

Блок

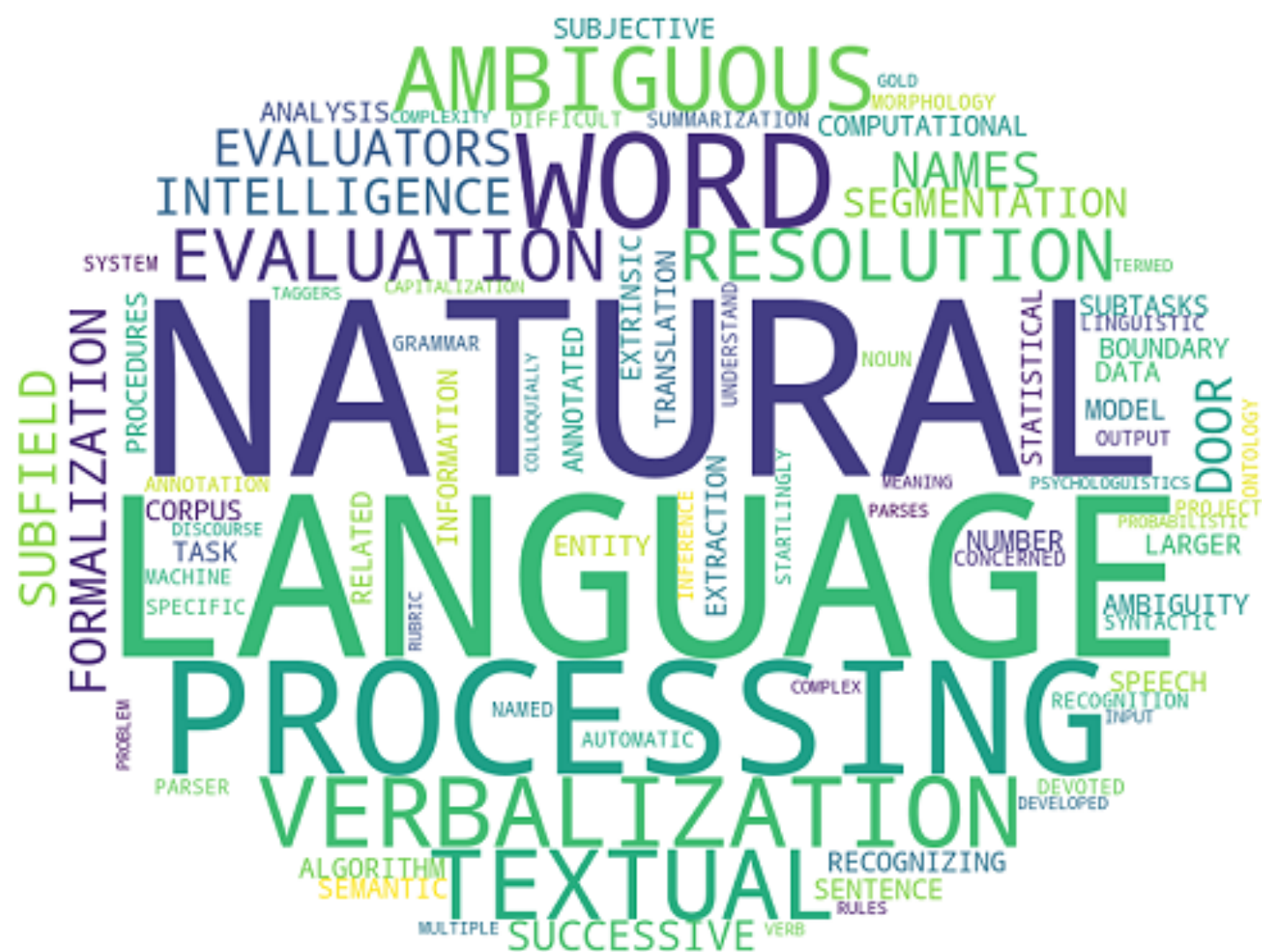
Feature Engineering

Занятие № 6

Работа с текстовыми данными

- Посмотреть на текст с другой стороны
- Узнаем, как текст превращается в математику

Цели занятия



Сколько слов в русском языке?

Сколько слов в русском языке?

Число лемм: ~500k

Число всех слов: >4.5 млн.

