

# **Векторное представление СЛОВ**

**Word embeddings. Introduction**

# Текст и МО

- Алгоритмы МО работают с числовыми данными (признаками)
- А мы работаем с текстами и словами
- Как превратить текст или слова в цифры?

P.S. до того, как превращать текст в чиселки, нам нужно сделать **предобработку**: нормализовать, токенизировать, лемматизировать, решить момент с пунктуацией



# Мешок слов

aka bag-of-words, векторное представление текста

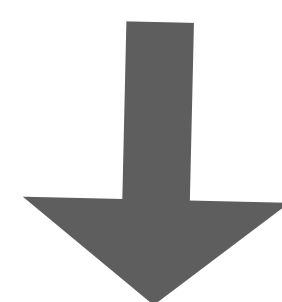
- У нас есть несколько документов (напр., письма со спамом/не спамом)
- делаем множество из слов во всех документах
- считаем и заносим в вектор число, сколько раз каждое слово встретилось в каждом документе.

Document	the	cat	sat	in	hat	with
<i>the cat sat</i>	1	1	1	0	0	0
<i>the cat sat in the hat</i>	2	1	1	1	1	0
<i>the cat with the hat</i>	2	1	0	0	1	1

# Мешок слов

aka bag-of-words

Document	the	cat	sat	in	hat	with
<i>the cat sat</i>	1	1	1	0	0	0
<i>the cat sat in the hat</i>	2	1	1	1	1	0
<i>the cat with the hat</i>	2	1	0	0	1	1



*the cat sat*: [1, 1, 1, 0, 0, 0]

*the cat sat in the hat*: [2, 1, 1, 1, 1, 0]

*the cat with the hat*: [2, 1, 0, 0, 1, 1]

# Bag of Words Example

## Document 1

The quick brown  
fox jumped over  
the lazy dog's  
back.

## Document 2

Now is the time  
for all good men  
to come to the  
aid of their party.

Term	Document 1	Document 2
aid	0	1
all	0	1
back	1	0
brown	1	0
come	0	1
dog	1	0
fox	1	0
good	0	1
jump	1	0
lazy	1	0
men	0	1
now	0	1
over	1	0
party	0	1
quick	1	0
their	0	1
time	0	1

## Stopword List

for
is
of
the
to

# Мешок слов

## aka bag-of-words

- Упрощенное векторное представление текста
- У нас есть несколько документов (напр., письма со спамом/не спамом), считаем и заносим в вектор число, сколько раз каждое слово встретилось в каждом документе.
- Можно считать по N-граммам
- Плюсы? Минусы?
- \*В каких задачах можно использовать?

# Задание

Посчитайте мешок слов для этих 3 текстов

good movie
not a good movie
did not like



# Задание

Посчитайте мешок слов для этих 3 текстов

good movie		good	movie	not	a	did	like
not a good movie							
did not like							

# Задание

- сначала делаем множество слов
- считаем кол-во каждого слова в каждом документе

good movie		<b>good</b>	<b>movie</b>	<b>not</b>	<b>a</b>	<b>did</b>	<b>like</b>
not a good movie	→	1	1	0	0	0	0
did not like		1	1	1	1	0	0
		0	0	1	0	1	1

# Tf-idf

## Мешок слов на стероидах

TF



Frequency of a word  
within the document

IDF



Frequency of a word  
across the documents

# Tf-idf

Мешок слов на стероидах

$$w_{x,y} = \text{tf}_{x,y} \times \log \left( \frac{N}{\text{df}_x} \right)$$

**TF-IDF**

Term  $x$  within document  $y$

$\text{tf}_{x,y}$  = frequency of  $x$  in  $y$

$\text{df}_x$  = number of documents containing  $x$

$N$  = total number of documents

# Берем «документы»

- $d_1$ : *"The sky is blue."*
- $d_2$ : *"The sun is bright today."*
- $d_3$ : *"The sun in the sky is bright."*
- $d_4$ : *"We can see the shining sun, the bright sun."*

# Выкидываем стоп-слова

- $d_1$ : "sky blue"
- $d_2$ : "sun bright today"
- $d_3$ : "sun sky bright"
- $d_4$ : "can see shining sun bright sun"

# Рассчитываем частоту слов (term frequency)

$$f_{t,d}$$

	blue	bright	can	see	shining	sky	sun	today
1	1	0	0	0	0	1	0	0
2	0	1	0	0	0	0	1	1
3	0	1	0	0	0	1	1	0
4	0	1	1	1	1	0	2	0



$$\text{tf}(t, d) = \frac{f_{t,d}}{\sum_{t'} f_{t',d}}$$

	blue	bright	can	see	shining	sky	sun	today
1	1/2	0	0	0	0	1/2	0	0
2	0	1/3	0	0	0	0	1/3	1/3
3	0	1/3	0	0	0	1/3	1/3	0
4	0	1/6	1/6	1/6	1/6	0	1/3	0



