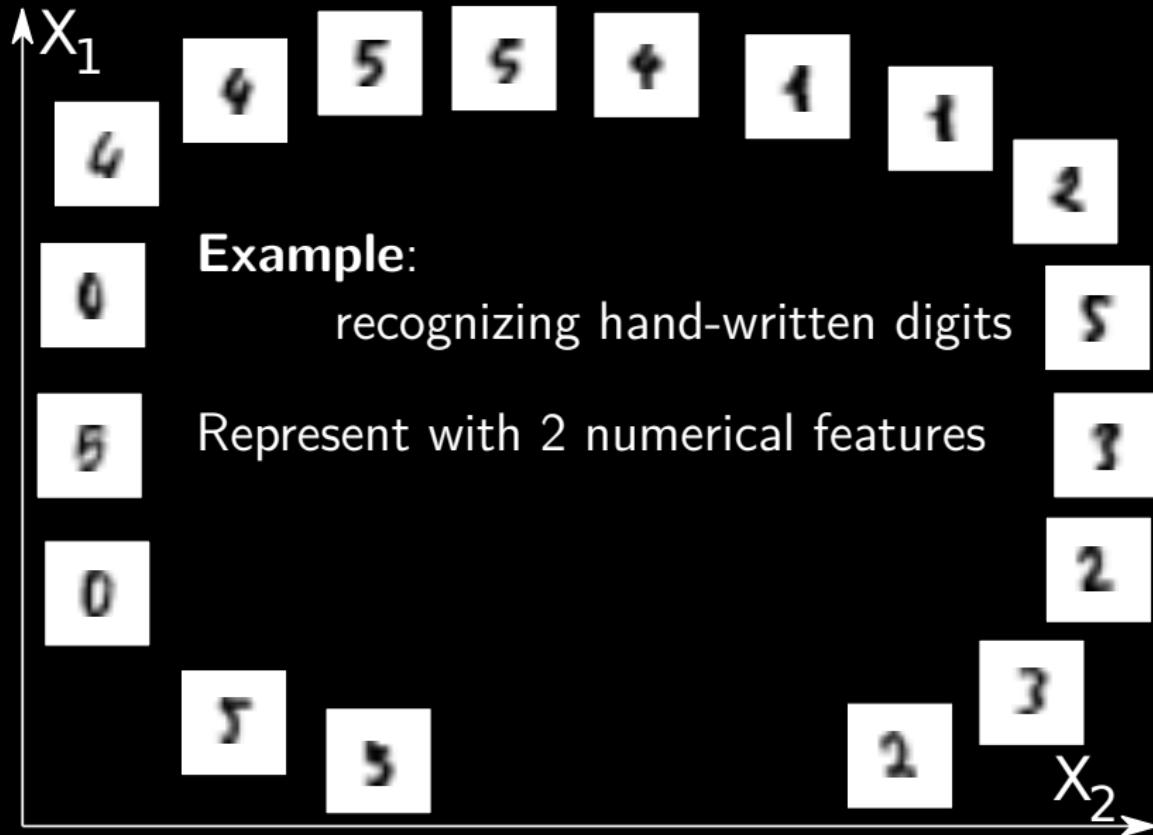
# 1 Machine learning in a nutshell: classification



