

Introduction to NLP

Scope of this lab

Not to bombard with NLP terminologies

Not to provide a NLP crash course in an hour!

But

to introduce you to NLP basics

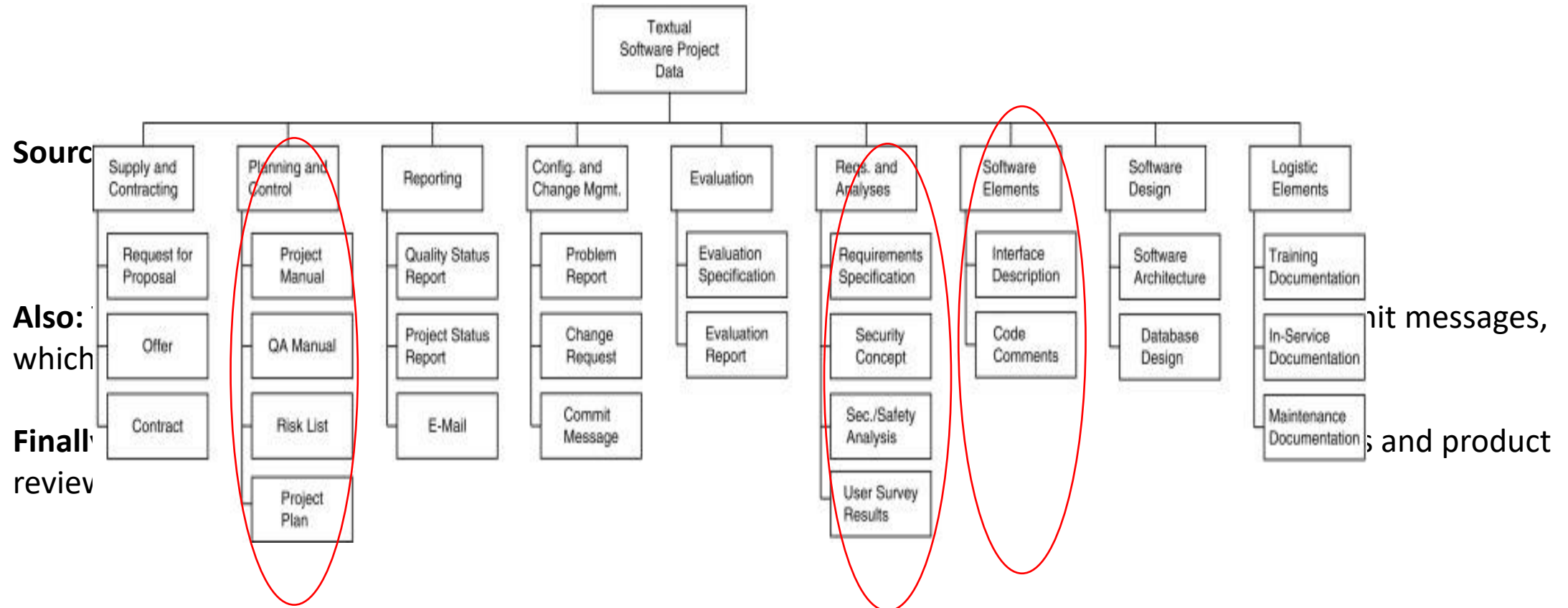
to explain the overall NLP platter that you have for your disposal

to provide a bit of hands-on to kindle your brain with what's in store
and where to get started for your actual work

What are the various industry application of NLP that come to your mind?