# Рассчитываем обратную частоту слов в документе (idf)

$f_{t,d}$

	blue	bright	can	see	shining	sky	sun	today
1	1	0	0	0	0	1	0	0
2	0	1	0	0	0	0	1	1
3	0	1	0	0	0	1	1	0
4	0	1	1	1	1	0	2	0
n_t	1	3	1	1	1	2	3	1



$N = 4$

$$\text{idf}(t, D) = \log_{10} \frac{N}{n_t}$$

	blue	bright	can	see	shining	sky	sun	today
	0.602	0.125	0.602	0.602	0.602	0.301	0.125	0.602

$\log_{10} \frac{4}{1} = 0.602$        $\log_{10} \frac{4}{3} = 0.125$



# Получаем матрицу Tf-idf

$tf(t, d)$

	blue	bright	can	see	shining	sky	sun	today
1	1/2	0	0	0	0	1/2	0	0
2	0	1/3	0	0	0	0	1/3	1/3
3	0	1/3	0	0	0	1/3	1/3	0
4	0	1/6	1/6	1/6	1/6	0	1/3	0

**x**

$idf(t, D)$

	blue	bright	can	see	shining	sky	sun	today
	0.602	0.125	0.602	0.602	0.602	0.301	0.125	0.602

$$tfidf(t, d, D) = tf(t, d) \cdot idf(t, D)$$

- TF-IDF: Multiply TF and IDF scores, use to rank importance of words within documents
- Most important word for each document is highlighted



	blue	bright	can	see	shining	sky	sun	today
1	<b>0.301</b>	0	0	0	0	0.151	0	0
2	0	0.0417	0	0	0	0	0.0417	<b>0.201</b>
3	0	0.0417	0	0	0	<b>0.100</b>	0.0417	0
4	0	0.0209	<b>0.100</b>	<b>0.100</b>	<b>0.100</b>	0	0.0417	0

# Tf-idf

## Мешок слов на стероидах

- показывает важность (вес) слова в документе
- занижает вес слишком частотных слов

# Задание

Посчитайте теперь TF-IDF для этих текстов

good movie
not a good movie
did not like

# Задание

Посчитайте теперь TF-IDF для этих текстов

	<b>good</b>	<b>movie</b>	<b>not</b>	<b>a</b>	<b>did</b>	<b>like</b>
good movie						
not a good movie						
did not like						

# Слова, а не тексты

- ONE слишком плох
- Мешок слов и TF-IDF являются представлениями текстов/ документов, но не конкретных слов. Т.е. работают с классификациями текстов.
- А что если для задачи нужно переводить в векторы именно слова?

# Word2vec

## Words to vectors

- Вектор-представление для каждого слова
- Что мы хотим от такого вектора?

# Word2vec

## Words to vectors

- Вектор-представление для каждого слова
- Что мы хотим от такого вектора?
- Отражал смысл слова (похожие слова имели бы похожие векторы)

# Лирическое отступление

Психологический тест «Большая пятерка»

Openness to experience	79 out of 100
Agreeableness	75 out of 100
Conscientiousness	42 out of 100
Negative emotionality	50 out of 100
Extraversion	58 out of 100



