

# ПРОКЛЯТИЕ РАЗМЕРНОСТИ. РЕГУЛЯРИЗАЦИЯ

## КВАНТИТАТИВНЫЙ АНАЛИЗ ТЕКСТА

---

Кирилл Александрович Маслинский

04.04.2022 / 06

НИУ ВШЭ Санкт-Петербург

# ПРОКЛЯТИЕ РАЗМЕРНОСТИ

---

## SPARSE DATA PROBLEM

---

Terms					
Docs	выгребать	выгребной	выгружать	выгрузка	выгрыза
1	0	0	0	0	
2	0	0	0	0	
3	0	0	0	0	
4	0	0	0	0	
5	0	0	0	0	

A document-term matrix (1530 documents, 13322 terms)

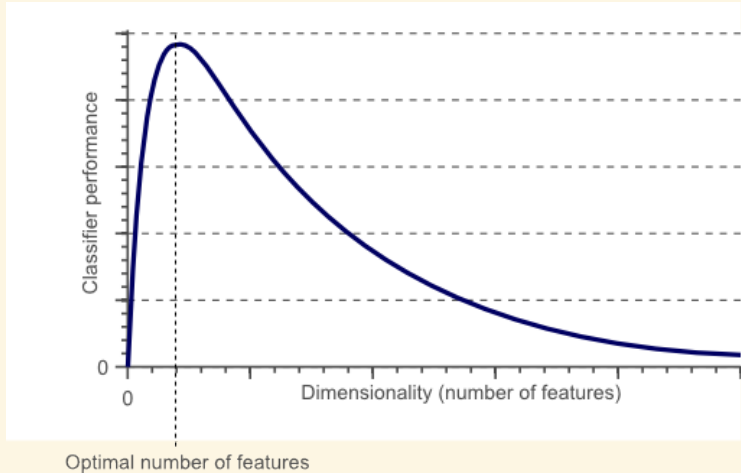
Non-/sparse entries: 68859/20313801

Sparsity : 100%

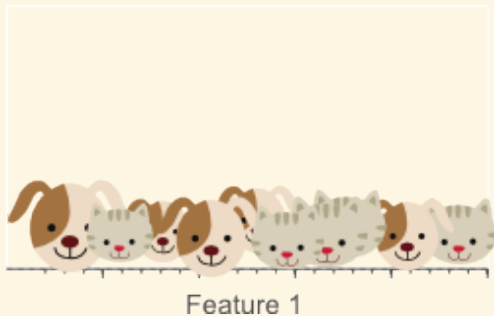
Maximal term length: 66

Weighting : term frequency (tf)