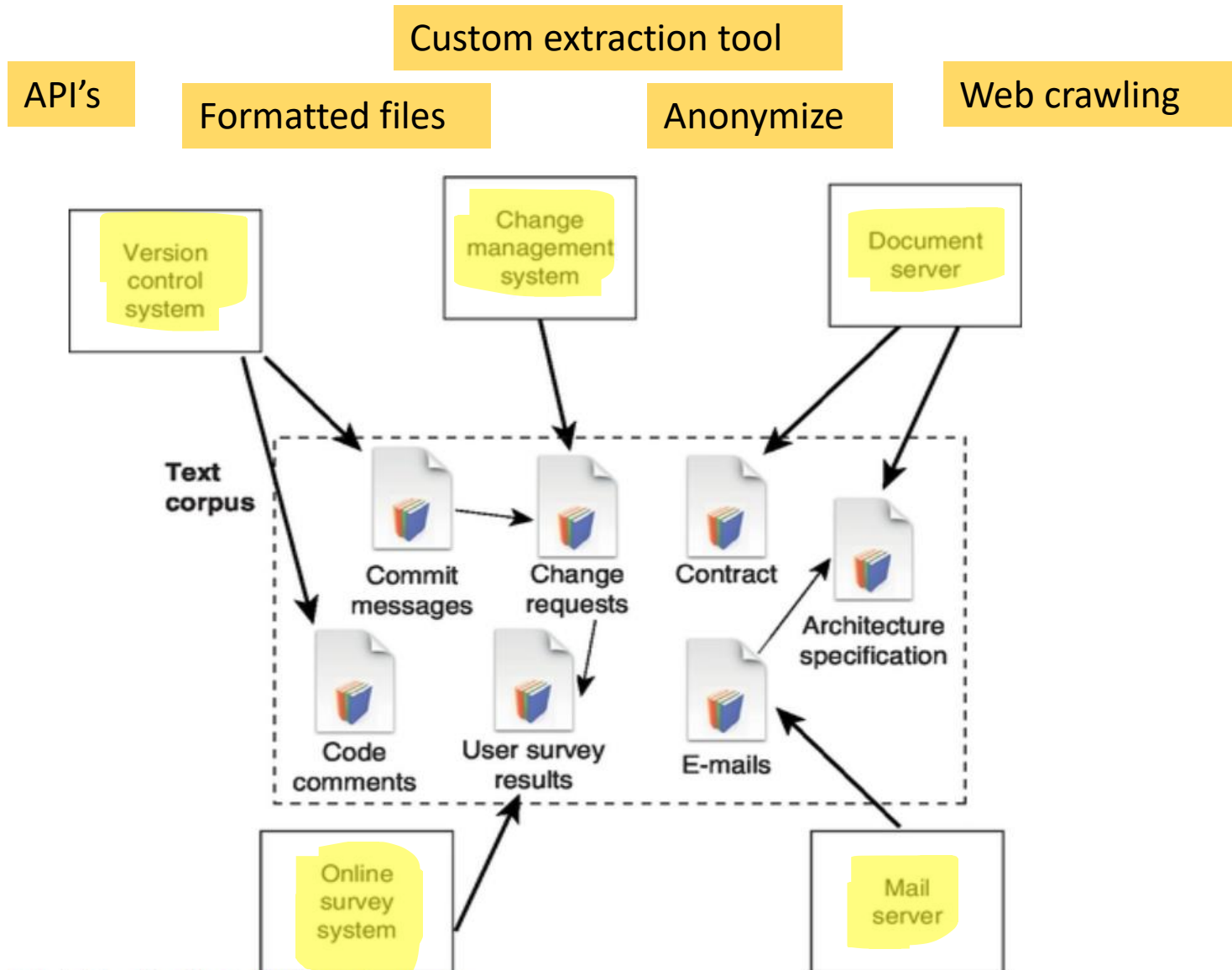
- Search engines
- Advanced text editors: Such as Grammarly.com
- Computational advertising
- Fraud detection
- Sentiment analysis
- Opinion mining
- Text summarization
- Context analysis

Textual data in Software Project





Textual Software Project Data and Retrieval [1]



3.2 Text collection from different sources.

Importance of NLP in Software Engineering

- Text Retrieval (TR) and NLP in software is one of the fastest growing areas of research in SE. [2]
- Exposing the SE community to these techniques and their applications in SE would help to fill a gap in their current background and allow them to immediately use TR and NLP to advance their research.
- In particular, for TR, approaches such as Vector Space Model, Latent Semantic Analysis, Latent Dirichlet Association, Language Models will be covered. NLP techniques covered will include part-of-speech tagging, stemming, stopword elimination, semantics analysis, sentiment analysis, etc.

