

Machine Learning Course
basic track

Machine Learning

Lecture 8: Feature engineering & feature importances

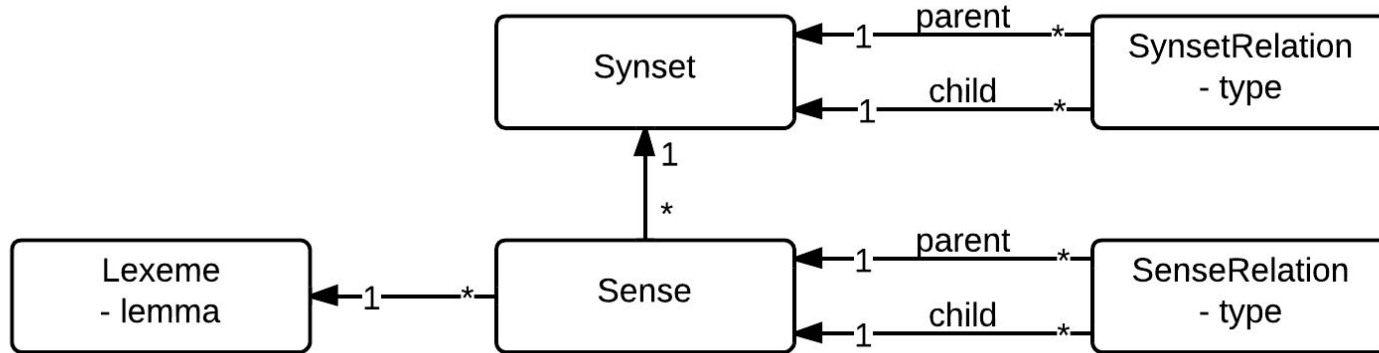
Harbour.Space University
February 2020

Radoslav Neychev

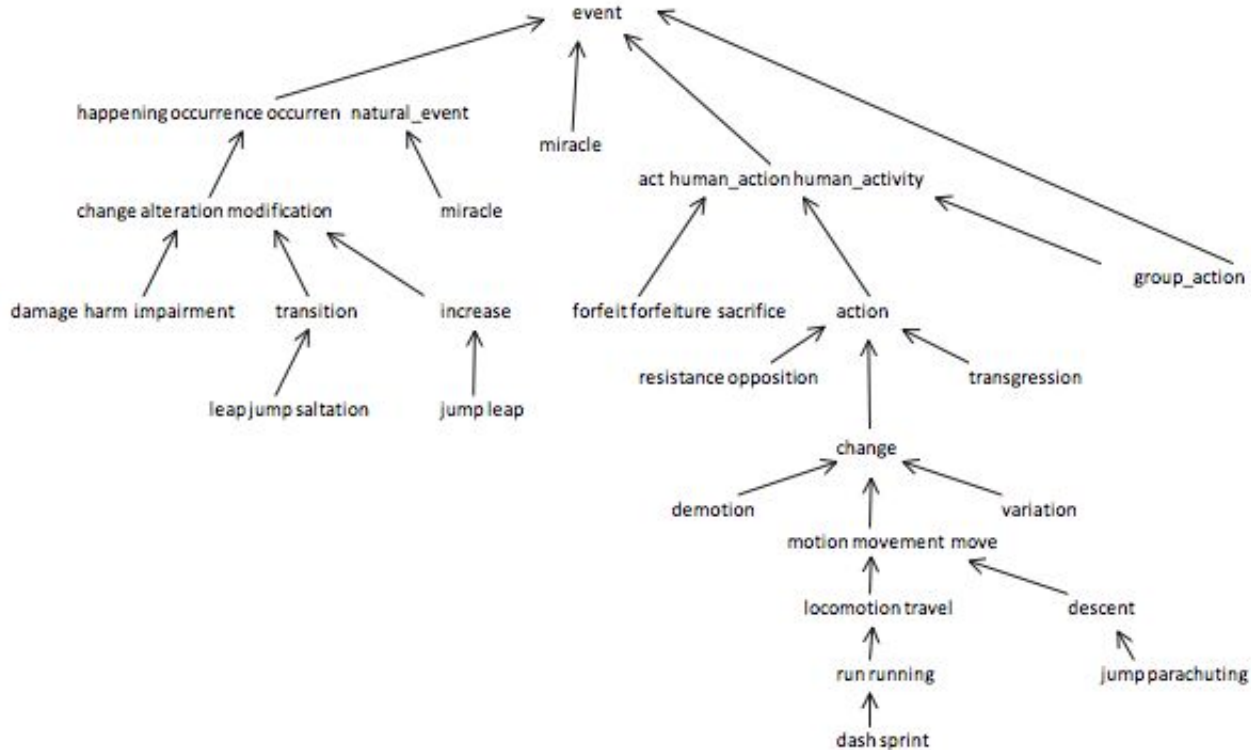
1. Word representations and categorical features
2. Missing data
3. Feature importances estimation

