

# Машинное обучение: классические методы и открытые вопросы

Евгений Соколов

Факультет компьютерных наук НИУ ВШЭ

О чём вообще это всё?

# Правилочный машинный перевод

- Как сделать сервис для перевода, если надо прямо сейчас?

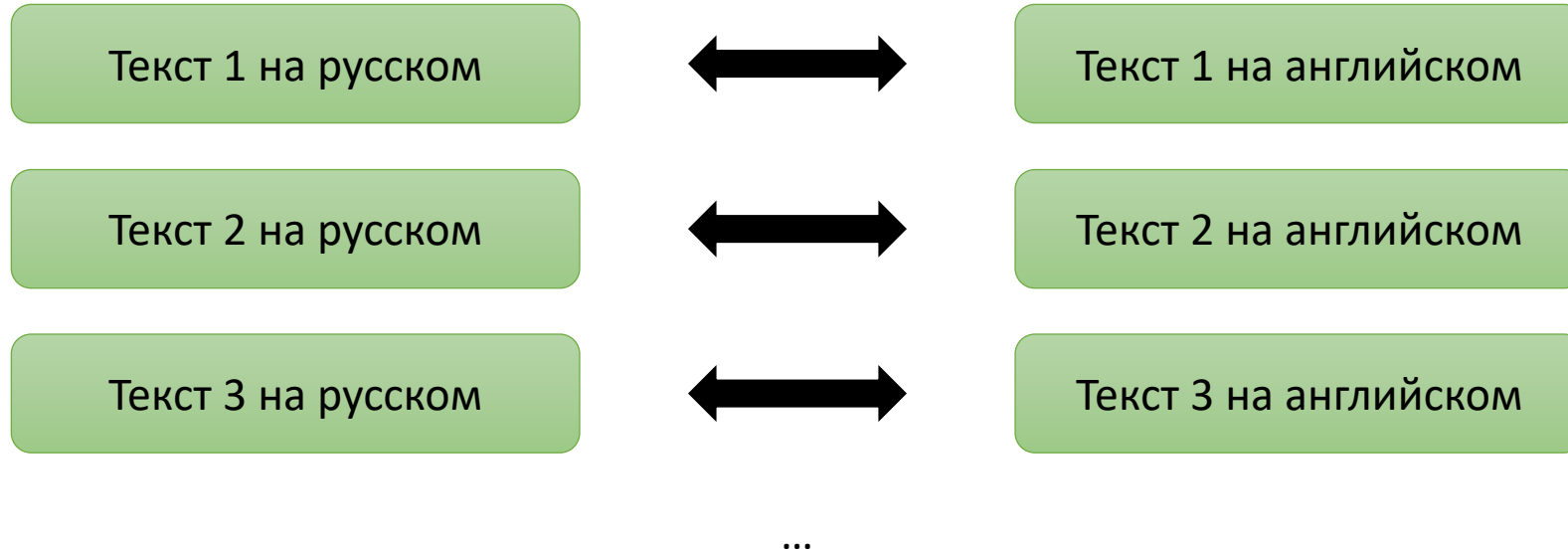
# Правилочный машинный перевод

- Как сделать сервис для перевода, если надо прямо сейчас?
- Перевод по словам
- Грамматические правила