Classification schema of NLP techniques



Part-of-Speech Tagging

The choice Tag Sets depends on the language and application
Example tag set sizes (for English)

- Brown corpus, 87 tags
- Penn treebank 45 tags
- BNC, 61 tags

Tag	Description	Example	Tag	Description	Example
CC	Coordin. Conjunction	<i>and, but, or</i>	SYM	Symbol	<i>+, %, &</i>
CD	Cardinal number	<i>one, two, three</i>	TO	“to”	<i>to</i>
DT	Determiner	<i>a, the</i>	UH	Interjection	<i>ah, oops</i>
EX	Existential ‘there’	<i>there</i>	VB	Verb, base form	<i>eat</i>
FW	Foreign word	<i>mea culpa</i>	VBD	Verb, past tense	<i>ate</i>
IN	Preposition/sub-conj	<i>of, in, by</i>	VBG	Verb, gerund	<i>eating</i>
JJ	Adjective	<i>yellow</i>	VCN	Verb, past participle	<i>eaten</i>
JJR	Adj., comparative	<i>bigger</i>	VBP	Verb, non-3sg pres	<i>eat</i>
JJS	Adj., superlative	<i>wildest</i>	VBZ	Verb, 3sg pres	<i>eats</i>
LS	List item marker	<i>1, 2, One</i>	WDT	Wh-determiner	<i>which, that</i>
MD	Modal	<i>can, should</i>	WP	Wh-pronoun	<i>what, who</i>
NN	Noun, sing. or mass	<i>llama</i>	WP\$	Possessive wh-	<i>whose</i>
NNS	Noun, plural	<i>llamas</i>	WRB	Wh-adverb	<i>how, where</i>
NNP	Proper noun, singular	<i>IBM</i>	\$	Dollar sign	<i>\$</i>
NNPS	Proper noun, plural	<i>Carolinas</i>	#	Pound sign	<i>#</i>
PDT	Predeterminer	<i>all, both</i>	“	Left quote	<i>(‘ or “)</i>
POS	Possessive ending	<i>’s</i>	”	Right quote	<i>(’ or ”)</i>
PRP	Personal pronoun	<i>I, you, he</i>	(Left parenthesis	<i>([, ({ , <)</i>
PRP\$	Possessive pronoun	<i>your, one’s</i>)	Right parenthesis	<i>([, ({ , <)</i>
RB	Adverb	<i>quickly, never</i>	,	Comma	<i>,</i>
RBR	Adverb, comparative	<i>faster</i>	.	Sentence-final punc	<i>(. ! ?)</i>
RBS	Adverb, superlative	<i>fastest</i>	:	Mid-sentence punc	<i>(: ; ... - -)</i>
RP	Particle	<i>up, off</i>			

Source : Penn Tree Bank Tagset - <http://coltekin.net/cagri/courses/snlp2017/slides/pos-tagging.pdf>

Chunking

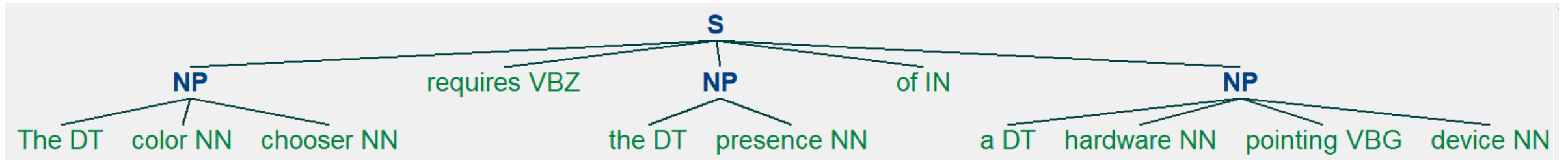
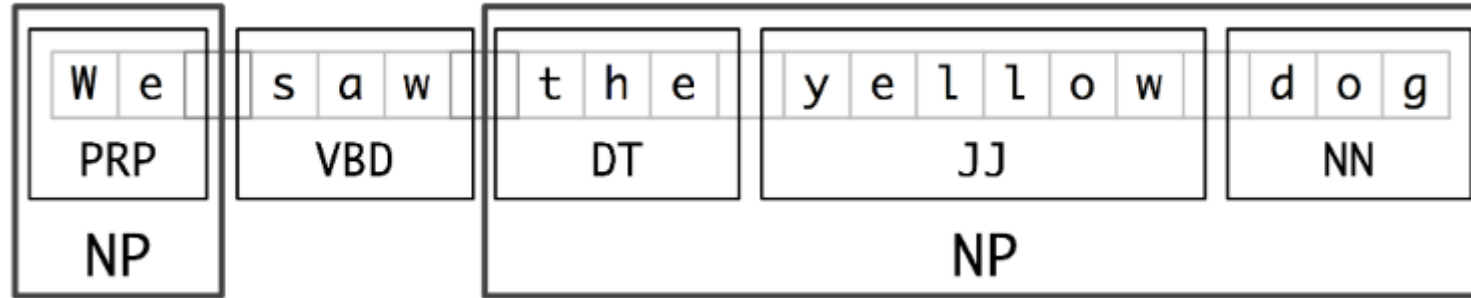
The color chooser requires **the presence** of a **hardware pointing device**

