# Lecture reviews — Week 07

## J.-C. Chappelier & M. Rajman

Laboratoire d'Intelligence Artificielle
Faculté I&C

**EPFL**

# Week 7 keypoints

- supervised/unsupervised

- preprocessing is key

- baseline methods:
  - classification: Naive Bayes, (Logistic regression,) KNN
  - clustering: K-means, dendrograms
  - dim. reduction: PCA, UMAP

- don't forget evaluation keypoints (see lesson 2)

# Week 7 – study case

*[handwritten: 6324]*

*[handwritten: a abacus abbey - - - - cat .... mouse .... Zulu Zygomatic]*

*[handwritten: $D_1$ : 37  0  0 - - - 10 ,.. 12 ... 0  0]*

Some financial company offers you to work on
"*fraud detection using Natural Language Technology applied to client documents*".

① Some preliminary work has already been performed by a former intern
who created document vectors based on an indexing set of 6'324 terms
and reduced them to vectors of size 100 using PCA.

Reviewing his/her work and report, you found a graph related to the corresponding
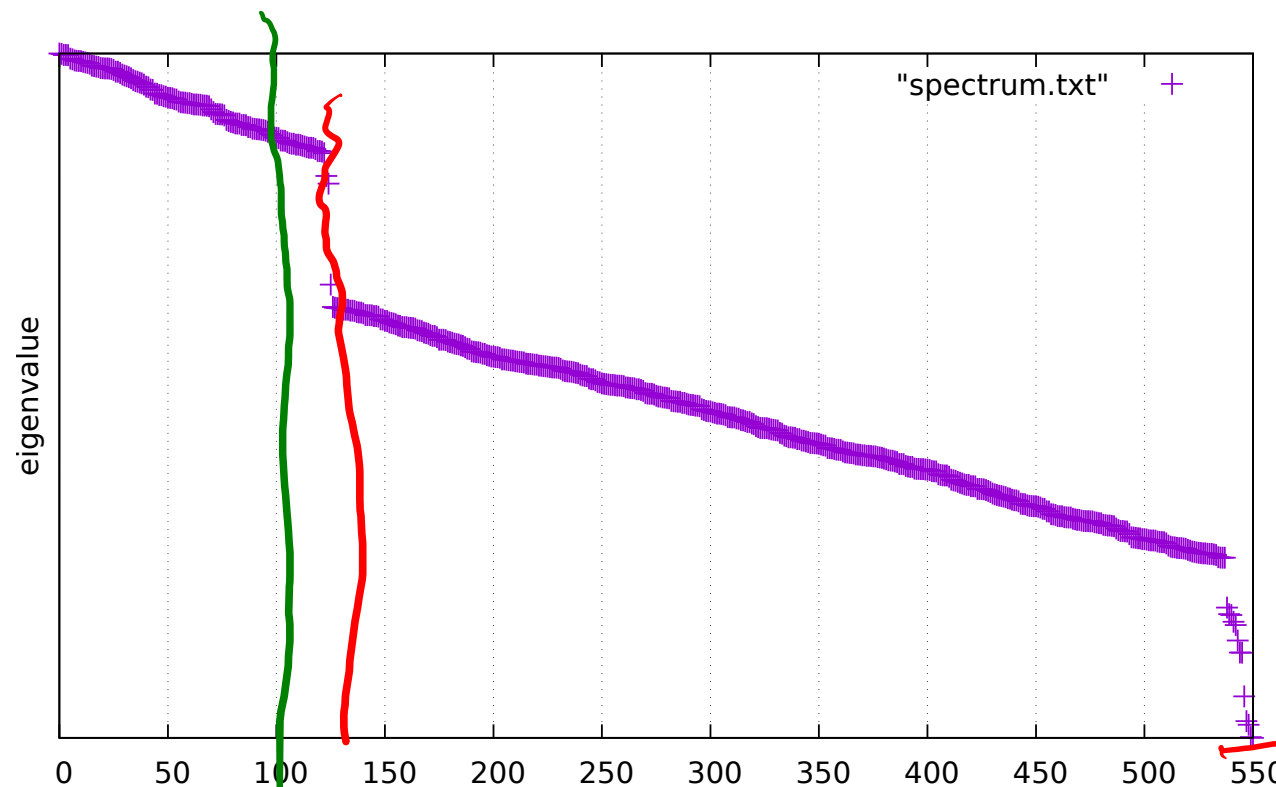singular values.
Next slide shows a (rescaled) zoom on the first 550 left-most points in that graph.
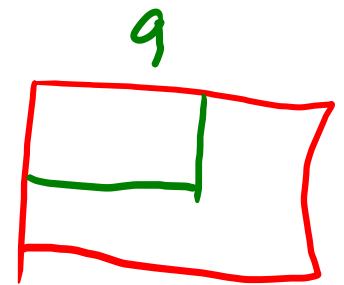
*[handwritten: $D_N$]*