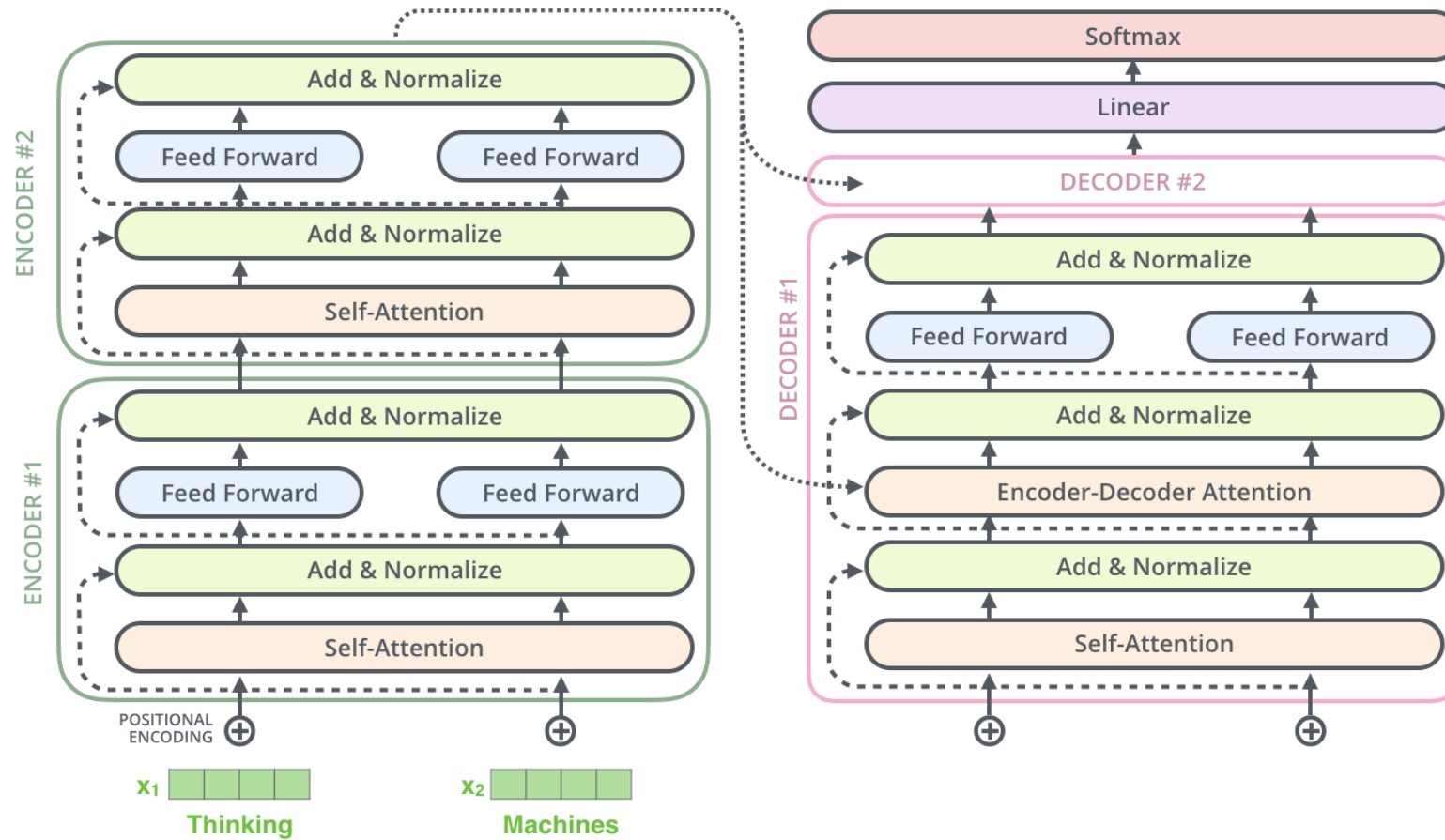
# Правилочный машинный перевод



# Идея: идти от данных, а не от правил



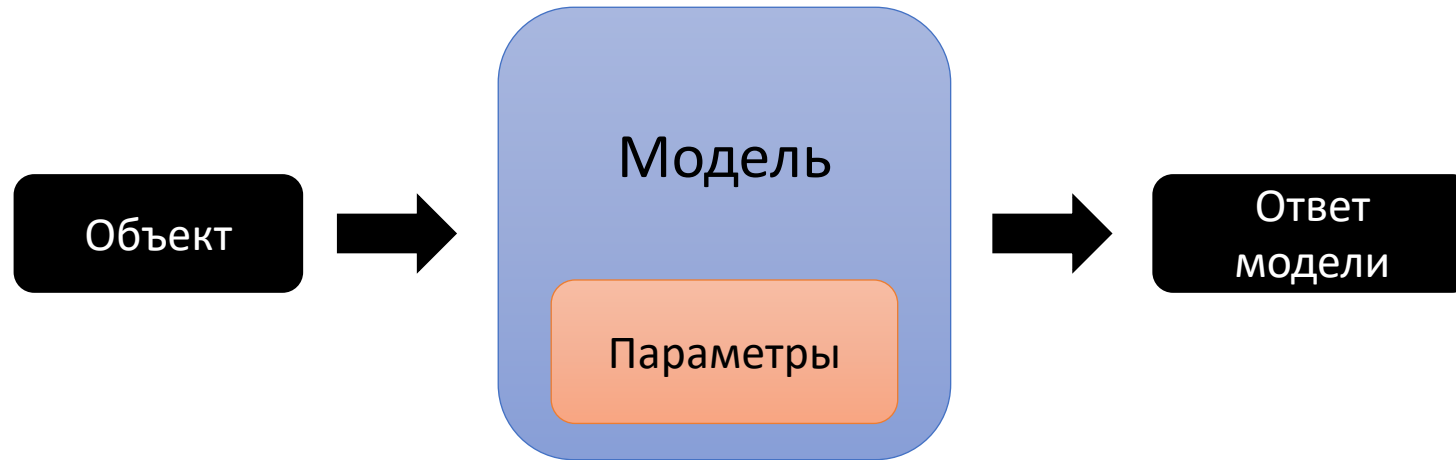
# Нейросетевой машинный перевод



Теперь более строго