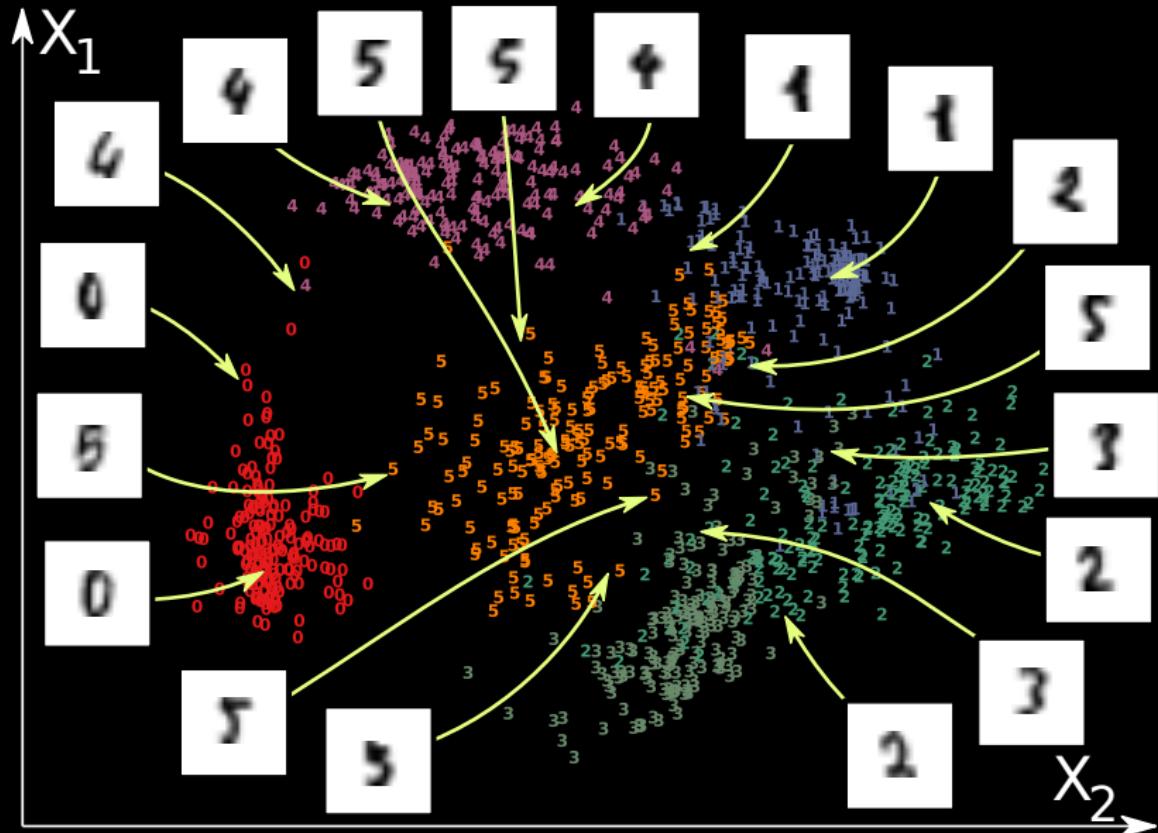
**Example:**

recognizing hand-written digits

# 1 Machine learning in a nutshell: classification

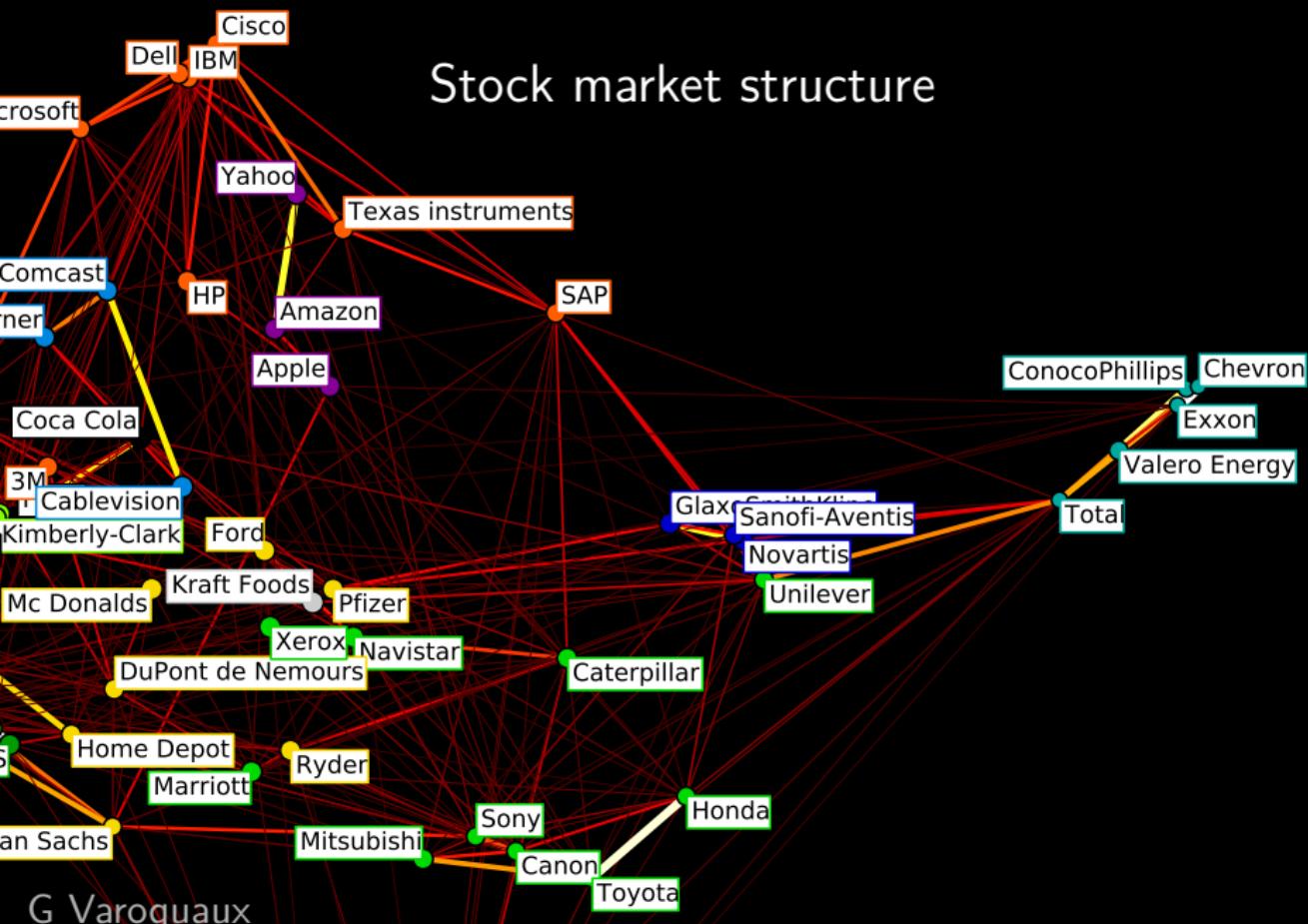


# 1 Machine learning in a nutshell: classification

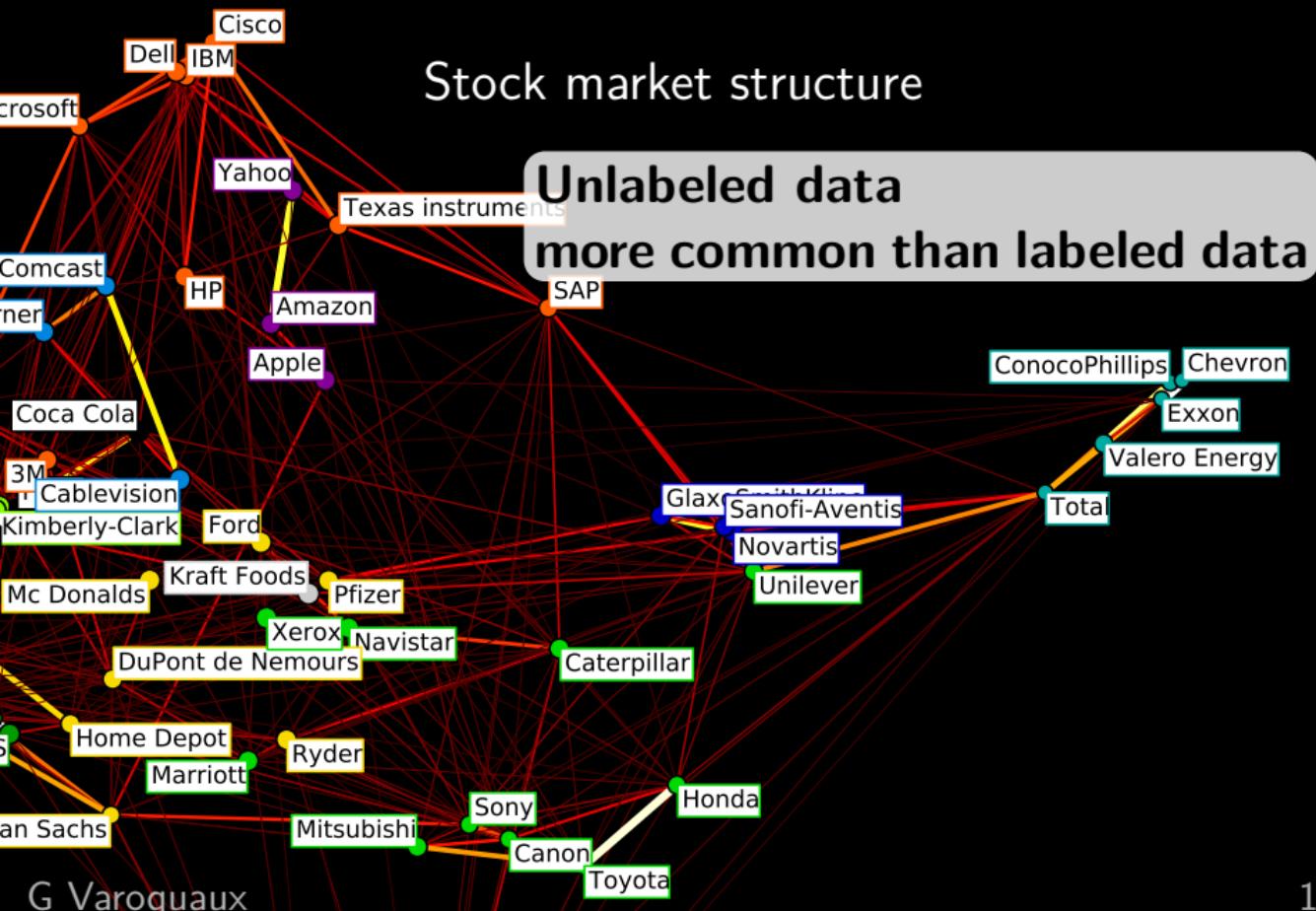


# 1 Machine learning in a nutshell: unsupervised

## Stock market structure

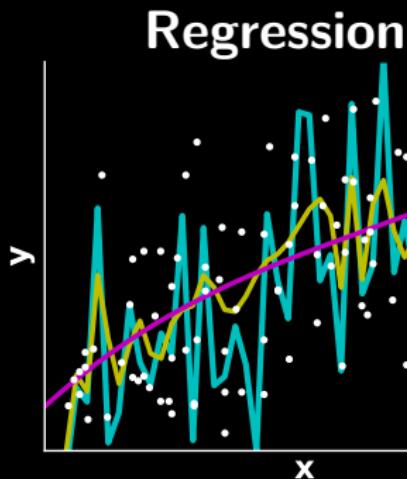


# 1 Machine learning in a nutshell: unsupervised



# Machine learning

Mathematics and algorithms for fitting predictive models



### Classification



**Notions of overfit and test error**

# Machine learning is everywhere

- Image recognition
- Marketing (click-through rate)