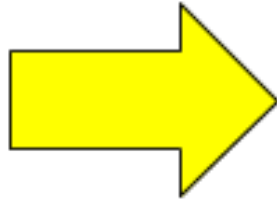
Слова и математика

Слова и математика

- Bag of words
 - Tf-idf
- Word2vec/Glove

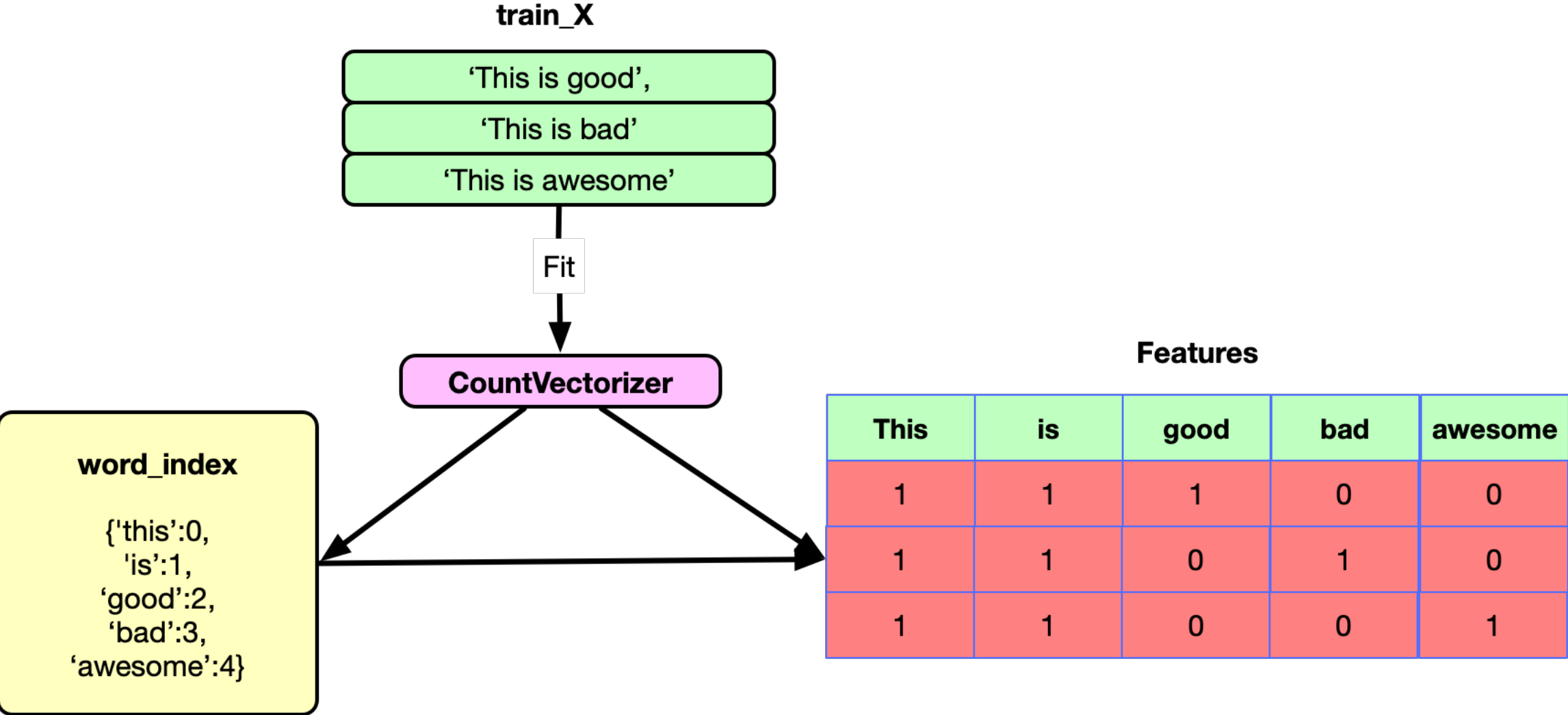
Bag of words

Color
Red
Red
Yellow
Green
Yellow

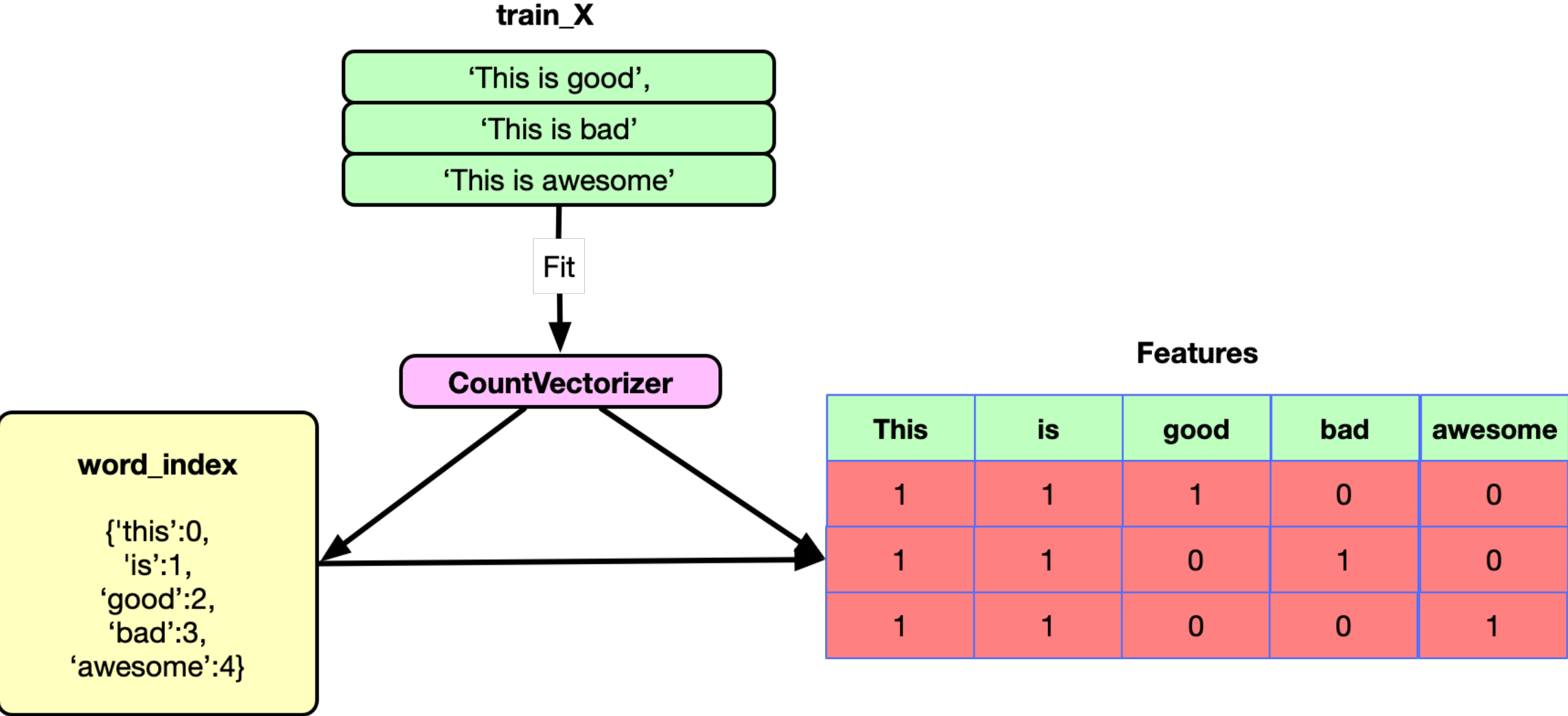


Red	Yellow	Green
1	0	0