# Большая пятерка (первая черта)

Масштабируем по шкале от 1 до -1

Extraversion

100

0

Introversion

Jay

Extraversion

38

Extraversion

1

-1

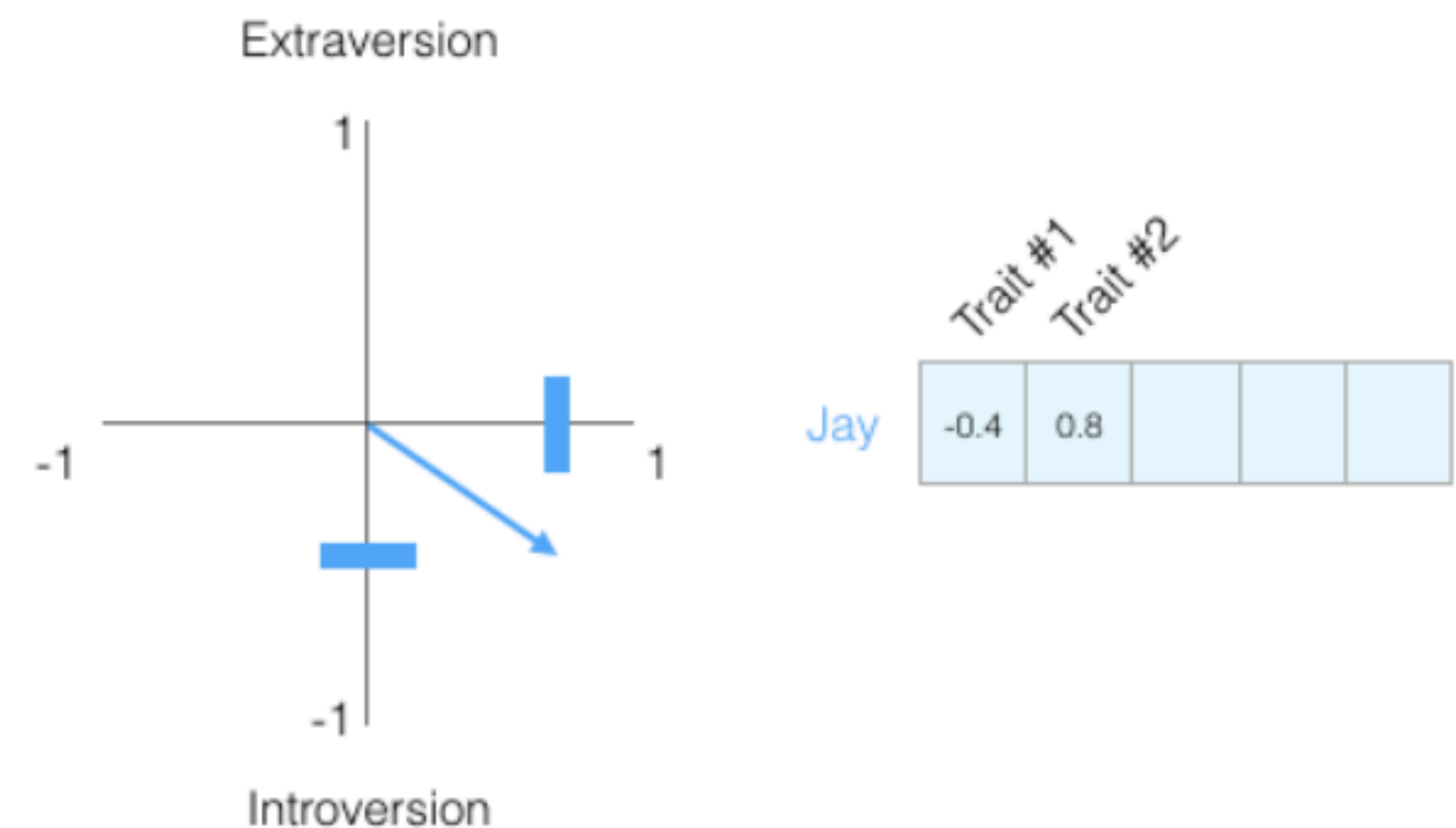
Introversion

Jay

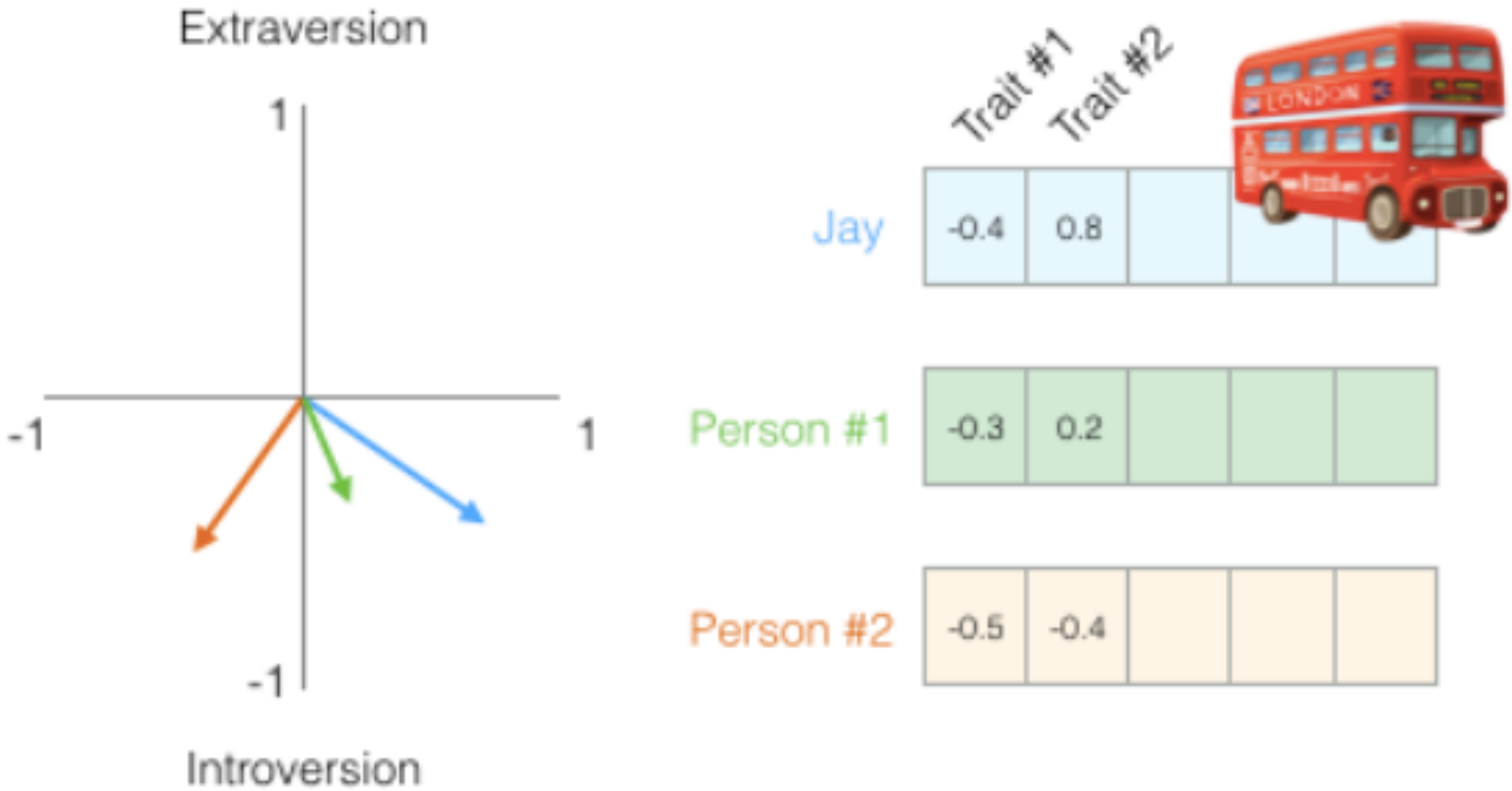
Extraversion

-0.4

# Большая пятерка (добавим еще черту)