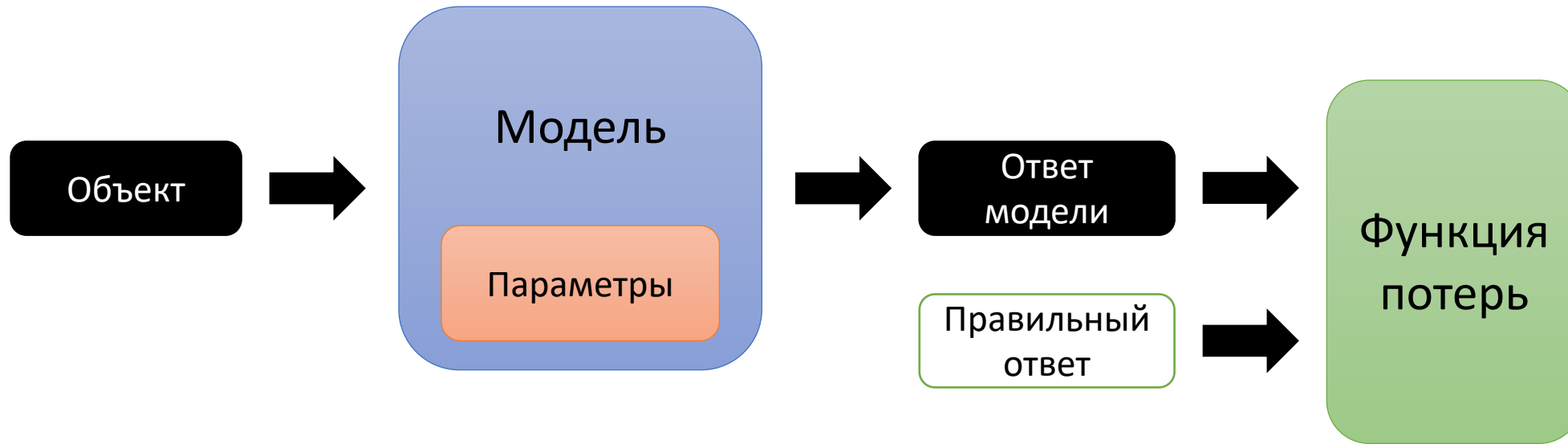


# Общая схема

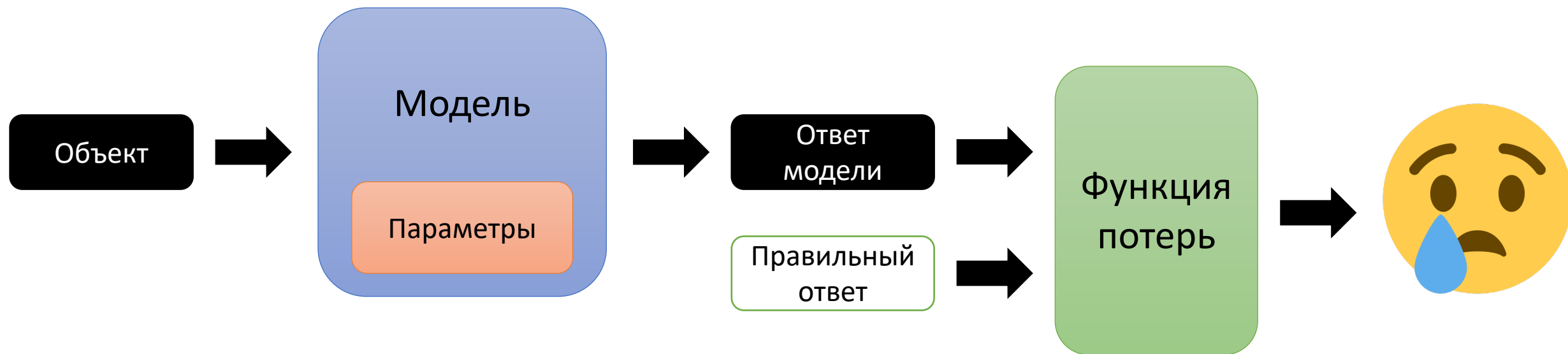


- Модель — формула или алгоритм для решения задачи
- Параметры — «ручки», влияющие на поведение модели