POS Tagging (averaged_perceptron_tagger)

1

```
[('The', 'DT'), ('color', 'NN'), ('chooser', 'NN'), ('requires', 'VBZ'), ('the', 'DT'), ('presence', 'NN'), ('of', 'IN'), ('a', 'DT'), ('hardware', 'NN'), ('pointing', 'VBG'), ('device', 'NN')]
```

```
grammar = r"""NP:{<DT>?<NN><VBG><NN>}
```



Named Entity Recognition

- Wikipedia: **Named-entity recognition (NER)** (also known as **entity identification**, **entity chunking** and **entity extraction**) is a subtask of [information extraction](#) that seeks to locate and classify [named entity](#) mentioned in [unstructured text](#) into pre-defined categories such as person names, organizations, locations, medical codes, time expressions, quantities, monetary values, percentages, etc.

Input: “When Michael Jordan was at the peak of his powers as an NBA superstar, his Chicago Bulls team were moving down the completion, winning six National Basketball Association titles”.

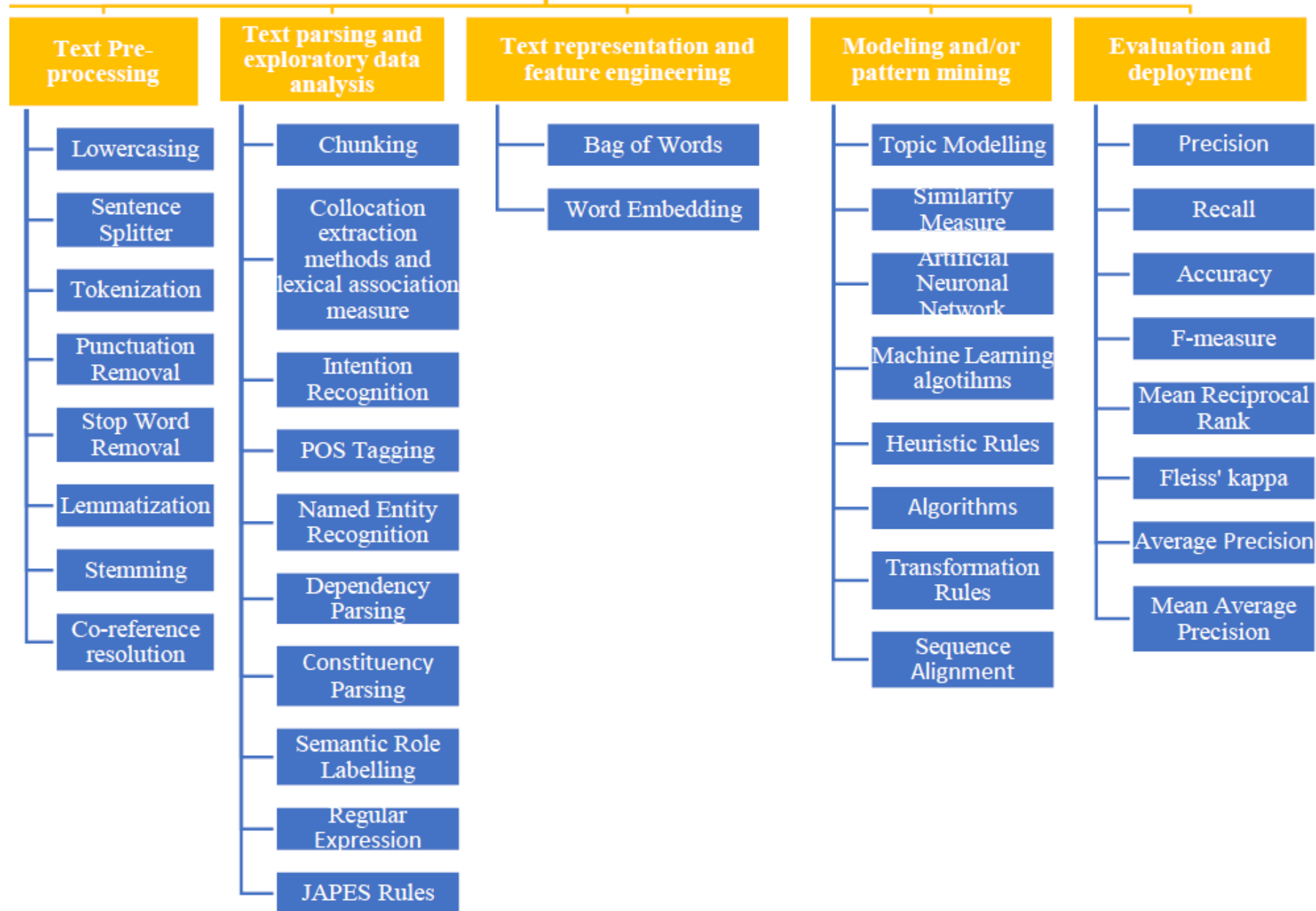
Output:

“Chicago Bulls”

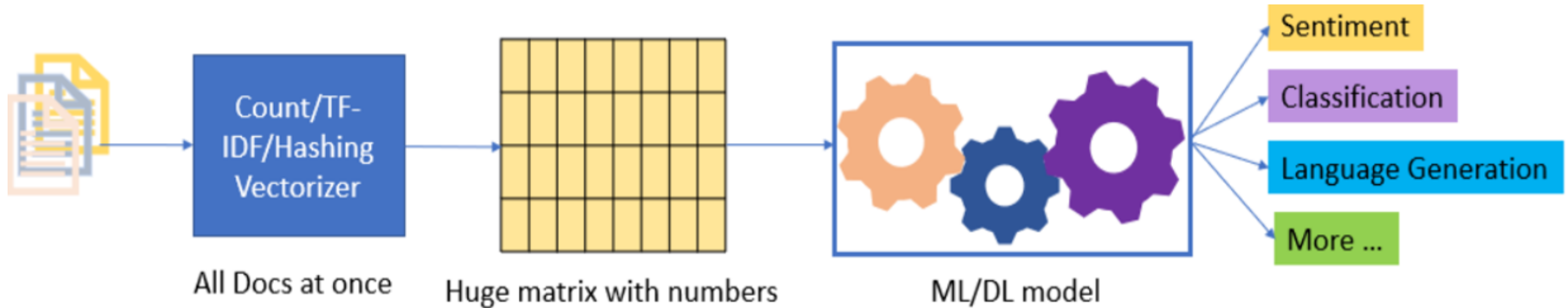
“Michael Jordan”

“National Basketball Association”