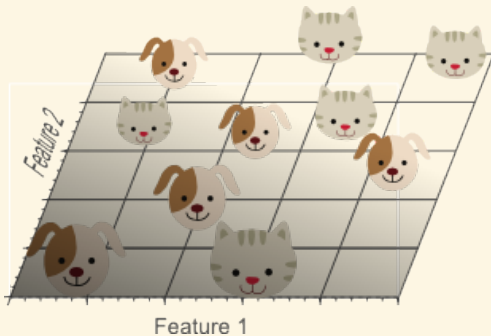
# HUGHES PHENOMENON



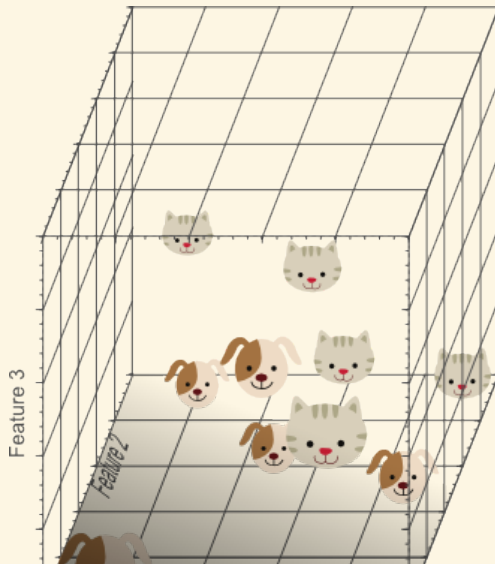
## ПРОКЛЯТИЕ РАЗМЕРНОСТИ НА КОШКАХ



## ПРОКЛЯТИЕ РАЗМЕРНОСТИ НА КОШКАХ

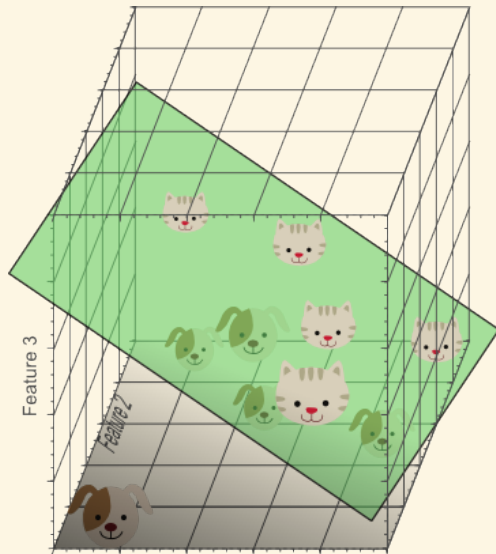


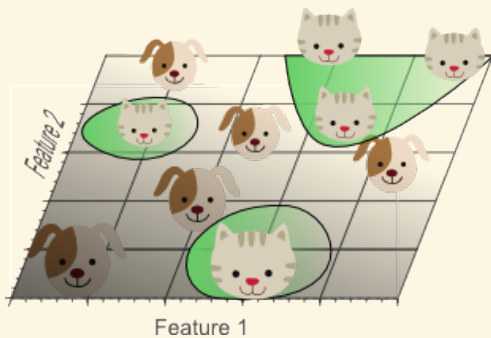
## ПРОКЛЯТИЕ РАЗМЕРНОСТИ НА КОШКАХ

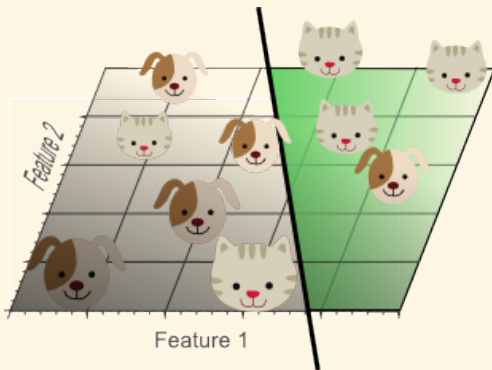




# ПРОКЛЯТИЕ РАЗМЕРНОСТИ НА КОШКАХ







## СНИЖЕНИЕ РАЗМЕРНОСТИ: ПРОСТЫЕ СПОСОБЫ

---

- Матрица терминов-документов очень большая и редкая
- Близкие по смыслу слова не обязательно встречаются в одних и тех же документах:
  - синонимия
  - полисемия
  - шум
- Нужно сократить размерность матрицы (сделать меньше столбцов).

- Матрица терминов-документов очень большая и редкая
- Близкие по смыслу слова не обязательно встречаются в одних и тех же документах:
  - синонимия
  - полисемия
  - шум
- Нужно сократить размерность матрицы (сделать меньше столбцов).

Простейший способ уменьшить число столбцов — просто **удалить лишние слова**:

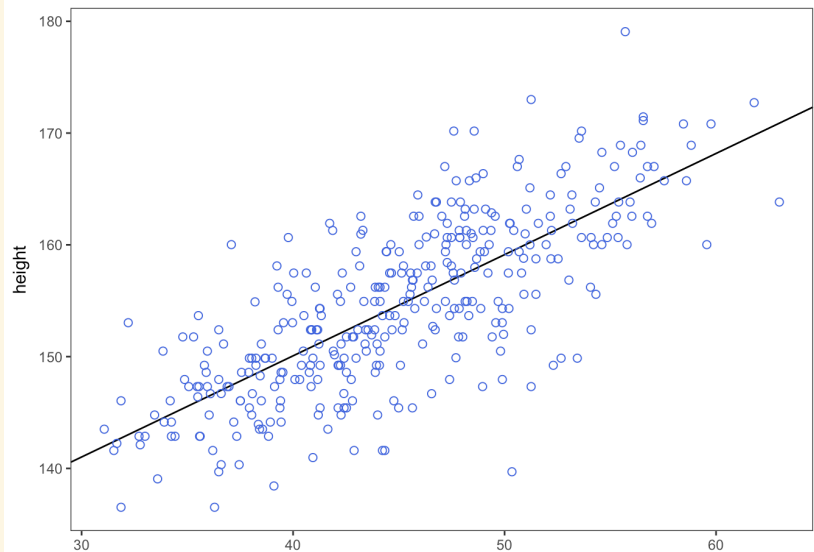
- Статический список:  
без более бы был была были было быть в вам вас весь  
во вот все всего всех вы где да даже для ...
- Динамический список:
  - Слишком частотные ( $N$  самых частотных)
  - Слишком редкие (порог: не менее чем в  $F$  документов)
  - Слишком короткие (меньше  $M$  букв)