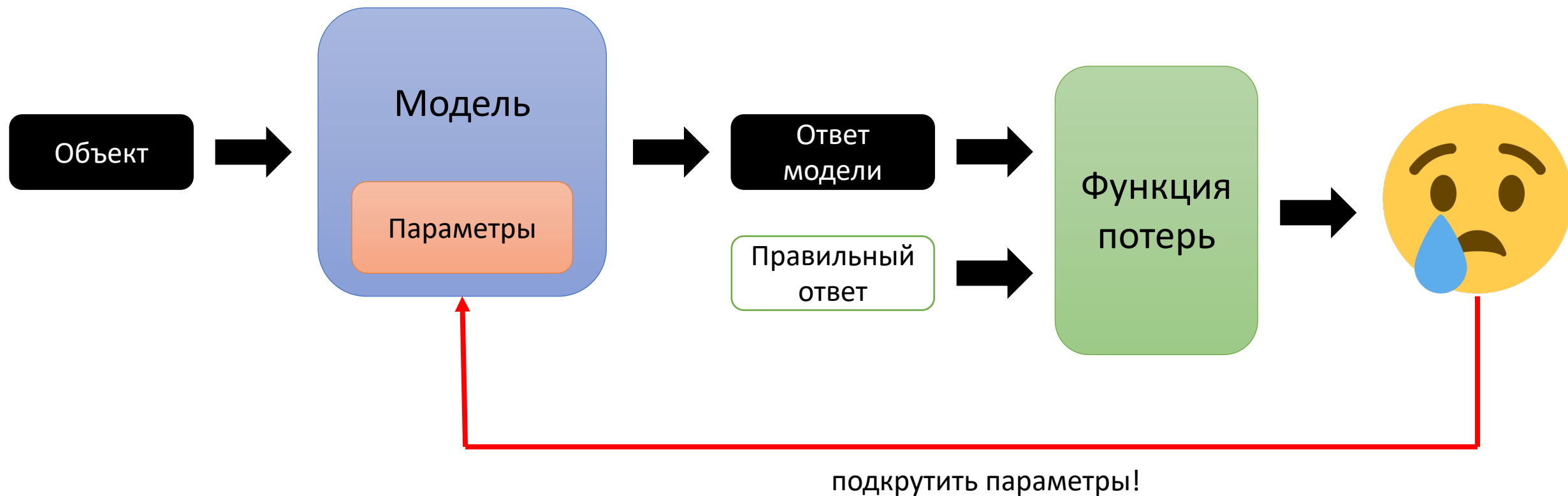
# Общая схема



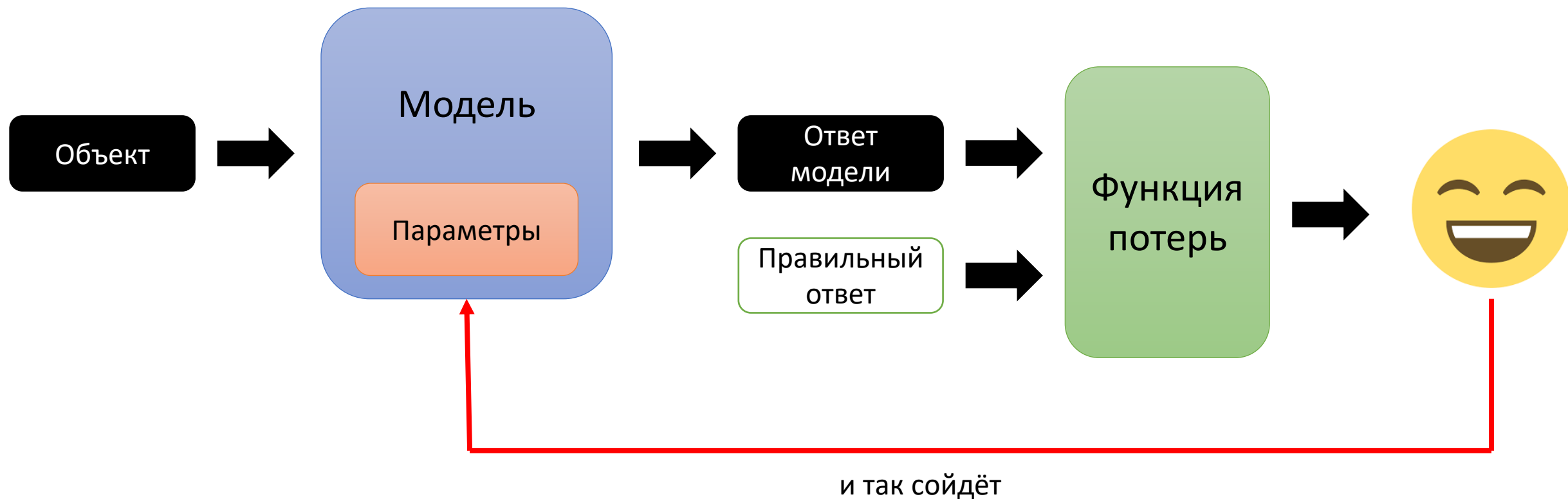
# Общая схема



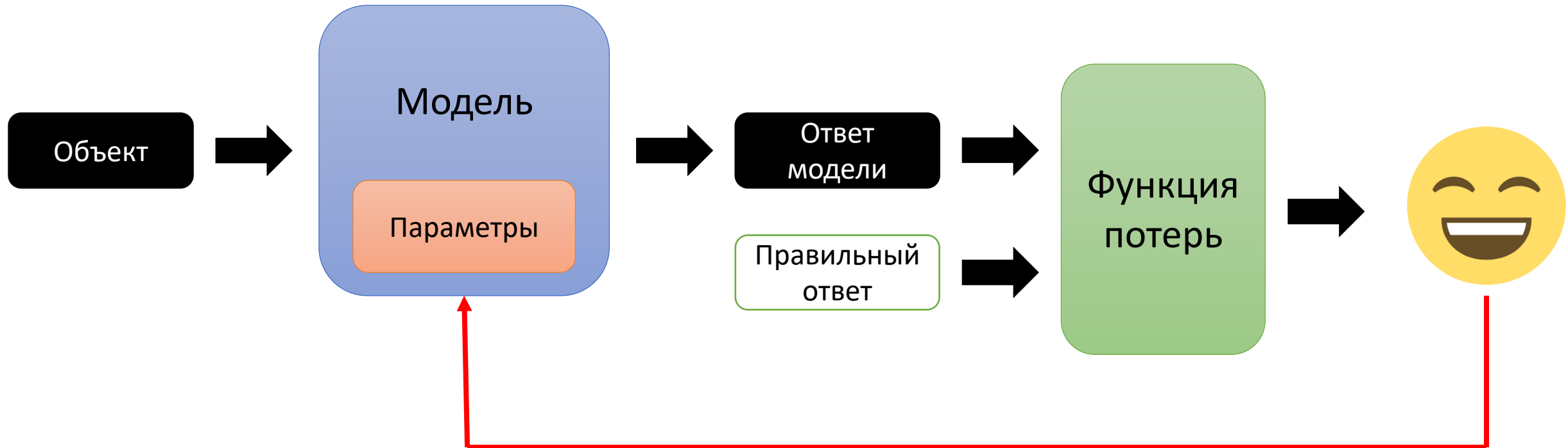