How to represent text in a computer?

Use a taxonomy like WordNet that has hypernyms (is-a) relationships and synonym sets



How to represent text in a computer: WordNet



Discrete representations: problems

- Missing new words
- Subjective
- Requires human labor to create and adapt
- Hard to compute accurate word similarity

Discrete representations: one-hot encoding

"a"	"abbreviations"		"zoology"	"zoom"
1	0		0	0
0	1		0	1
0	0		0	0
⋮	⋮	⋮	⋮	⋮
0	0		0	0
0	0		1	0
0	0		0	1

$$s(Q, D) = \sum_w tf_{w,Q} \cdot \frac{tf_{w,D}}{tf_{w,D} + \frac{k|D|}{avg|D|}} \cdot \log \frac{|C|}{df_w}$$

If word is repeated in the query, it's probably important
 Repetitions of query words in the document → good
 Rare words more important
 The more query words we match, the better. Σ over the vocabulary
 Repetitions of same word less important than different words. Except in very long documents

TF - term frequency

IDF - Inversed Document Frequency

TF-IDF: make it simple

$$\text{tf}(\text{"this"}, d_1) = \frac{1}{5} = 0.2$$

$$\text{tf}(\text{"this"}, d_2) = \frac{1}{7} \approx 0.14$$

$$\text{idf}(\text{"this"}, D) = \log\left(\frac{2}{2}\right) = 0$$



$$\text{tfidf}(\text{"this"}, d_1, D) = 0.2 \times 0 = 0$$

$$\text{tfidf}(\text{"this"}, d_2, D) = 0.14 \times 0 = 0$$



Word 'this' is not very
informative

Document 1		Document 2	
Term	Term Count	Term	Term Count
this	1	this	1
is	1	is	1
a	2	another	2
sample	1	example	3

One of the most successful ideas of statistical NLP:

“You shall know a word by the company it keeps”

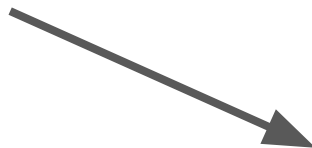
(J. R. Firth 1957: 11)

