

# ОБУЧЕНИЕ МОДЕЛЕЙ НА ТЕКСТАХ

---

# ПОДГОТОВКА ВЫБОРКИ

---



- Удаление редких и популярных слов

# ПОДГОТОВКА ВЫБОРКИ

---

- › Удаление редких и популярных слов
- › Признаки:  $n$ -граммы + счётчики/TF-IDF

# ПОДГОТОВКА ВЫБОРКИ

---



➤ Число признаков —  $10^3 - 10^4$  и больше

# ПОДГОТОВКА ВЫБОРКИ

---

- › Число признаков —  $10^3 - 10^4$  и больше
- › Можно пробовать отбор признаков и понижение размерности

- Случайный лес — низкая скорость обучения

- › Случайный лес — низкая скорость обучения
- › Градиентный бустинг — проблемы из-за маленькой глубины деревьев

- › Случайный лес — низкая скорость обучения
- › Градиентный бустинг — проблемы из-за маленькой глубины деревьев
- › Наивный байесовский классификатор



- › Случайный лес — низкая скорость обучения
- › Градиентный бустинг — проблемы из-за маленькой глубины деревьев
- › Наивный байесовский классификатор
- › Линейные модели используются чаще всего

- › Стохастический градиентный спуск позволяет читать с диска по одному объекту

› Стохастический градиентный спуск позволяет читать с диска по одному объекту

› пока не выполнен критерий останова:

$t$  = следующий текст

для всех слов  $x$  в  $t$  :

$$w_{h(x)} = w_{h(x)} - \alpha \nabla_{w_{h(x)}} Q(w)$$

➤  $n$ -граммы и мешок слов

- ›  $n$ -граммы и мешок слов
- › Линейные модели и хэширование