# Общая схема



# Общая схема



# Общая схема

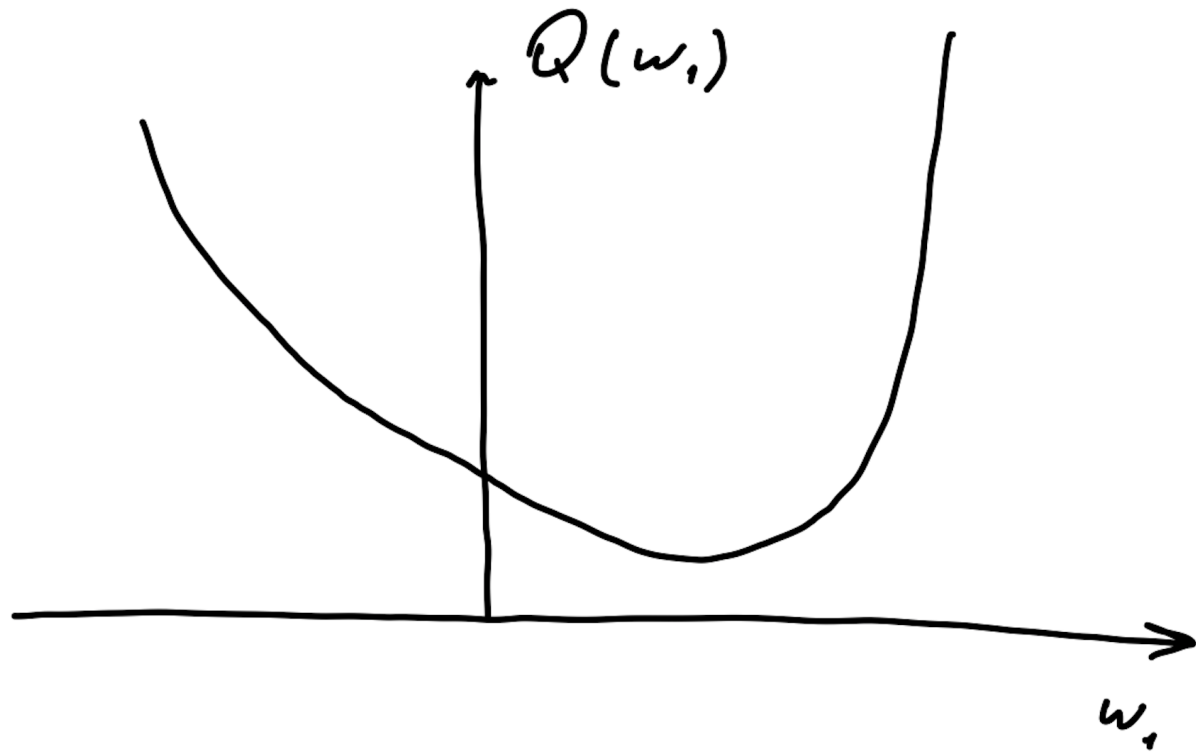


По очереди проделываем для всей выборки, пока ошибка уменьшается

