

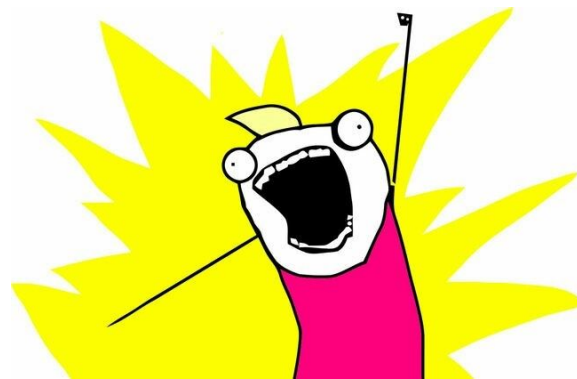
# Тематическое моделирование

... и как сделать свой Brand Analytics

# Задача

- Выделить в коллекции документов темы
- Мы заранее не знаем, сколько этих тем и какие они
- Приятный бонус - векторная модель (об этом попозже)

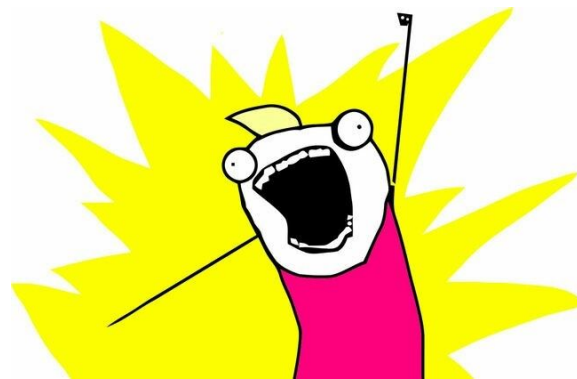
Что мы хотим?

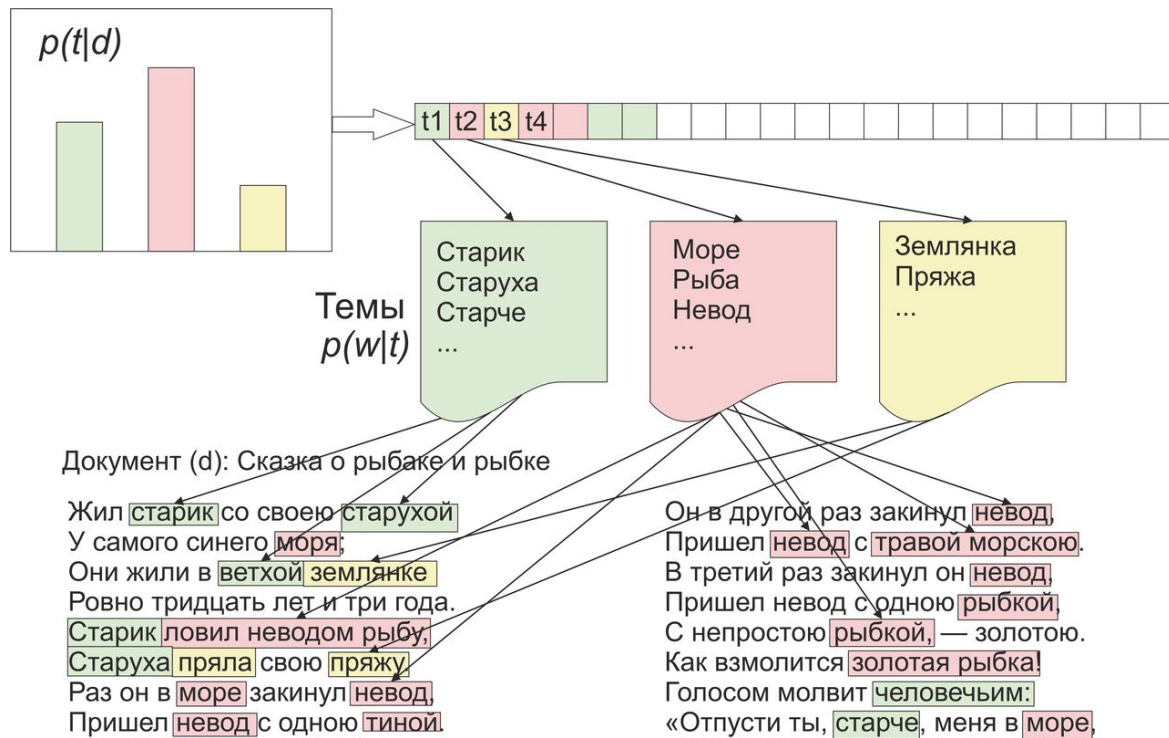


# Цель

- Визуализация
- Эксплоративный анализ
- Кластеризация
- И снова векторная модель