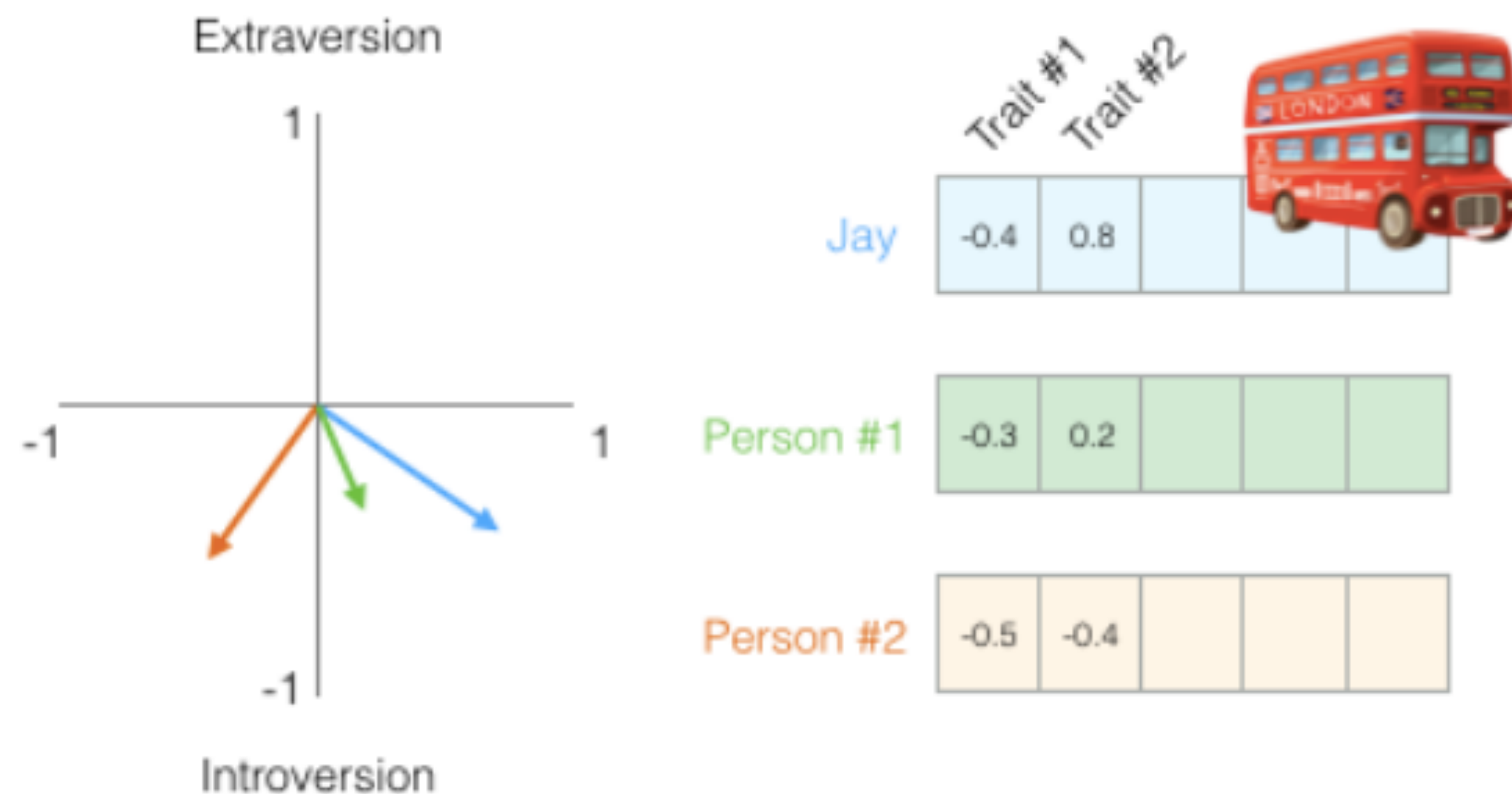


# Большая пятерка



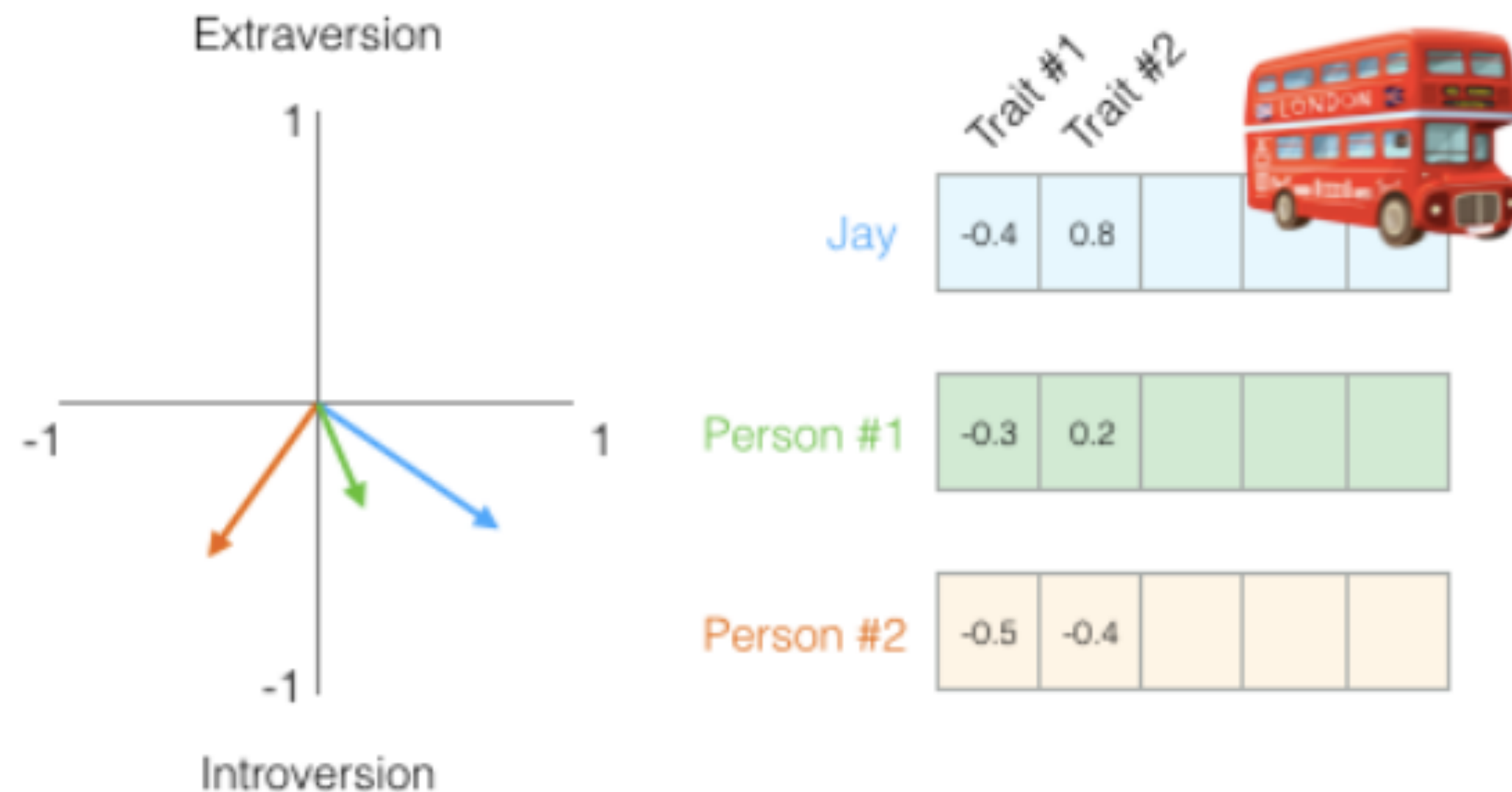
Еще 2 человека ответили на 2 вопроса теста  
Кто больше похож на синего?