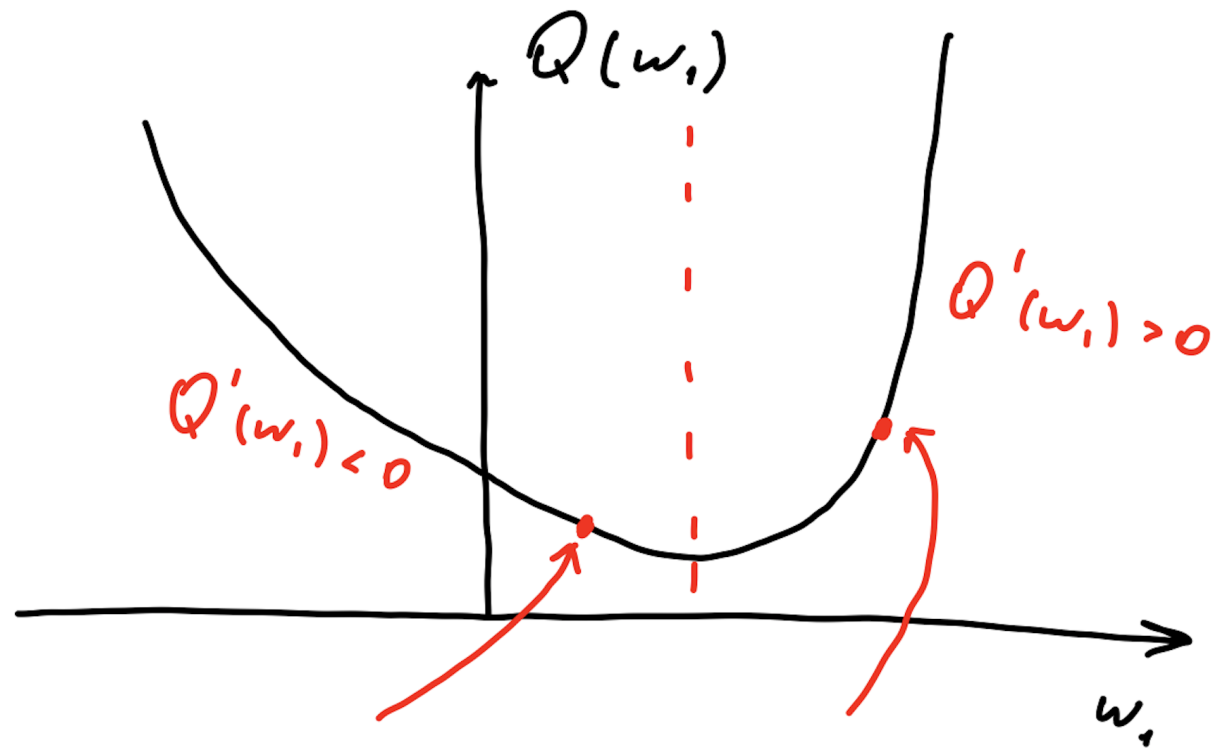
# Более формально

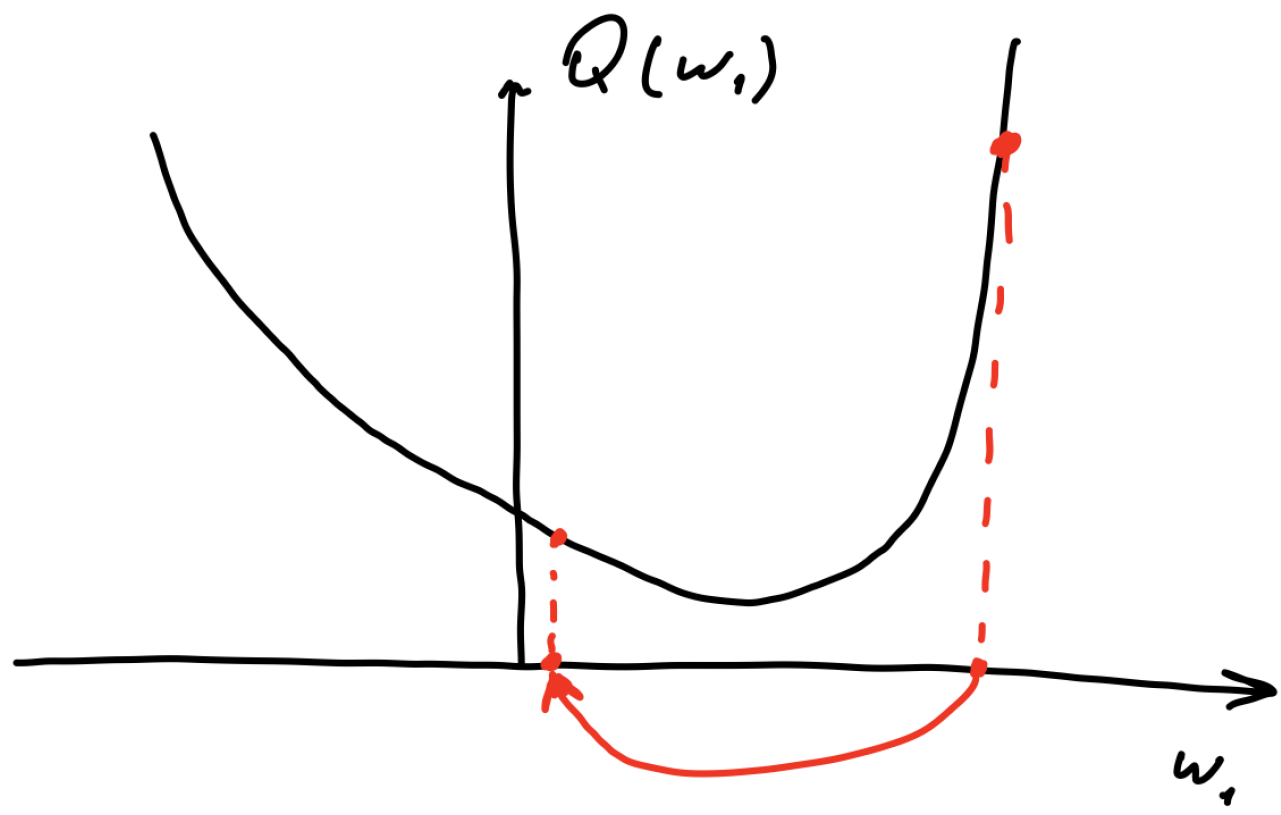
$$Q(w) = \sum_{i=1}^{\ell} L(y_i, a(x_i, w)) \rightarrow \min_w$$