1	0	0
0	1	0
0	0	1

Bag of words



Bag of words



Tf-IDF

$$w_{x,y} = tf_{x,y} \times \log \left(\frac{N}{df_x} \right)$$

TF-IDF

Term x within document y

$tf_{x,y}$ = frequency of x in y

df_x = number of documents containing x

N = total number of documents

Tf-IDF

$$w_{x,y} = tf_{x,y} \times \log \left(\frac{N}{df_x} \right)$$

TF-IDF

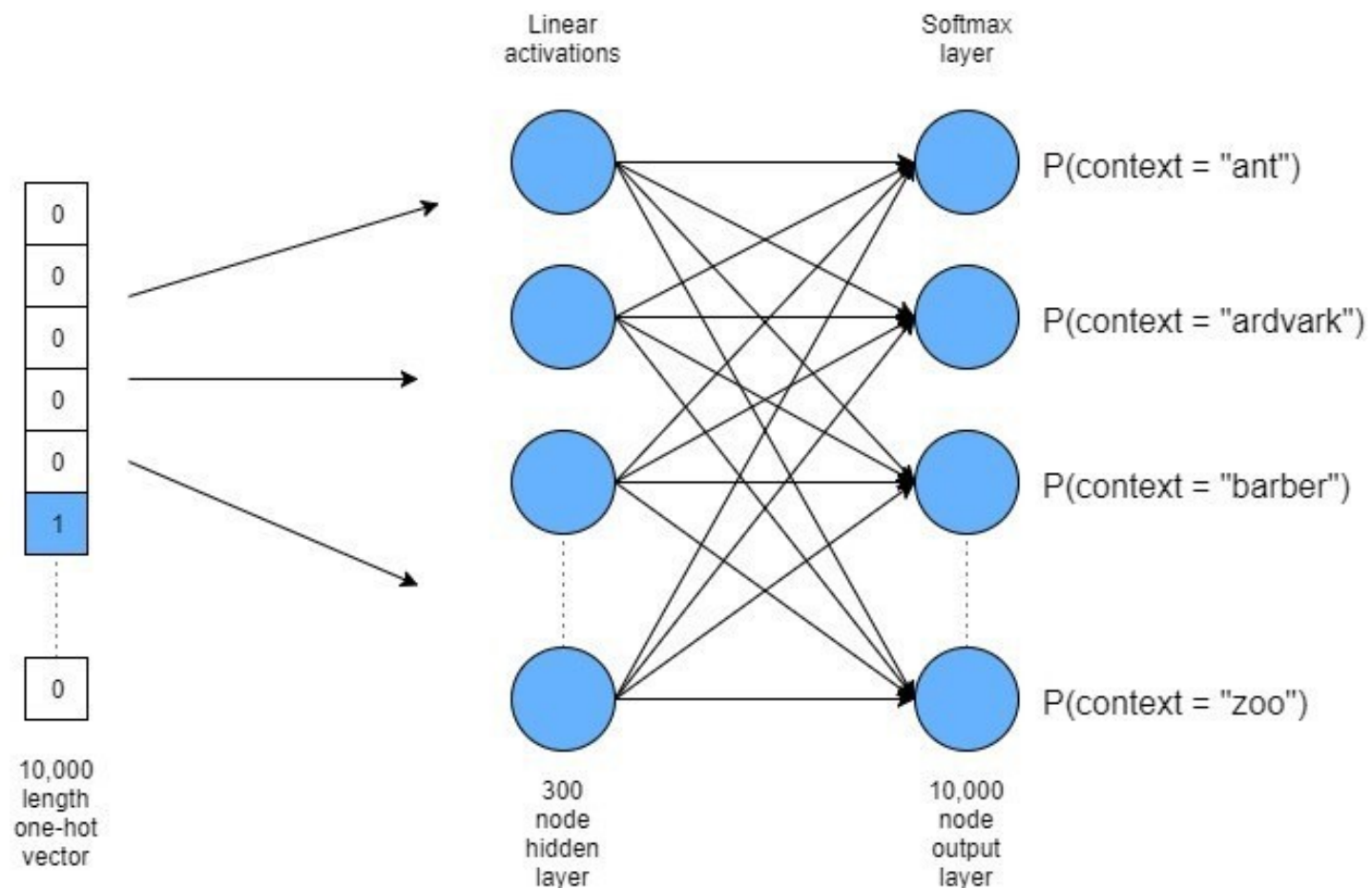
Term x within document y

$tf_{x,y}$ = frequency of x in y

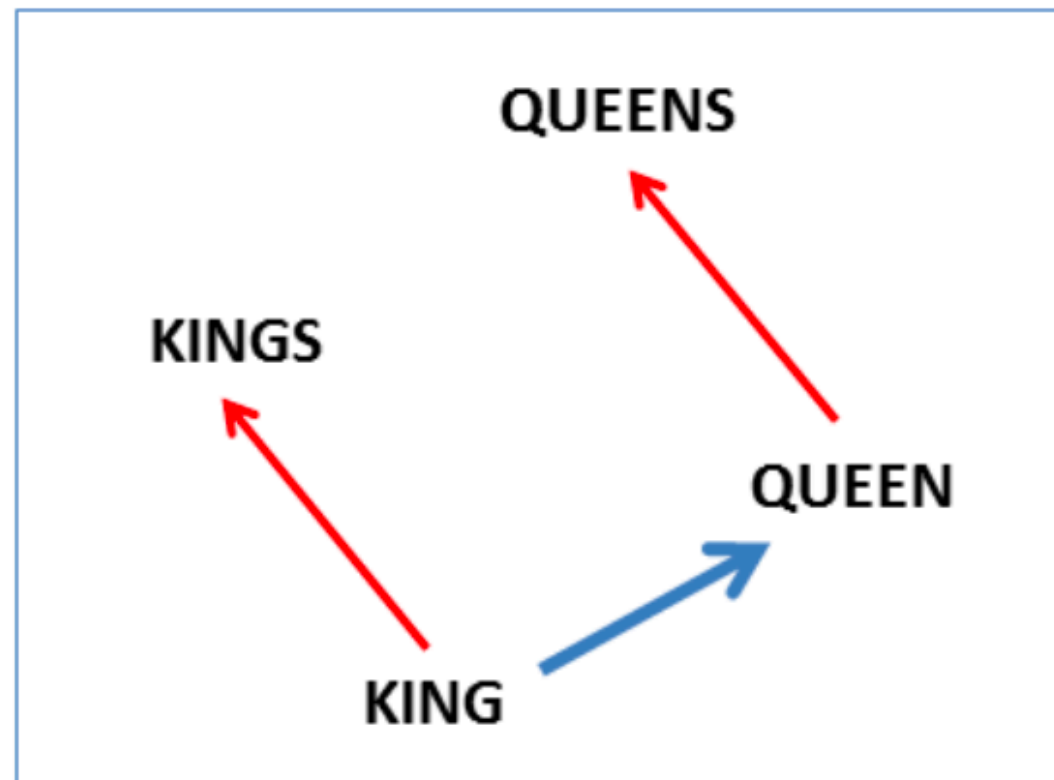
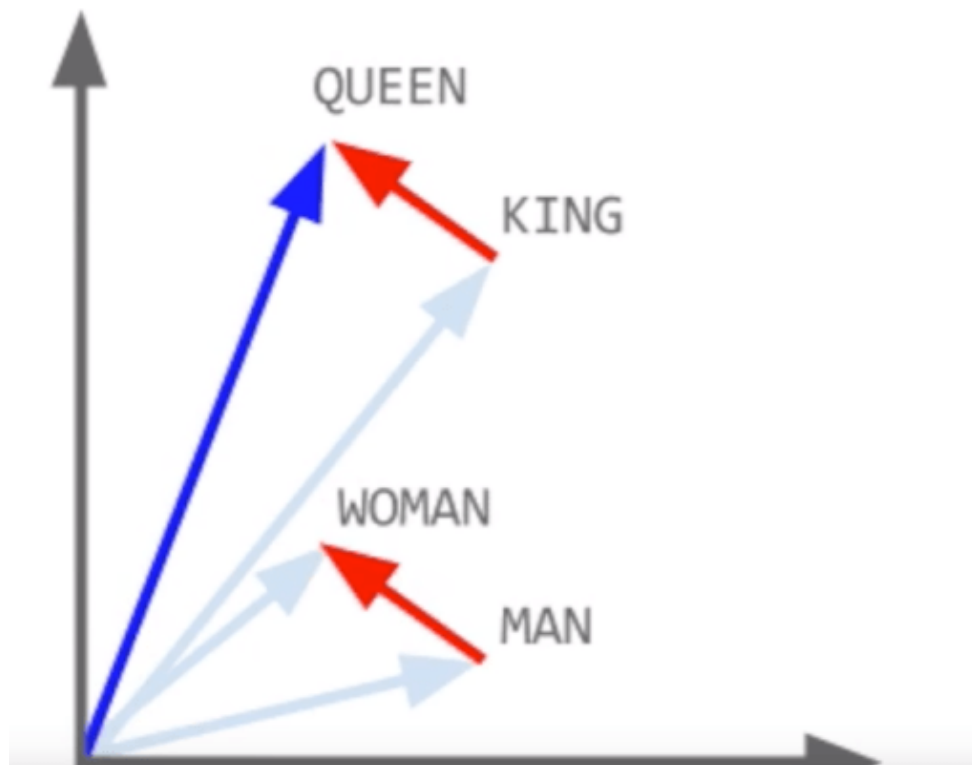
df_x = number of documents containing x