Words cooccurrences

Finding N-grams in a text



Word-document
cooccurrence matrix

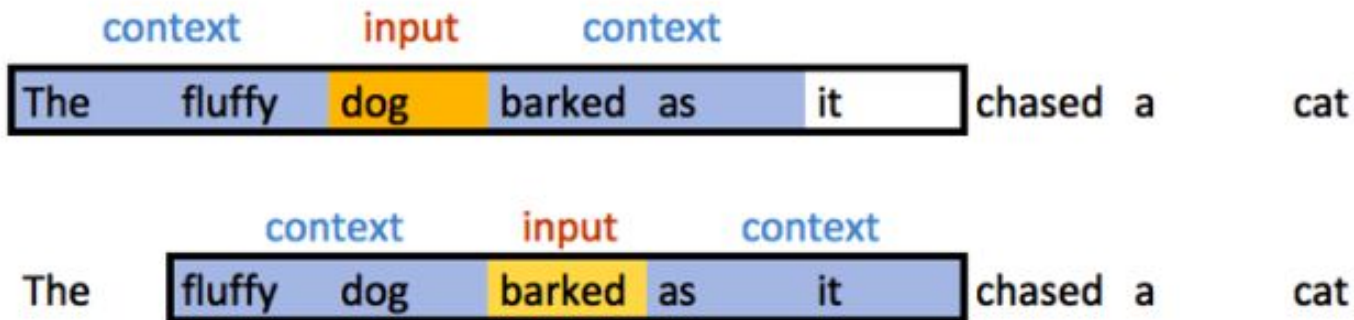


Window around
each word

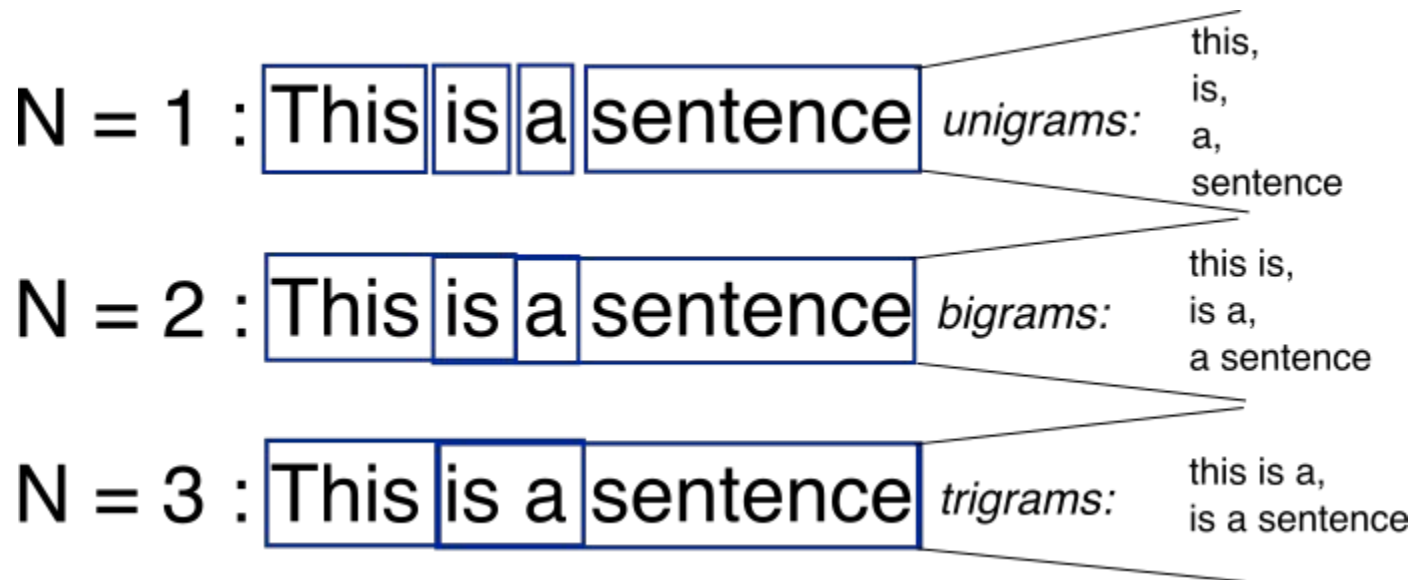
Word-document cooccurrence matrix

$$X = \begin{matrix} & \begin{matrix} I & like & enjoy & deep & learning & NLP & flying & . \end{matrix} \\ \begin{matrix} I \\ like \\ enjoy \\ deep \\ learning \\ NLP \\ flying \\ . \end{matrix} & \left[\begin{array}{cccccccc} 0 & 2 & 1 & 0 & 0 & 0 & 0 & 0 \\ 2 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \end{array} \right] \end{matrix}$$