# Week 7 – study case
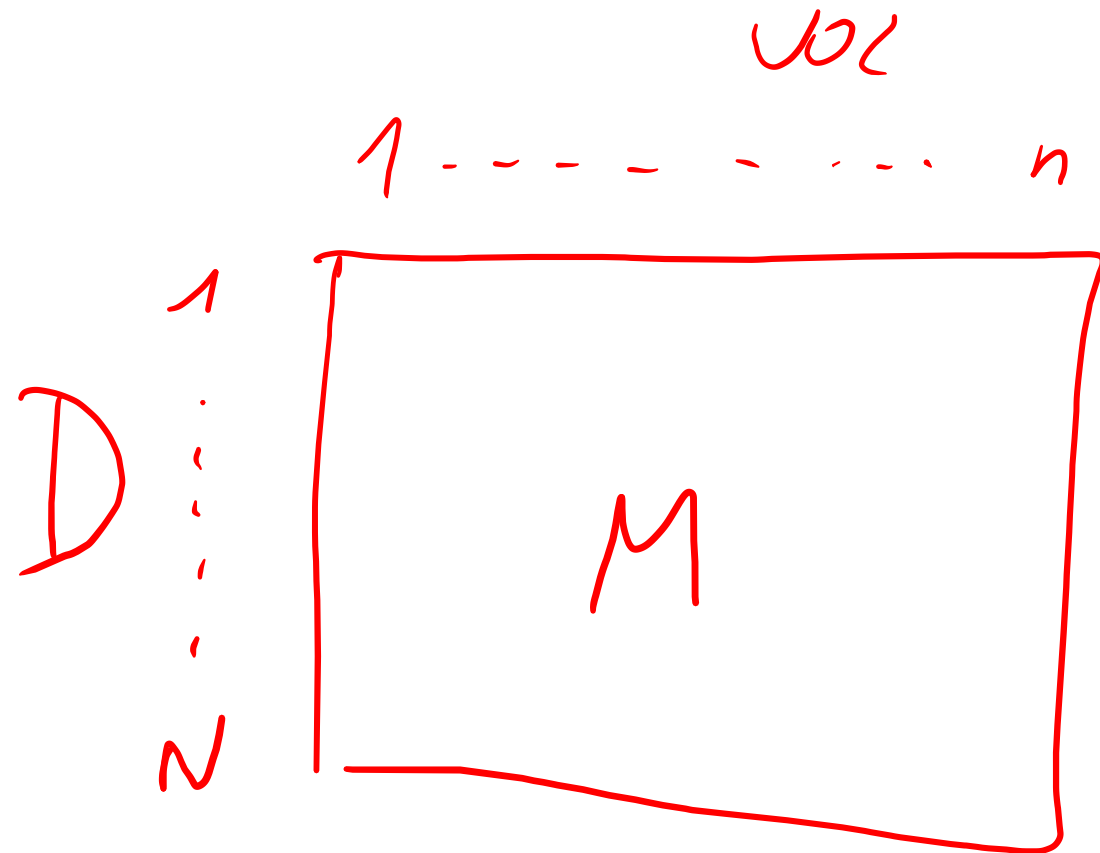


a) What is the abscissa (x-value, horizontal axis) of the right-most point in the original complete graph (not reported here)?

b) What do you think about the intern's methodology for selecting the dimension of the vector space? Would you have performed differently? If yes, how?

$$MM^t \rightarrow U \Lambda V^t$$

# Week 7 – study case

② Before considering more sophisticated Deep-Learning methods, you wisely decide to start with a simple baseline, namely a Naive Bayes model (on the former representation).

*features $f_i$*

*→ document after PCA : vector dim 130*

**a)** Based on your former answer, what is the input of the Naive Bayes module? What is the output? *→ fraud / not fraud* What are the parameters? What is needed for training such a model?

$P(fraud)$

$P(f_i | fraud)$

*+ same with ⌐fraud*

**b)** Concretely, what probability should be computed as an output from the (very simple excerpt of) client document:

*My salary is about 10'000 CHF and I don't pay any tax.*

*1) preprocessing ?*

*2) project ⟶ vector $f_i$ 130*

*3) classif : $P(fraud) \cdot \prod_{i=1}^{130} P(f_i | fraud)$*

*same for ⌐fraud.*

# Week 7 – study case

③ From your first analysis of the baseline results, you realize that single tokens do not adequately capture dependencies that clearly appear at the syntactic level (for instance the one between "*don't*" and "*pay*" in the former example). Using some syntactic parser, you are able to transform the former example sentence

$P(\neg fraud)$

*My salary is about 10'000 CHF and I don't pay any tax.*

into:

$P(f_i' \mid fraud)$

*SALARY-10K-RANGE* / *not_pay* / *tax*     $f'$

**a)** What probability would then be computed as the resulting output by the Naive Bayes model in such a case?

**b)** Compared to former Naive Bayes model, what is the main fundamental reason why you can reasonably expect the results to be better?

$$P(c) \cdot \prod_{i=1}^{3} P(f_i' \mid c)$$

$c \in \{ fraud, \neg fraud \}$