- Movie / music recommendation
- Medical data
- Logistic chains (eg supermarkets)
- Language translation
- Detecting industrial failures



# We built a machine learning library

In Python, in 2010



# Real statisticians use R

- And real astronomers use IRAF
- Real economists use Gauss
- Real coders use C assembler
- Real experiments are controlled in Labview
- Real Bayesians use BUGS stan
- Real text processing is done in Perl
- Real Deep learner is best done with torch (Lua)
- And medical doctors only trust SPSS



## *Python, what else?*

- General purpose
- Interactive language
- Easy to read / write

**Fantastic scientific  
ecosystem**

# 1 scikit-learn goals

Machine learning for all

No specific application domain

No requirements in machine learning

High-quality Pythonic software library

Interfaces designed for users

Community-driven development

BSD licensed, very diverse contributors

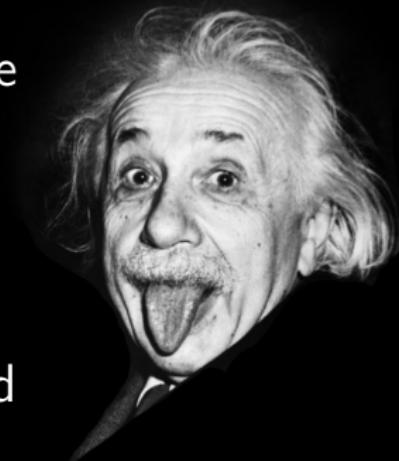
<http://scikit-learn.org>

## Machine learning research

- Conceptual complexity is not an issue
- New and bleeding edge is better
- Simple problems are old science

## In the field

- Tried and tested (aka boring) is good
- Little sophistication from the user
- API is more important than maths



Solving simple problems matters

Solving them really well matters a lot