



**University of
Zurich** ^{UZH}

Institute of Computational Linguistics

Introduction to Machine Learning

Lesson 5

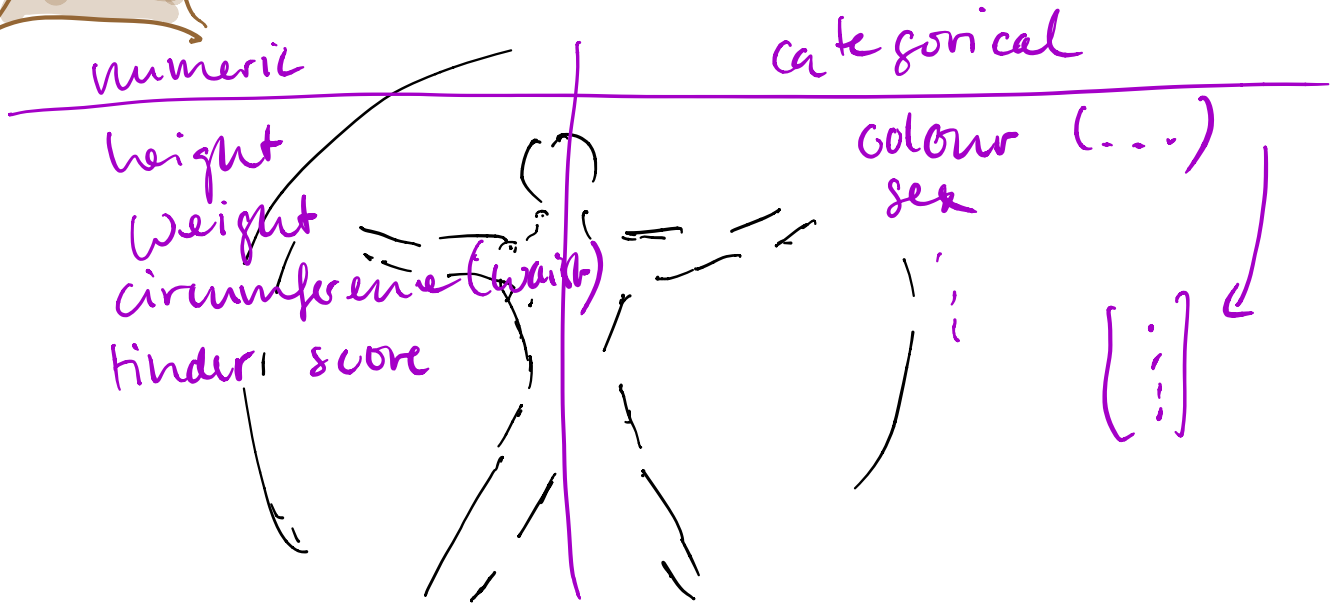
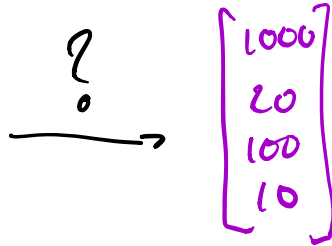
Mathias Müller, Phillip Ströbel

Now

- Feature extraction
- Feature extraction and preprocessing for text

Features?

- Features *represent* objects



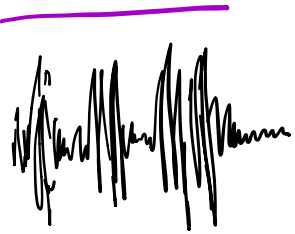
Feature extraction



Feature extraction



$$\begin{bmatrix} 1050 \\ 20 \\ 60 \end{bmatrix}$$



$$\begin{bmatrix} 100 \\ 2 \\ 5 \end{bmatrix}$$

and text?

"Super kal: fragil lithisch exzialigefisch"



$$\begin{bmatrix} 2 \\ 8 \\ 7 \end{bmatrix}$$

Why is raw text not a good representation?