N = total number of documents

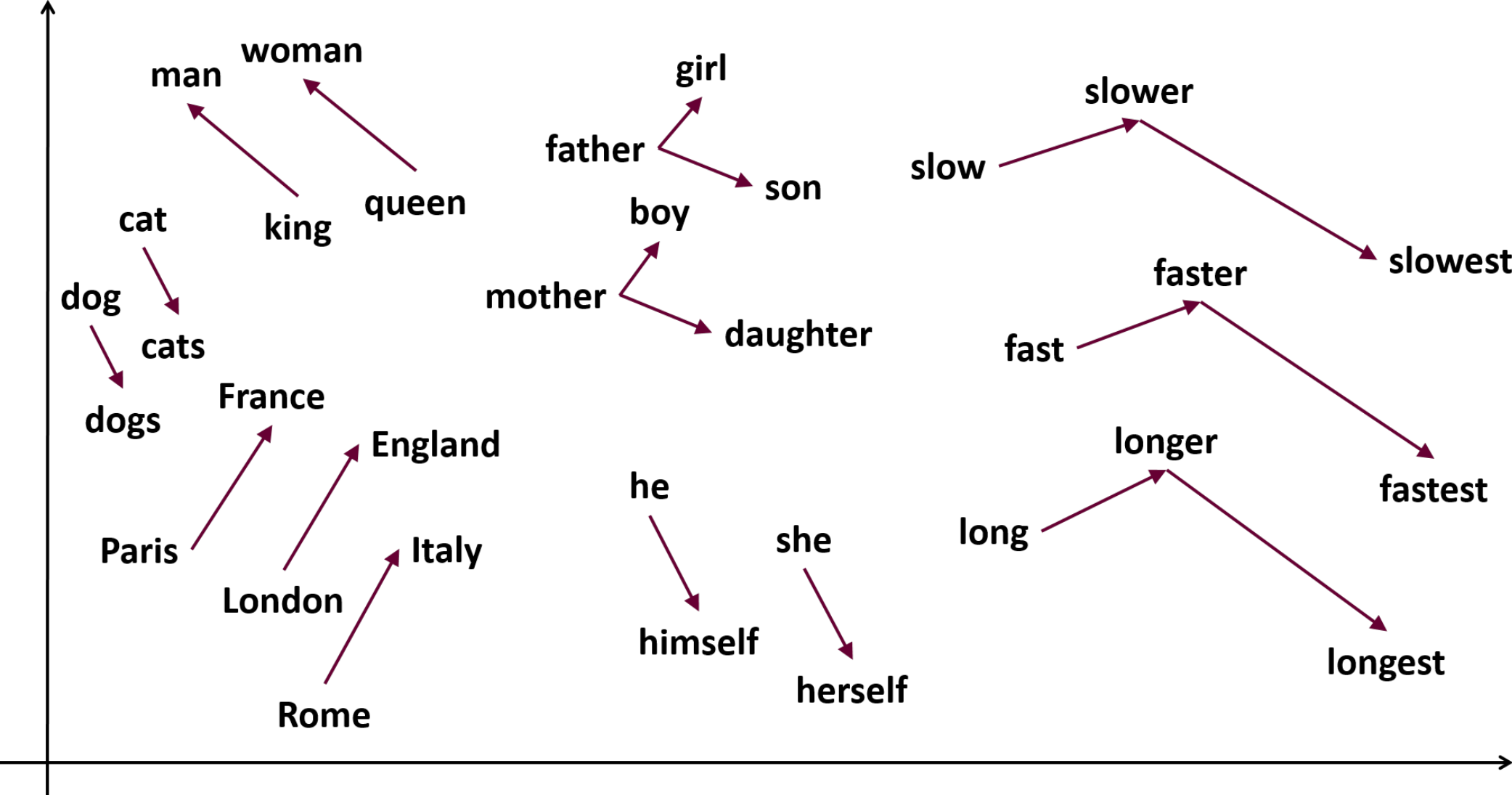
Word2vec



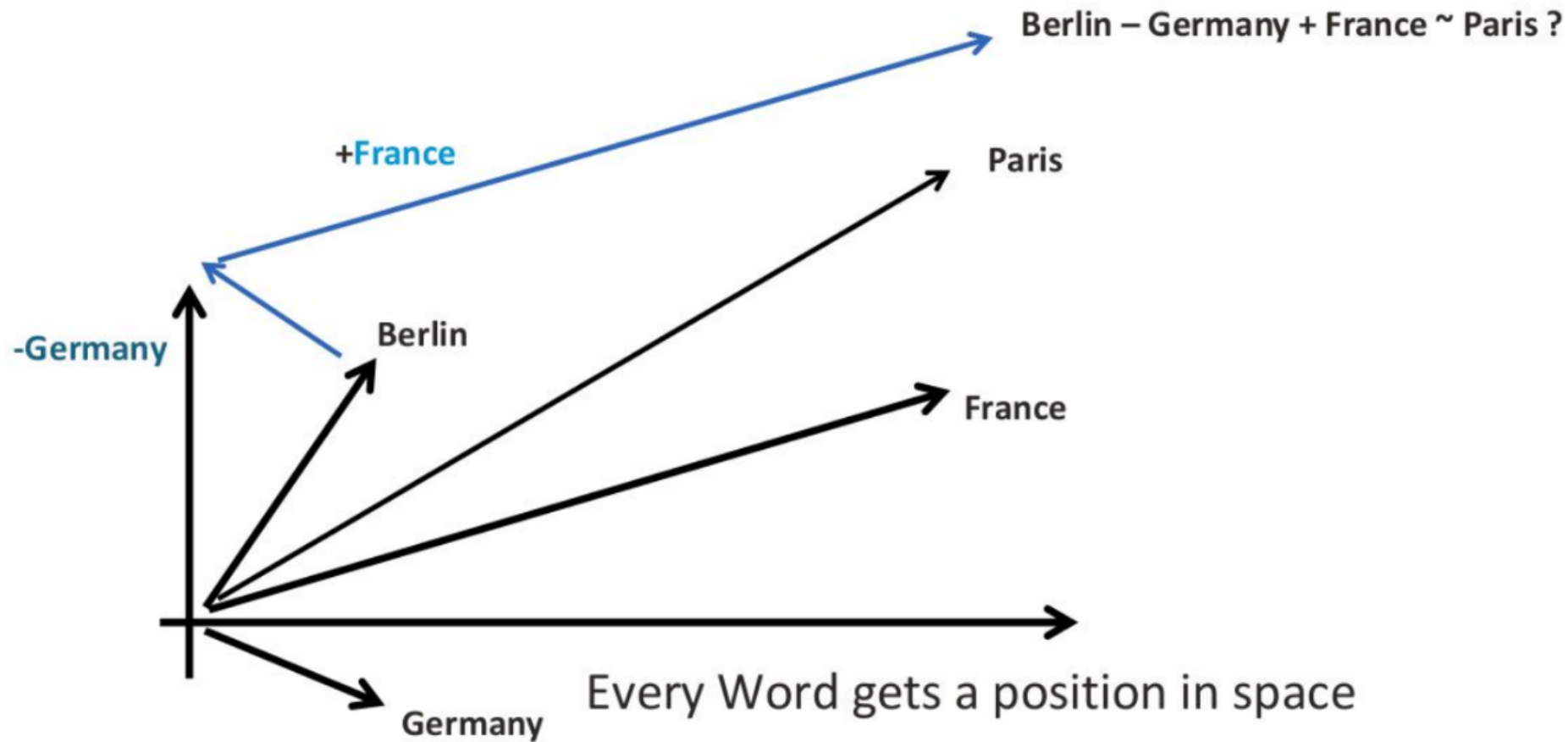
Word2vec



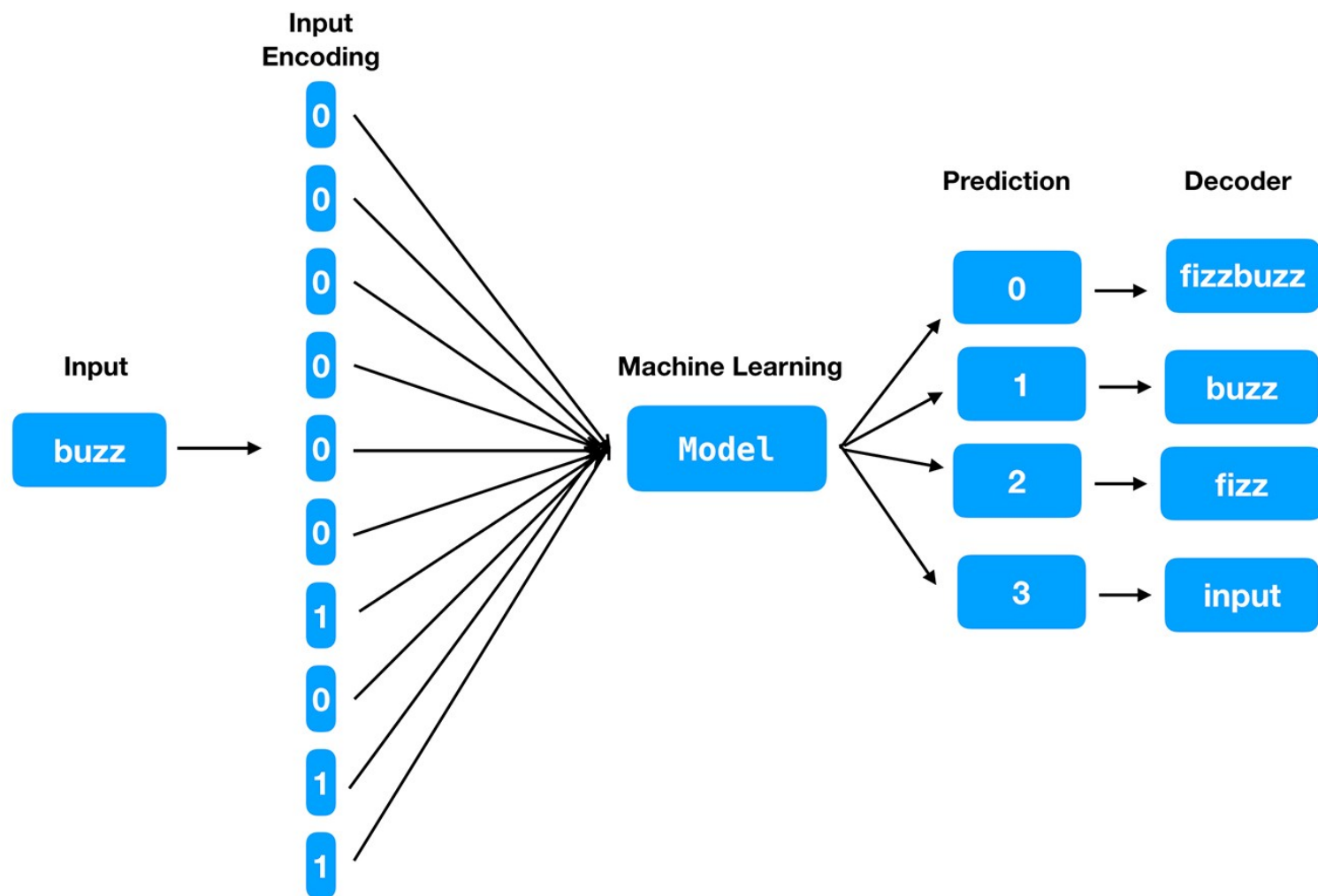
Word2vec



Word2vec



Char2vec