Words cooccurrences: sliding window



Words cooccurrences: n-grams



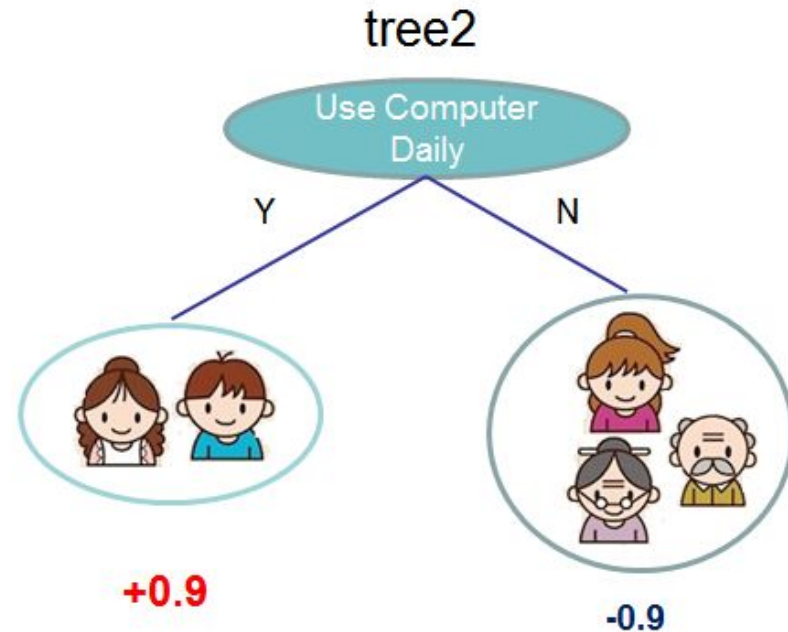
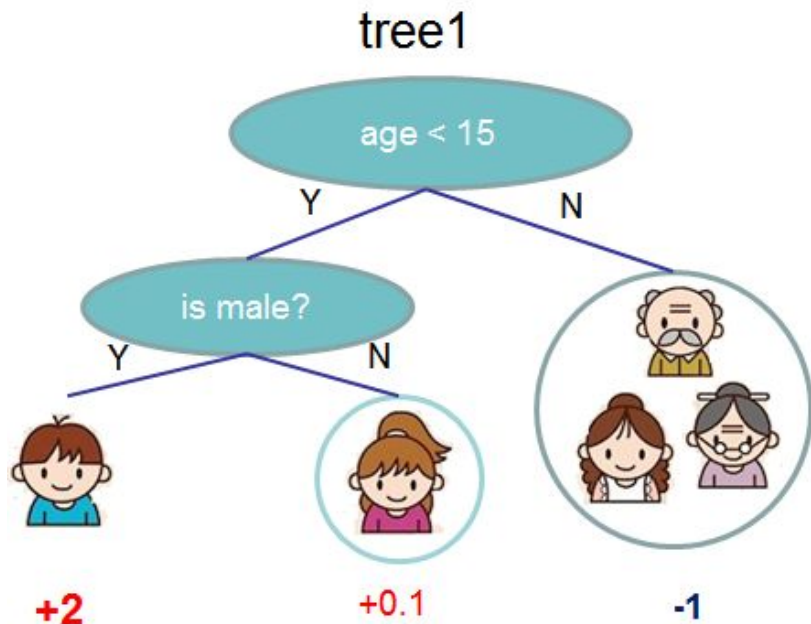
Cooccurrence vectors: problems

- Increase in size with vocabulary
- Very high dimensional: require a lot of storage
- Subsequent classification models have sparsity issues



Models are less robust

Feature importance estimation



$$f(\text{boy icon}) = 2 + 0.9 = 2.9$$

$$f(\text{old man icon}) = -1 - 0.9 = -1.9$$

Feature importance estimation

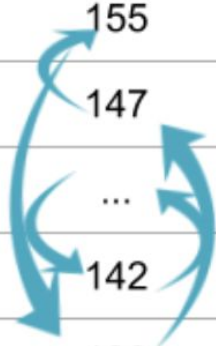
1. Permutation importance
2. Partial Dependence Plots (PDP)
3. Tree specific:
 - a. Gain
 - b. Frequency (Split Count)
 - c. Cover (weighted Split Count)
4. Shap

Permutation importance

Height at age 20 (cm)	Height at age 10 (cm)	...	Socks owned at age 10
182	155	...	20
175	147	...	10
...
156	142	...	8
153	130	...	24

Permutation importance

Height at age 20 (cm)	Height at age 10 (cm)	...	Socks owned at age 10
182	155	...	20
175	147	...	10
...
156	142	...	8
153	130	...	24