# Большая пятерка



- В 3-мерном пространстве можно определить на глаз
- Но лучше не полагаться на это и пользоваться субъективным инструментом
- Есть идеи?

# Большая пятерка



- В 3-мерном пространстве можно определить на глаз
- Но лучше не полагаться на это и пользоваться субъективным инструментом
- Косинусная близость (угол между векторами)

$$\text{cosine\_similarity}(\text{Jay}, \text{Person \#1}) = 0.87 \quad \checkmark$$

$$\text{cosine\_similarity}(\text{Jay}, \text{Person \#2}) = -0.20$$

# Многомерное пространство

	Trait #1	Trait #2	Trait #3	Trait #4	Trait #5
Jay	-0.4	0.8	0.5	-0.2	0.3
Person #1	-0.3	0.2	0.3	-0.4	0.9
Person #2	-0.5	-0.4	-0.2	0.7	-0.1

- Больше 3-мерного пространства не можем визуализировать
- Теперь точно сравниваем только с помощью косинусной близости

$$\text{cosine\_similarity}(\text{Jay}, \text{Person \#1}) = 0.66 \quad \checkmark$$

$$\text{cosine\_similarity}(\text{Jay}, \text{Person \#2}) = -0.37$$

# Переходим к словам

Почти word2vec

**Дистрибутивная гипотеза:** смысл слова - в его контексте, то есть в словах, с которым оно чаще всего встречается.

Откуда взять цифры для векторов всех слов?

# Переходим к словам

Почти word2vec

**Дистрибутивная гипотеза:** смысл слова - в его контексте, то есть в словах, с которым оно чаще всего встречается.

Откуда взять цифры для векторов всех слов?

	редис	картошка	кот	...	собака
редис	-	5	1		0
картошка	5	-	0		1
кот	1	0	-		6
...					
собака	0	1	6		-



# Эмбе́ддинг

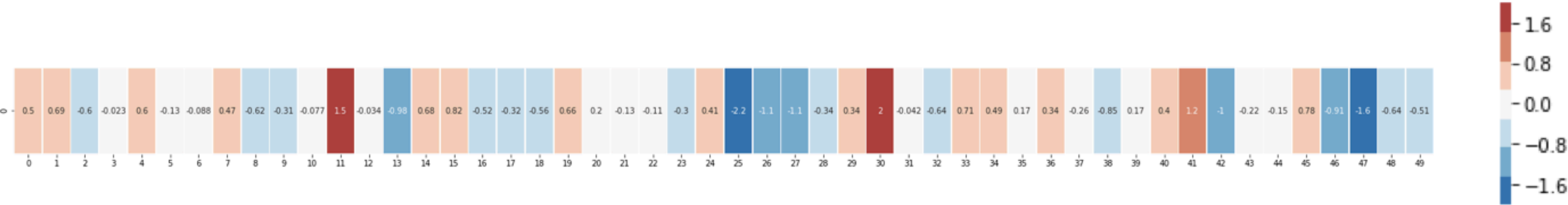
## Векторное представление слова

Вектор слова «король»