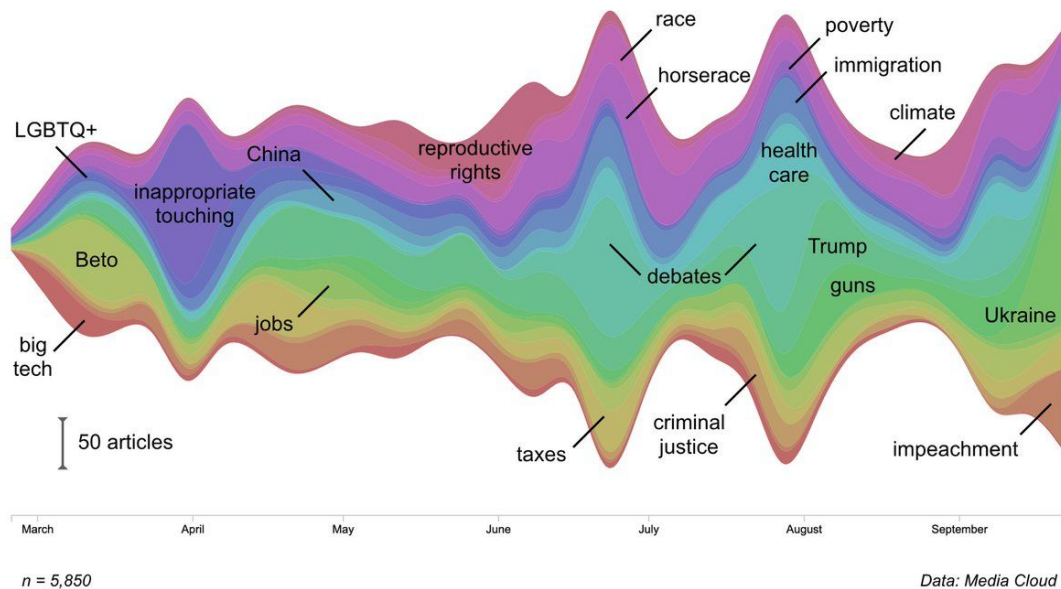
Зачем мы это хотим?





## How media are setting the 2020 agenda

A topic analysis of news articles published by 28 outlets since March 2019 mentioning Joe Biden, Bernie Sanders, Elizabeth Warren, Kamala Harris, Pete Buttigieg, Beto O'Rourke, Cory Booker, Kirsten Gillibrand, Amy Klobuchar, or Tulsi Gabbard



<https://www.trendsmap.com/twitter/tweet/1194233155642941440>



# Тематические модели: много трехбуквенных аббревиатур

- **LSA** (Latent semantic analysis, латентно-семантический анализ)
- **pLSA** (probabilistic latent semantic analysis, вероятностный латентно-семантический анализ)
- **LDA** (Latent Dirichlet Allocation, латентное размещение Дирихле)
- **BigARTM** (ARTM — Additive Regularization for Topic Modeling)

# LSA

- Самая ранняя модель
- Предполагает наличие скрытого (латентного) параметра - темы
- Наиболее популярная реализация - сингулярное разложение (SVD) матрицы термины-документы



# Матрица терми-документы

1. Медведев прибыл с официальным визитом во Вьетнам.
2. В Крыму объявили экстренное предупреждение из-за шторма.
3. В Пушкинской галерее пройдет выставка Рафаэля.
4. Премьер-министр России встретился с коллегами на саммите во Вьетнаме.
5. На презентации Apple продемонстрировали новый iPad Pro.
6. Работы Рафаэля впервые привезут в Москву.
7. В понедельник в Крыму ожидаются сильные дожди и штормовой ветер.
8. В Калифорнии стартовала ежегодная презентация компании Apple.

	D1	D2	D3	D4	D5	D6	D7	D8
<b>apple</b>	0	0	0	0	1	0	0	1
<b>ipad</b>	0	0	0	0	1	0	0	1
<b>pro</b>	0	0	0	0	1	0	0	1
<b>визит</b>	1	0	0	0	0	0	0	0
<b>вьетнам</b>	1	0	0	1	0	0	0	0
<b>галерея</b>	0	0	1	0	0	0	0	0
<b>ежегодный</b>	0	0	0	1	0	0	0	1
<b>крым</b>	0	1	0	0	0	0	1	0
...	...	...	...	...	...	...	...	...

Матрицу можно заполнить tf-idf весами!

# SVD

The diagram illustrates the Singular Value Decomposition (SVD) of a matrix  $A$ . Matrix  $A$  is shown as a pink rectangle with dimensions  $n \times d$ . It is equal to the product of three matrices:  $\hat{U}$  (pink rectangle,  $n \times r$ ),  $\Sigma$  (blue rectangle,  $n \times d$ ), and  $\hat{V}^T$  (pink rectangle,  $r \times d$ ). The matrix  $\Sigma$  is depicted as a blue rectangle with a pink top-left corner of size  $r \times r$  containing the symbol  $\hat{\Sigma}$ . Below the matrices, their full dimensions are listed:  $U$  is  $n \times n$ ,  $\Sigma$  is  $n \times d$ , and  $V^T$  is  $d \times d$ . The text 'Исходная матрица терм-документы' is positioned below matrix  $A$ .