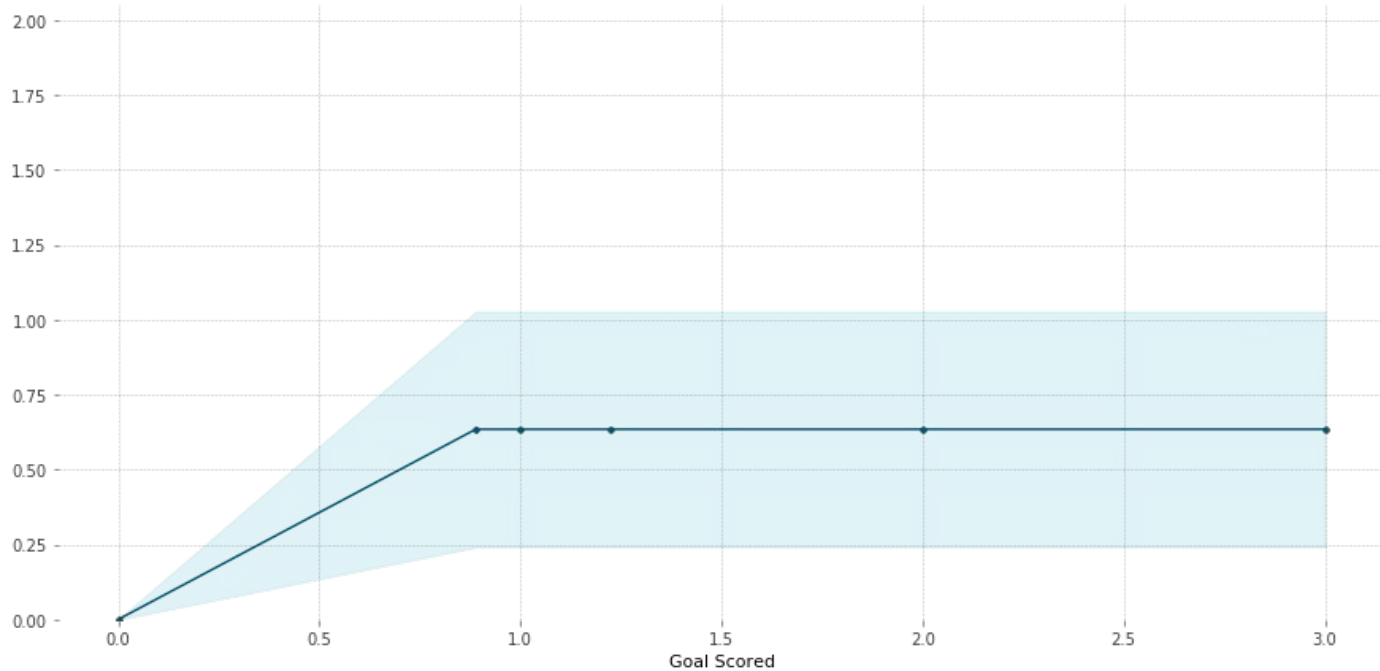
Train model

Observe changes caused by feature random permutations

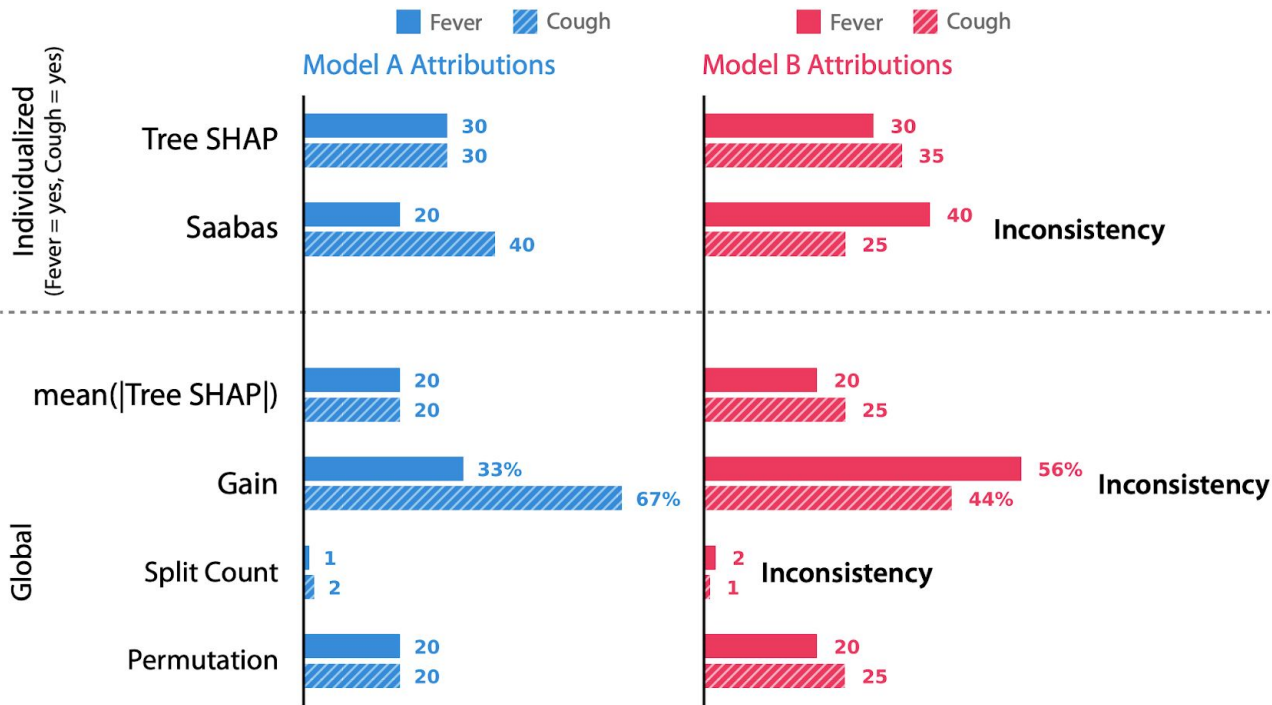
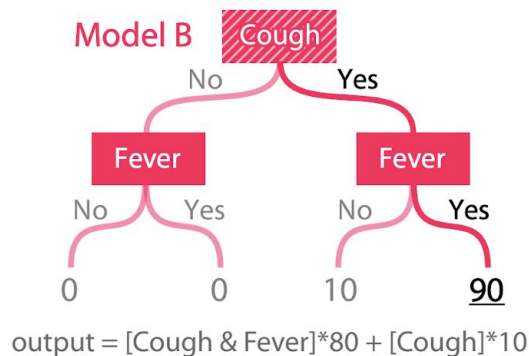
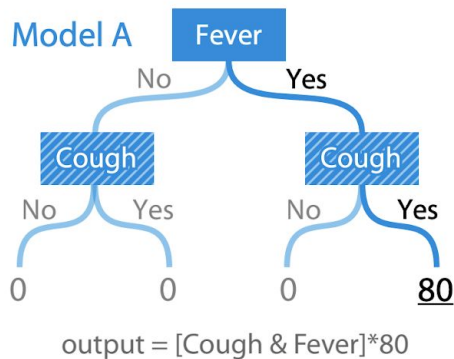
Partial Dependence Plots

PDP for feature "Goal Scored"

Number of unique grid points: 6



Importance estimation problems



Consider i -th feature. Shap value will be

$$\phi_i(p) = \sum_{S \subseteq N/\{i\}} \frac{|S|!(n - |S| - 1)!}{n!} (p(S \cup \{i\}) - p(S))$$

where $p(S \cup \{i\})$ is model prediction on feature subset S with i -th feature added.

Consider i -th feature. Shap value will be

$$\phi_i(p) = \sum_{S \subseteq N/\{i\}} \frac{|S|!(n - |S| - 1)!}{n!} (p(S \cup \{i\}) - p(S))$$

where $p(S \cup \{i\})$ is model prediction on feature subset S with i -th feature added.

SHAP values are the only consistent and locally accurate individualized feature attributions

1. Remember the bias-variance decomposition
2. Consider using SHAP values to estimate feature importances.