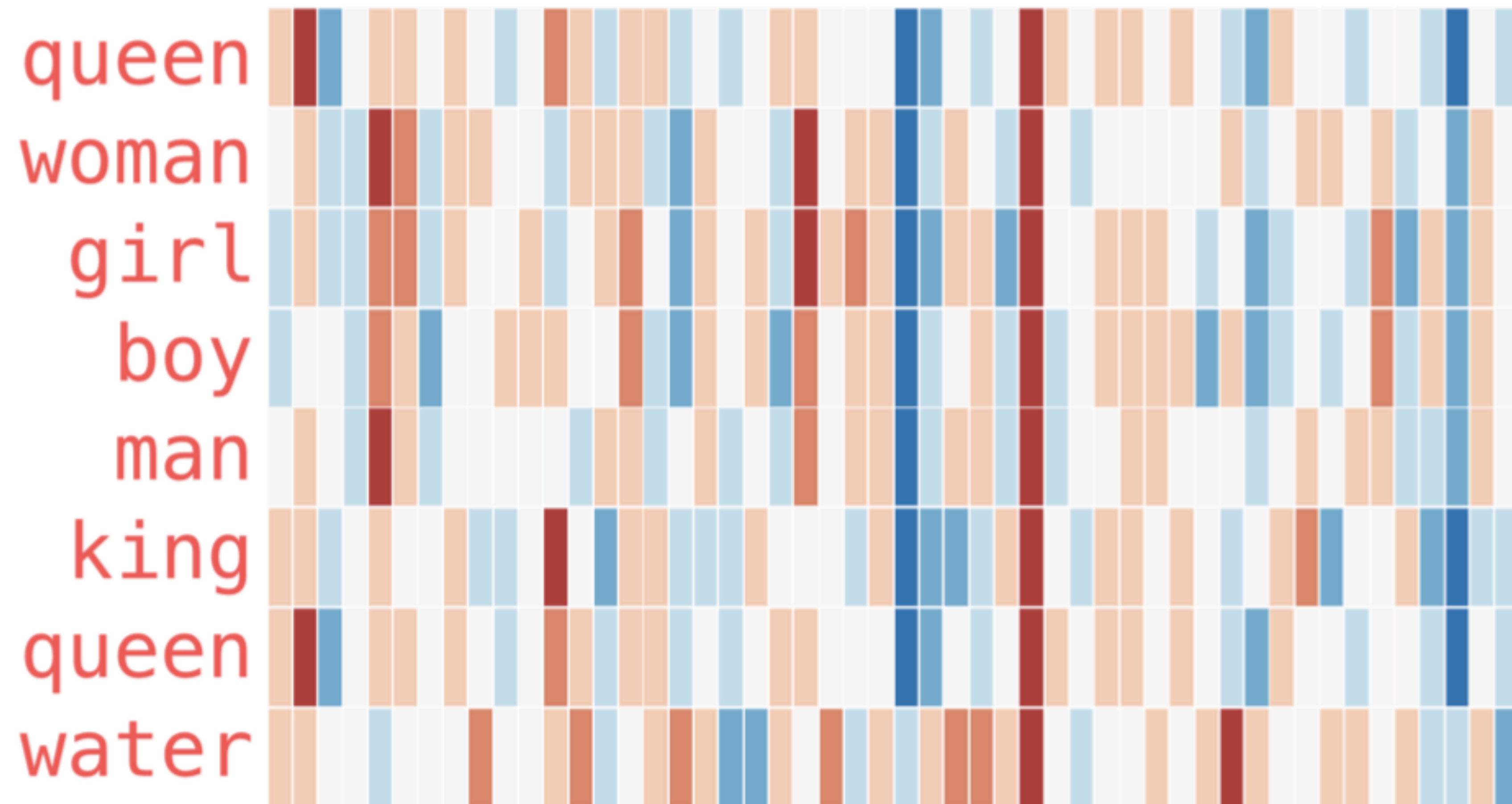
```
[ 0.50451 , 0.68607 , -0.59517 , -0.022801, 0.60046 , -0.13498 , -0.08813 ,  
0.47377 , -0.61798 , -0.31012 , -0.076666, 1.493 , -0.034189, -0.98173 ,  
0.68229 , 0.81722 , -0.51874 , -0.31503 , -0.55809 , 0.66421 , 0.1961 ,  
-0.13495 , -0.11476 , -0.30344 , 0.41177 , -2.223 , -1.0756 , -1.0783 ,  
-0.34354 , 0.33505 , 1.9927 , -0.04234 , -0.64319 , 0.71125 , 0.49159 ,  
0.16754 , 0.34344 , -0.25663 , -0.8523 , 0.1661 , 0.40102 , 1.1685 , -1.0137  
, -0.21585 , -0.15155 , 0.78321 , -0.91241 , -1.6106 , -0.64426 , -0.51042 ]
```