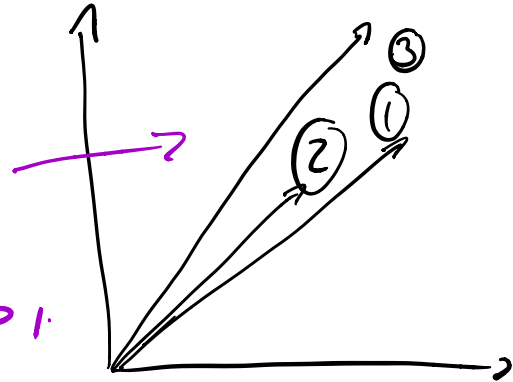
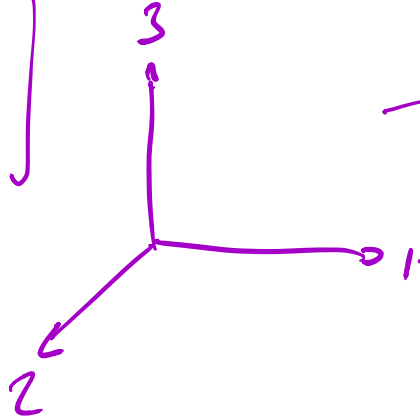
concepts
→ core concepts (all)

• (What are good features, then?)

- ① The game of life is a game of everlasting learning.
- ② The unexamined life is not worth living!
- ③ Never stop learning.

①

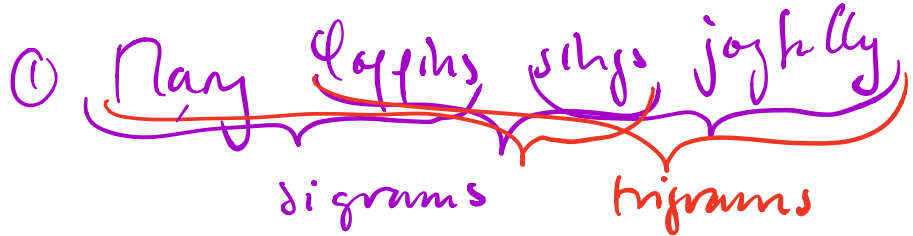
$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$



Kinds of features from text

Based on counts:

- word counts
- ngram counts
- TF/IDF counts



after preprocessing!

?

Kinds of preprocessing for text features

- normalization
- tokenization *urls, aren't*
- lemmatization? stemming?
- some linguistic analysis? *→ parsing, POS tags*
- stopwords *⇒ function words list*

① *DET* The | *NN* game | *det* of | *NN* life | *VERB* is | a | game | of | everlasting | learning ✓
subj *root*

② The *unexamined* life *is* (not) *worth* *living* ✓
stemming → *liv*

③ *Never* *stop* *learning* ✓
lemmatization → *live*

Frequency counts (“vectorization”)

- count how often each word appears in a text

- ① The game of life is a game of everlasting learning.
- ② The unexamined life is not worth living.
- ③ Never stop learning.

tokens = 20

index: The game of life is a everlasting learning unexamined not worth living never stop

①	[1	2	2	1	1	1]
②	[1	0	0	1	1	0]
③	[0]

$|V| = 14$

TF/IDF vectorization

- importance weighting for frequency counts

- ① The game of life is a game of everlasting learning.
- ② The unexamined life is not worth living!
- ③ Never stop learning.

TF

$$\frac{\# \text{ of occurrences}}{\text{length of document}}$$

- how often a term occurs in a document

① The game of life is a everlasting learning.
[0.1 0.2 0.2 0.1 0.1 0.1 0.1 0.1 ... 0]

② The unexamined life is not worth living.
[0.14 _____ 1 ... 0]

③ Never stop learning.
[0.3 0.3 0.3 ...]

IDF

$$\log \left(\frac{\# \text{ of documents}}{\# \text{ of documents where } t \text{ occurs}} \right)$$

- inverse of how many documents contain the term

$$\Rightarrow \text{learning} \\ \text{IDF}(\text{learning}) = \log \left(\frac{3}{2} \right) = 0.4$$

TF/IDF scores

- multiply TF and IDF

$$\textcircled{1} \quad \text{TF}(\text{learning}) = 0.1 \cdot 0.4 = 0.04$$

$$\textcircled{3} \quad \text{TF}(\text{learning}) = 0.3 \cdot 0.4 = 0.12$$

$$\text{IDF}(\text{learning}) = 0.4$$

in doc $\textcircled{3}$ learning is more important!

because $0.12 > 0.04$

$\textcircled{1}$ The game life is a everlasting learning.

[TF/IDF —————]

~~0.1~~
↓
0.04

TF/IDF (cont.)

Sequential nature / order preserving

- count features discard ordering

(A) "Mary Loppins sings joyfully"
= "Loppins joyfully Mary sings"

n-grams

- ① The game of life is a game of everlasting learning.
- ② The unexamined life is not worth living!
- ③ Never stop learning.

Summary

- Feature extraction is a necessary step for text data
- text is frequently **normalized** as a preprocessing step before extracting features
- **vectorization** is a common feature extraction technique

Practical

- Notebook 5.ipynb