Classification schema of NLP techniques



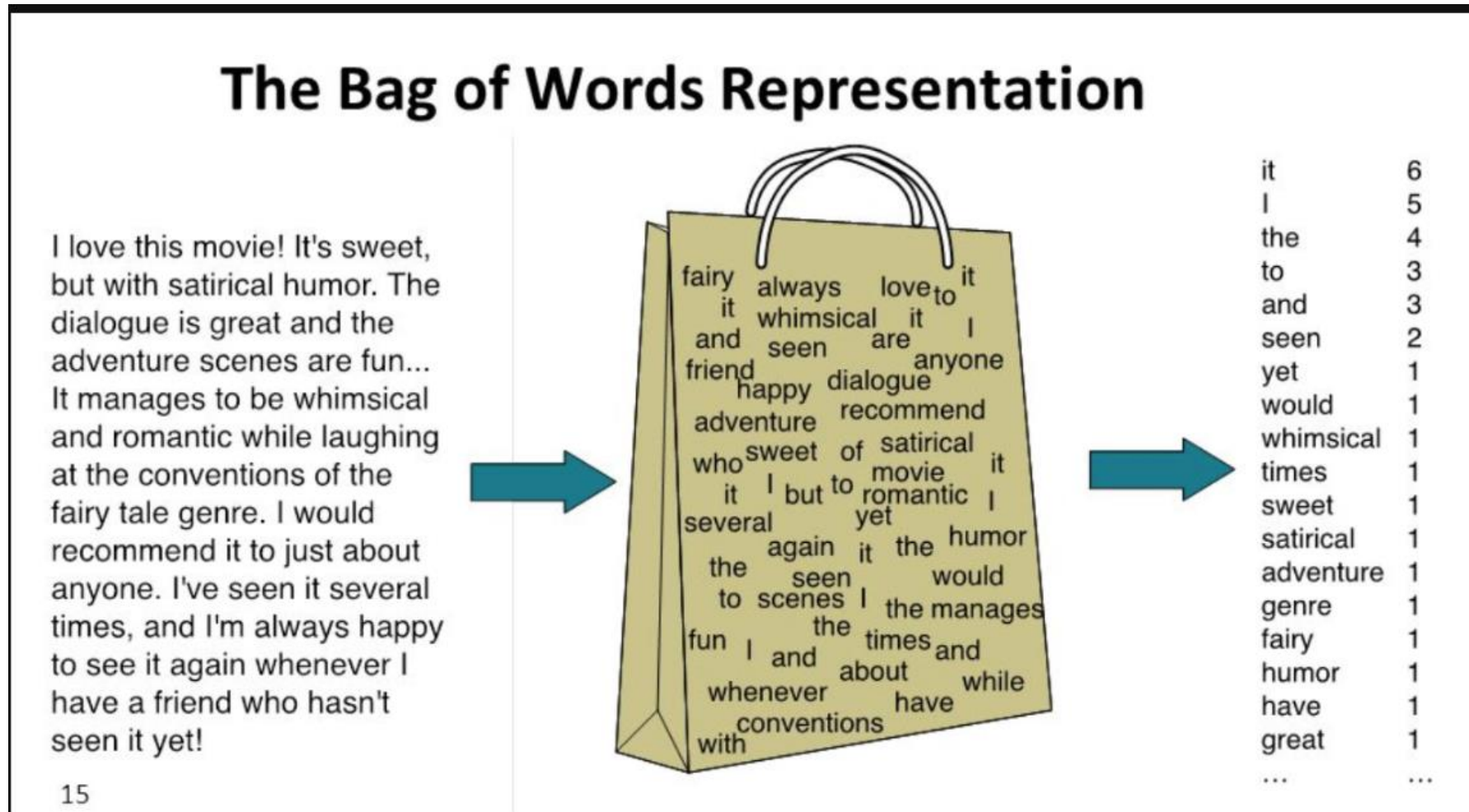
NLP and Machine learning



Source: <https://www.datajango.com/deep-learning-and-nlp/>

Bag of words

Methods that are used for natural language processing to represent documents where the order of words (grammar) is not important.



	MARY	IS	HUNGRY	HAPPY	FOR	APPLES	NOT	JOHN	HE	
"Mary is hungry for apples."	1	1	1	0	1	1	0	0	0	→ [1, 1, 1, 0, 1, 1, 0, 0, 0]
"John is happy he is not hungry for apples."	0	2	1	1	1	1	1	1	1	→ [0, 2, 1, 1, 1, 1, 1, 1, 1]

Source: <https://blog.insightdatascience.com/how-to-solve-90-of-nlp-problems-a-step-by-step-guide-fda605278e4e>

References

1. Bird, Christian, Tim Menzies, and Thomas Zimmermann, eds. The art and science of analyzing software data. Elsevier, 2015.
2. Arnaudova, Venera, et al. "The use of text retrieval and natural language processing in software engineering." Proceedings of the 37th International Conference on Software Engineering-Volume 2. 2015.
3. Juergens, Elmar, et al. "Do code clones matter?." 2009 IEEE 31st International Conference on Software Engineering. IEEE, 2009.
4. Juergens, Elmar, et al. "Can clone detection support quality assessments of requirements specifications?." Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering-Volume 2. ACM, 2010.
5. <https://medium.com/square-corner-blog/topic-modeling-optimizing-for-human-interpretability-48a81f6ce0ed>
6. Stevens, Keith, Kegelmeyer, Philip, Andrzejewski, David, and Buttler, David. Exploring topic coherence over many models and many topics. In EMNLP, 2012.