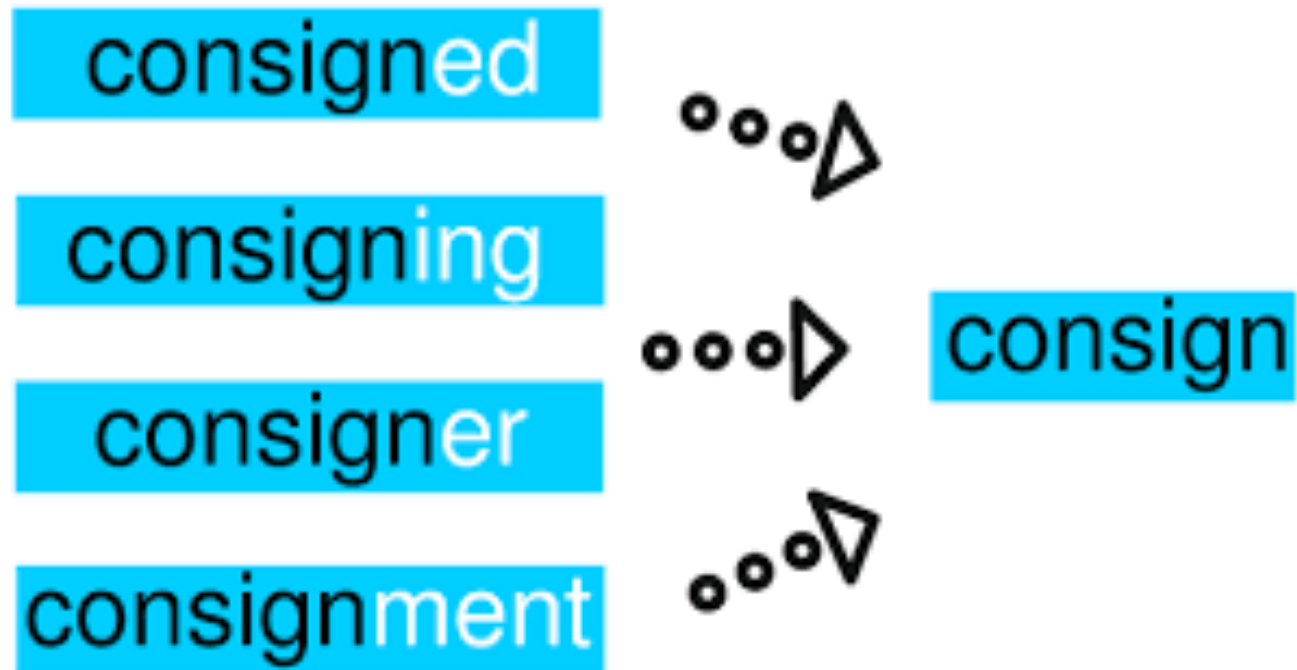


Слов много, куда девать?

Стемминг



Стемминг



Стемминг



Лемматизация

обезьяны ⇒ обезьяна
искал ⇒ искать
любезных ⇒ любезный

Лемматизация

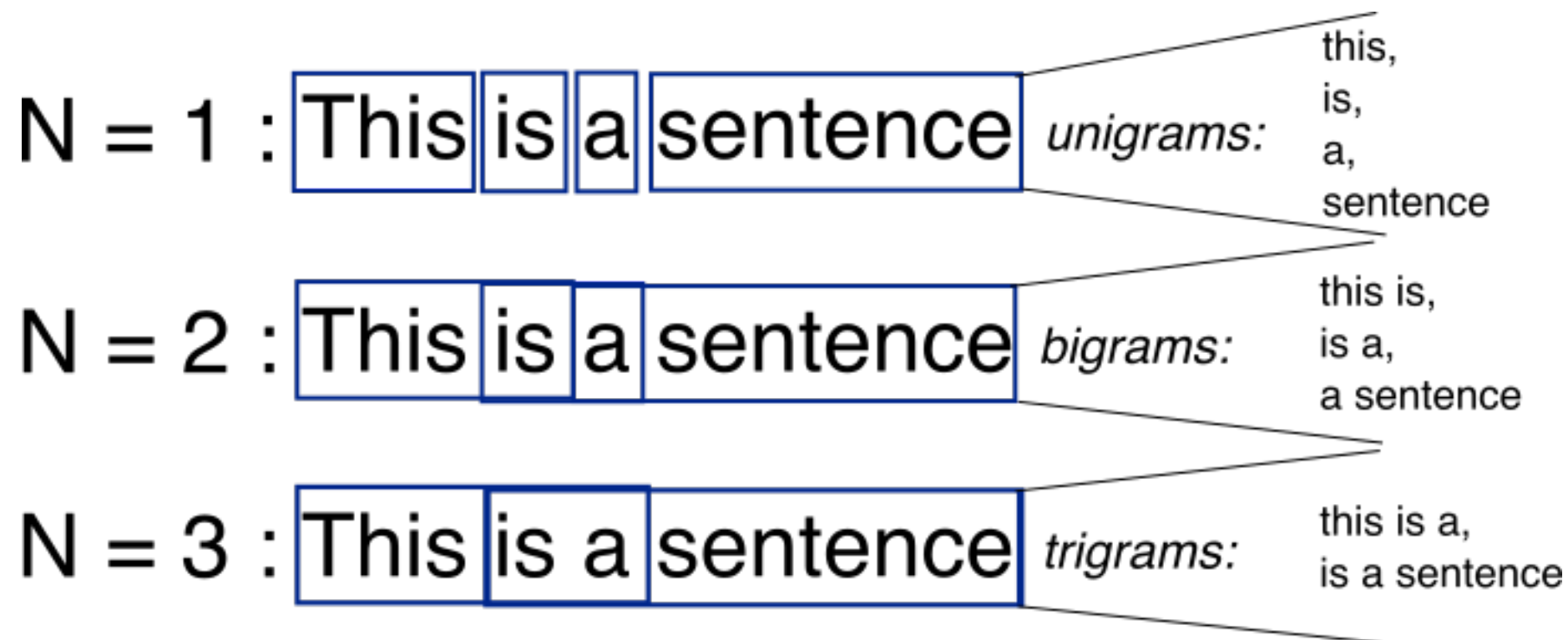
Пушистая	кошка	очень	сладко	спит	на	мягком	кожаном	диване	,	поймав	и	съев	маленькую	противную	мышку	с	длинным	хвостиком
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
пушистый	кошка	очень	сладко сладкий	спать	на	мягкий	кожаный	диван	,	поймав	и	съев	маленький	противный	мышка	с	длинный	хвостик

СТОП-слова

- Союзы
- Предлоги
- Частицы
- Высокочастотные слова
- И т.д. и т.п.

**А как можно определить
взаимодействие слов?**

N-gram



ВОПРОСЫ