Исходная матрица терм-документы

$U$   
 $n \times n$

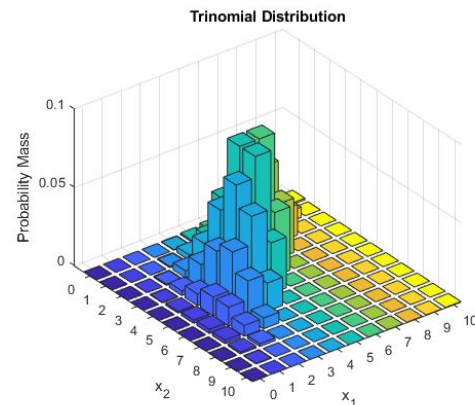
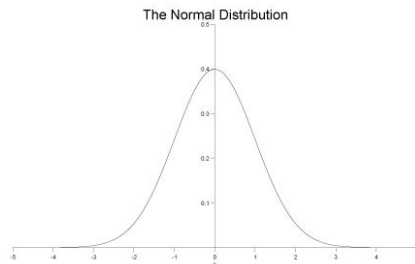
$\Sigma$   
 $n \times d$

$V^T$   
 $d \times d$

[Наглядная gifка](#)

# pLSA

- LSA предполагает, что распределение слов и документов нормальное, pLSA - мультиномиальное
- Совместная встречаемость термина и документа ( $w, d$ ) моделируется как сочетание независимых мультиномиальных распределений, где  $s$  — тема. Количество тем — это гиперпараметр, который выбирается до начала анализа



# LDA

- генеративная вероятностная модель
- документ - это набор случайных скрытых тем, где каждая тема определяется распределением слов, при этом каждое слово в конкретном документе можно отнести к одной из его тем
- в качестве априорного распределения для тем используется распределение Дирихле

# А что там с векторной моделью?

- Тематическое моделирование - еще один способ (помимо подходов дистрибутивной семантики) получить векторное представление слов и документов, при этом полученные вектора тем будут **интерпретируемыми**

Topic 1		Topic 2		Topic 3	
term	weight	term	weight	term	weight
game	0.014	space	0.021	drive	0.021
team	0.011	nasa	0.006	card	0.015
hockey	0.009	earth	0.006	system	0.013
play	0.008	henry	0.005	scsi	0.012
games	0.007	launch	0.004	hard	0.011

# Что на практике

1. Выбираем количество тем  $N$
2. Строим тематическую модель с  $N$  темами
3. Оцениваем качество (в этом, в том числе, помогает визуализация)
4. Пробуем другое число тем  $N$

Перерыв!

# Brand Analytics: от названия бренда до аналитических отчетов





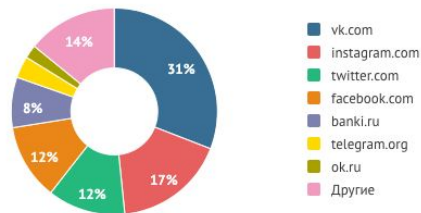
# Что умеет сервис?



# Что это значит на самом деле:

1. Используем открытые API соцсетей (твиттер, VK и т.д.) и медиа-порталов, чтобы вытаскивать все тексты с упоминанием бренда

scraping      regular  
expressions



# Что это значит на самом деле:

## 2. Предобрабатываем текст

scraping

regular  
expressions

text normalization

# Что это значит на самом деле:

3. Частотный анализ: количество упоминаний, ключевые слова и словосочетания, облака слов и т.д.

scraping      regular expressions      text normalization

tf-idf      keywords extraction

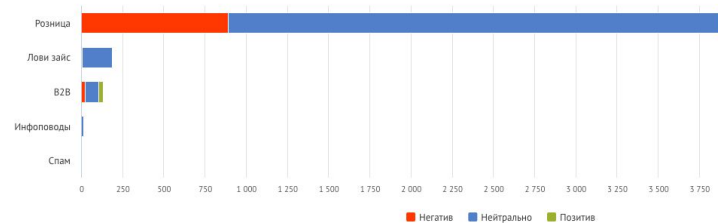
8 584 ↑ 13.7%  
сообщений

6 247 ↑ 21.0%  
авторов



# Что это значит на самом деле:

## 4. Строим тематические модели



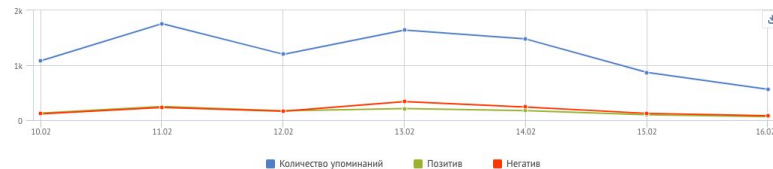
scraping      regular expressions      text normalization

tf-idf      keywords extraction

topic modeling

# Что это значит на самом деле:

## 5. Делаем анализ тональности



scraping      regular expressions      text normalization

tf-idf      keywords extraction

topic modeling      word embeddings

text classification