# АЛГОРИТМ КЛАССИФИКАЦИИ: ЛОГИСТИЧЕСКАЯ РЕГРЕССИЯ

---



# ПРОСТАЯ ЛИНЕЙНАЯ РЕГРЕССИЯ



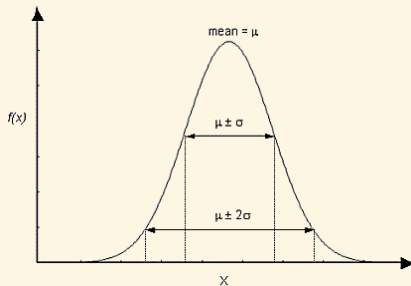
# ЛИНЕЙНАЯ РЕГРЕССИЯ: ГЕНЕРАТИВНАЯ ФОРМУЛИРОВКА

$$\begin{aligned} H_i &\sim \mathcal{N}(\mu_i, \sigma) \\ \mu_i &= \alpha + \beta W_i \end{aligned} \quad (1)$$

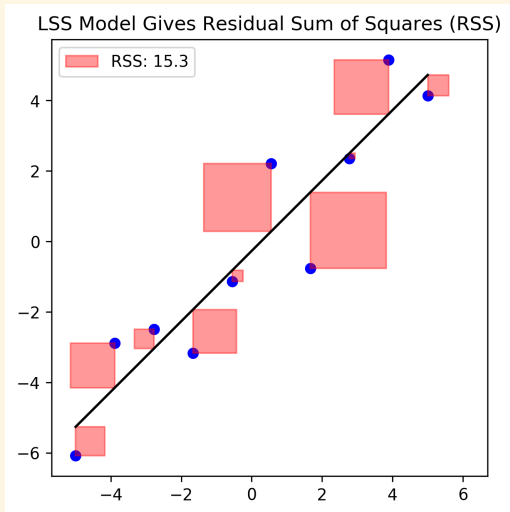
где

$H_i$  рост индивида  $i$

$W_i$  вес индивида  $i$



## RSS — ФУНКЦИЯ ПОТЕРЬ



# ЛОГИСТИЧЕСКАЯ РЕГРЕССИЯ: БИНАРНЫЙ КЛАССИФИКАТОР



$$Y_i \sim \text{Binomial}(1, p_i)$$

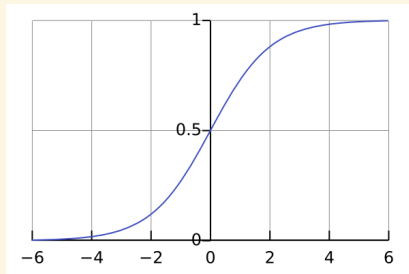
$$\text{logit}(p_i) = \alpha + \beta_1 x_1 + \beta_2 x_2, \quad (2)$$

$$\text{logit}(p_i) = \log \frac{p_i}{1 - p_i}$$

где

$Y_i$  класс индивида  $i \in \{0,1\}$

$p_i$  вероятность позитивного класса (1) для индивида  $i$



# ЛОГИСТИЧЕСКАЯ РЕГРЕССИЯ: ГЕНЕРАТИВНАЯ ФОРМУЛИРОВКА

$$Y_i \sim \text{Binomial}(1, p_i)$$

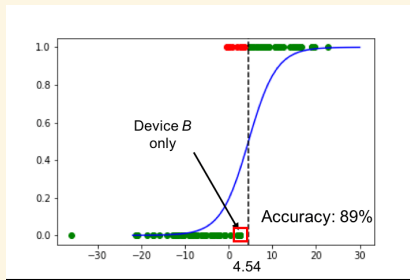
$$\text{logit}(p_i) = \alpha + \beta_1 x_1 + \beta_2 x_2, \quad (2)$$

$$\text{logit}(p_i) = \log \frac{p_i}{1 - p_i}$$

где

$Y_i$  класс индивида  $i \in \{0,1\}$

$p_i$  вероятность позитивного класса (1) для индивида  $i$



$$LS = -(y \log(p) + (1 - y) \log(1 - p)) \quad (3)$$

Примеры:

$$y = 1; p = 0.8$$

$$\begin{aligned} LS &= -(1 * \log(0.8) + (1 - 1) * \log(0.2)) = \\ &= 0.22 \end{aligned} \quad (4)$$

$$y = 0; p = 0.8$$

$$\begin{aligned} LS &= -(0 * \log(0.8) + (1 - 0) * \log(0.2)) = \\ &= 1.6 \end{aligned} \quad (5)$$