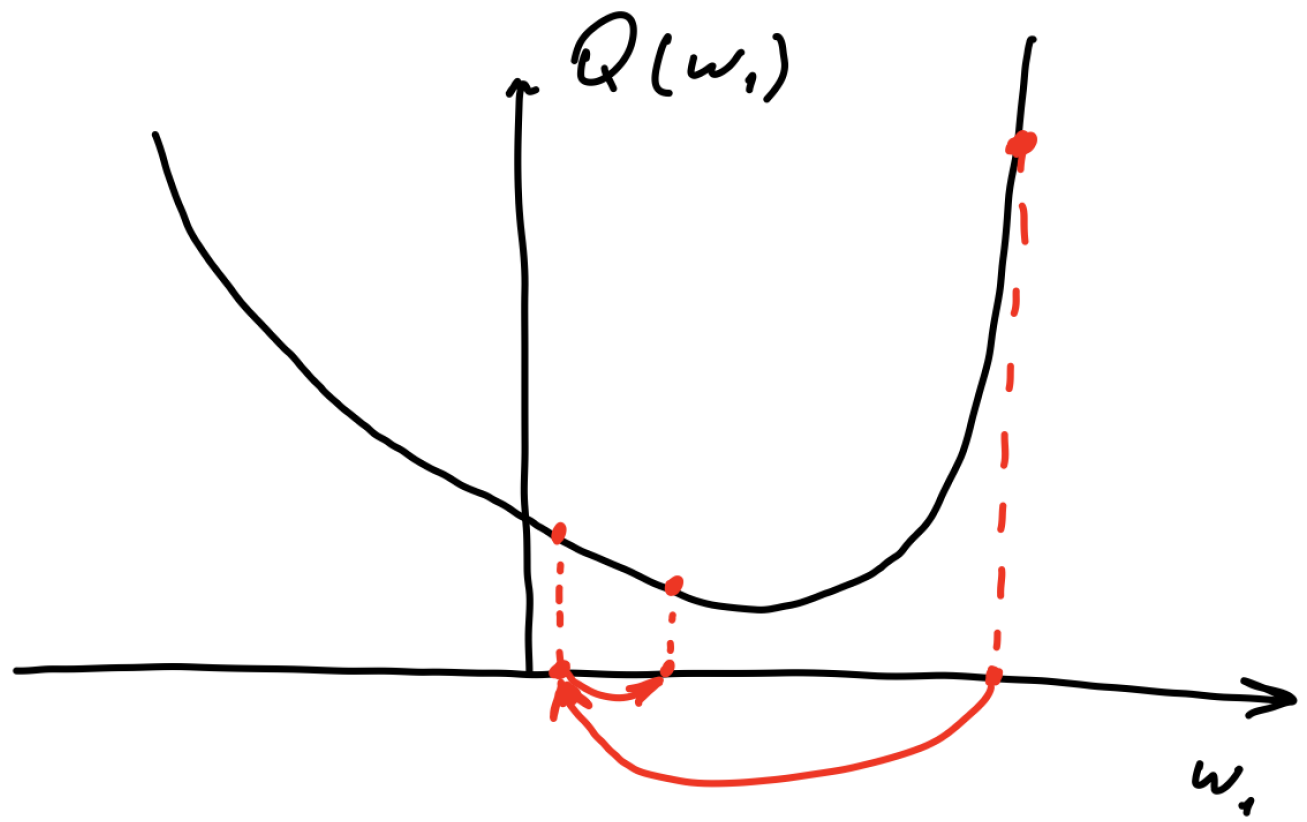
- $x_i, y_i$  — объект и правильный ответ
- $a(x, w)$  — модель с параметрами  $w$
- $L(y, z)$  — функция потерь