# Эмбеддинг «Король»

```
[ 0.50451 , 0.68607 , -0.59517 , -0.022801, 0.60046 , -0.13498 , -0.08813 ,  
0.47377 , -0.61798 , -0.31012 , -0.076666, 1.493 , -0.034189, -0.98173 ,  
0.68229 , 0.81722 , -0.51874 , -0.31503 , -0.55809 , 0.66421 , 0.1961 ,  
-0.13495 , -0.11476 , -0.30344 , 0.41177 , -2.223 , -1.0756 , -1.0783 ,  
-0.34354 , 0.33505 , 1.9927 , -0.04234 , -0.64319 , 0.71125 , 0.49159 ,  
0.16754 , 0.34344 , -0.25663 , -0.8523 , 0.1661 , 0.40102 , 1.1685 , -1.0137  
, -0.21585 , -0.15155 , 0.78321 , -0.91241 , -1.6106 , -0.64426 , -0.51042 ]
```

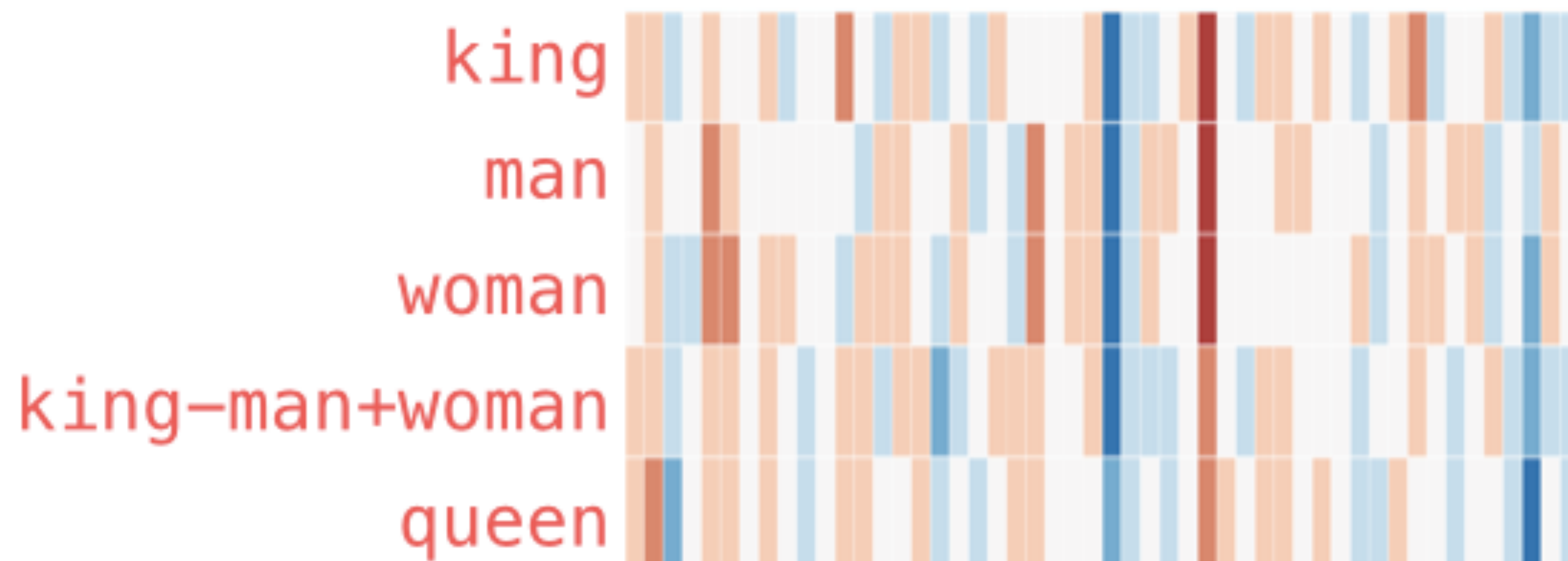


# Сравните эмбе́ддинги



# Эмбеддинги

king – man + woman ≈ queen



Полученный вектор от вычисления «король–мужчина+женщина» не совсем равен «королеве», но это наиболее близкий результат из 400 000 вложений слов в наборе данных

# Полезные материалы

1. <https://habr.com/ru/post/446530/>
2. <https://sysblok.ru/knowhow/word2vec-pokazhi-mne-svoj-kontekst-i-ja-skazhu-kto-ty/>
3. <https://rusvectors.org/ru/>

Блокнот с word2vec и теорией:

- [https://github.com/hse-ds/iad-deep-learning/blob/master/2021/seminars/sem07/sem07\\_solution.ipynb](https://github.com/hse-ds/iad-deep-learning/blob/master/2021/seminars/sem07/sem07_solution.ipynb)