$$LS = -(y \log(p) + (1 - y) \log(1 - p)) \quad (3)$$

Примеры:

$$y = 1; p = 0.8$$

$$\begin{aligned} LS &= -(1 * \log(0.8) + (1 - 1) * \log(0.2)) = \\ &= 0.22 \end{aligned} \quad (4)$$

$$y = 0; p = 0.8$$

$$\begin{aligned} LS &= -(0 * \log(0.8) + (1 - 0) * \log(0.2)) = \\ &= 1.6 \end{aligned} \quad (5)$$



$$LS = -(y \log(p) + (1 - y) \log(1 - p)) \quad (3)$$

Примеры:

$$y = 1; p = 0.8$$

$$\begin{aligned} LS &= -(1 * \log(0.8) + (1 - 1) * \log(0.2)) = \\ &= 0.22 \end{aligned} \quad (4)$$

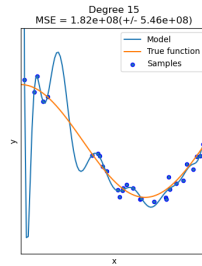
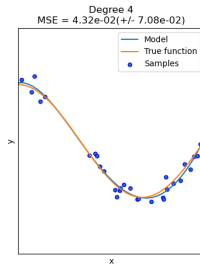
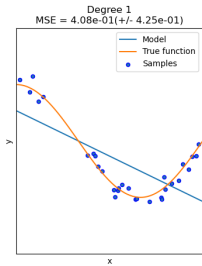
$$y = 0; p = 0.8$$

$$\begin{aligned} LS &= -(0 * \log(0.8) + (1 - 0) * \log(0.2)) = \\ &= 1.6 \end{aligned} \quad (5)$$

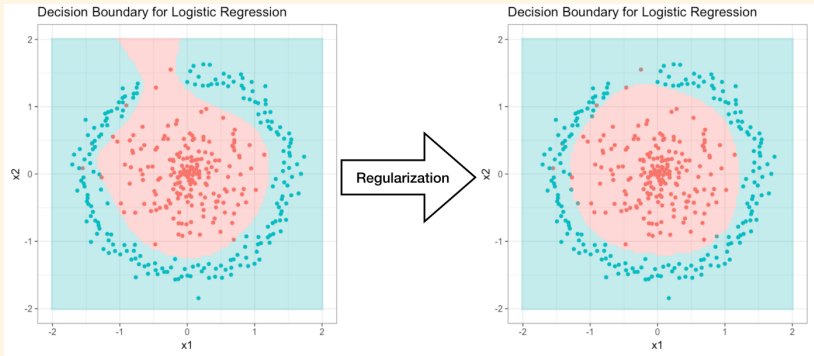
# СНИЖЕНИЕ РАЗМЕРНОСТИ: РЕГУЛЯРИЗАЦИЯ

---

# UNDERFITTING/OVERFITTING



# РЕГУЛЯРИЗАЦИЯ



Обычная регрессия:

$$Loss = Error(y, \hat{y}) \quad (6)$$

L1 Loss (LASSO regression):

$$Loss = Error(y, \hat{y}) + \lambda \sum |\beta_j| \quad (7)$$

L2 Loss (RIDGE regression):

$$Loss = Error(y, \hat{y}) + \lambda \sum \beta_j^2 \quad (8)$$

Обычная регрессия:

$$Loss = Error(y, \hat{y}) \quad (6)$$

L1 Loss (LASSO regression):

$$Loss = Error(y, \hat{y}) + \lambda \sum |\beta_j| \quad (7)$$

L2 Loss (RIDGE regression):

$$Loss = Error(y, \hat{y}) + \lambda \sum \beta_j^2 \quad (8)$$

Обычная регрессия:

$$Loss = Error(y, \hat{y}) \quad (6)$$

L1 Loss (LASSO regression):

$$Loss = Error(y, \hat{y}) + \lambda \sum |\beta_j| \quad (7)$$

L2 Loss (RIDGE regression):

$$Loss = Error(y, \hat{y}) + \lambda \sum \beta_j^2 \quad (8)$$

# SCEPTICAL HAMSTER

21:37

<

skeptischer Hamster zu verk...

☆

↑



**1 skeptischer Hamster zu verkaufen**

**20 €**

📍

 25899 Niebüll >

Art

Hamster

Er guckt einen skeptisch an, als würde man nichts richtig machen.

Es macht mich wahnsinnig, ich kann diesen vorwurfsvollen Blick nicht länger ertragen.

Sein Name ist Olaf.