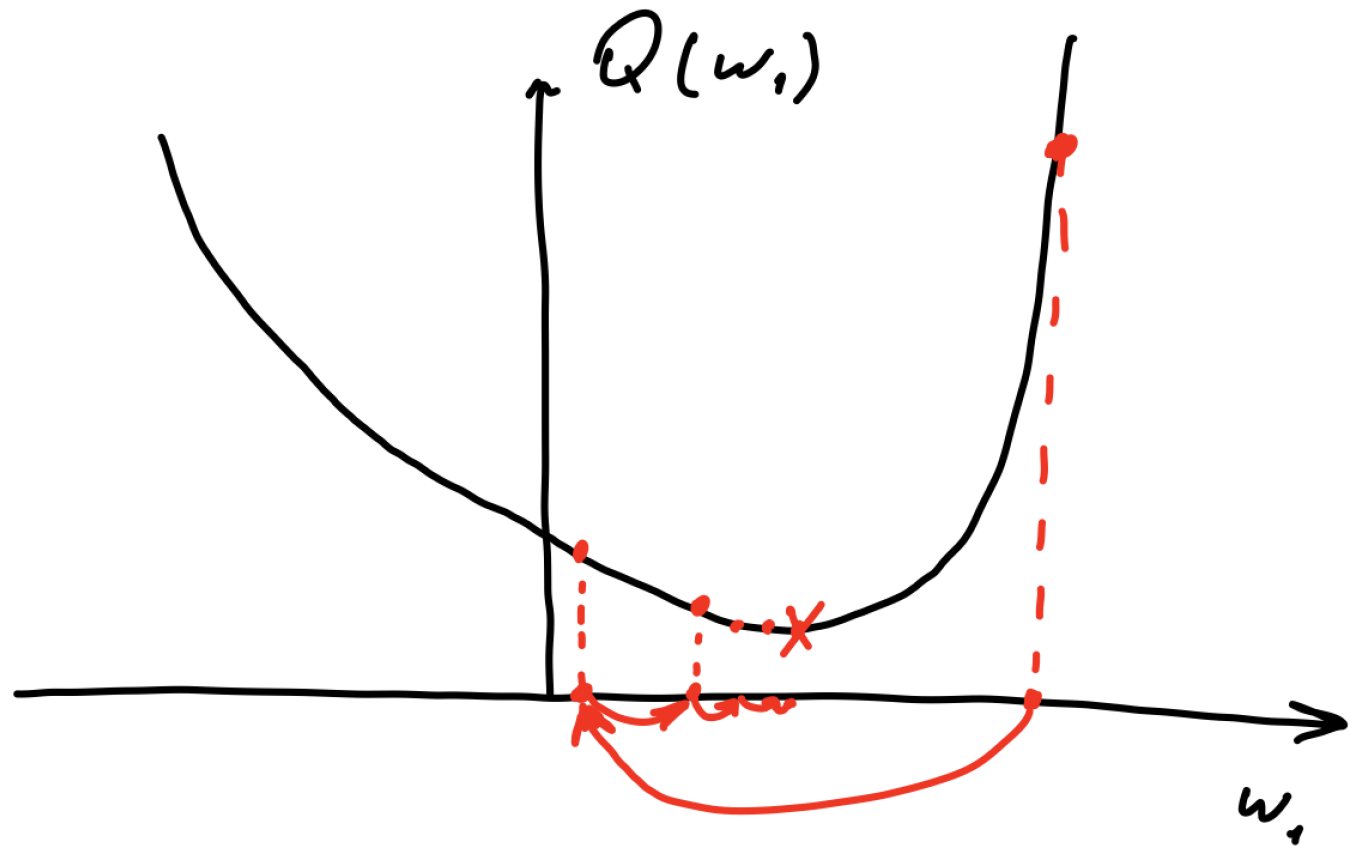


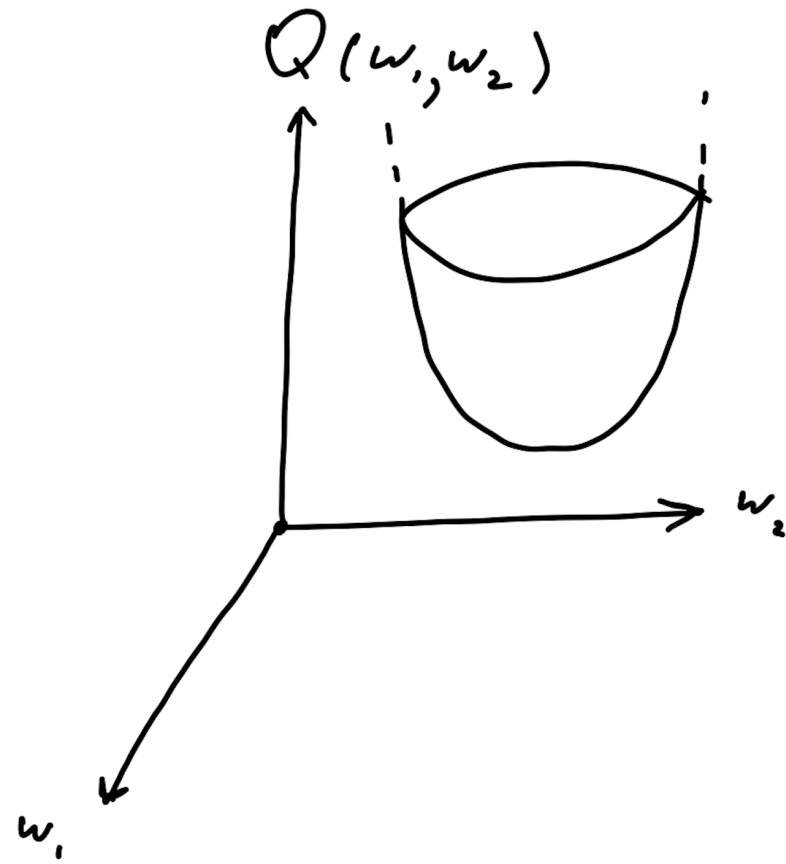


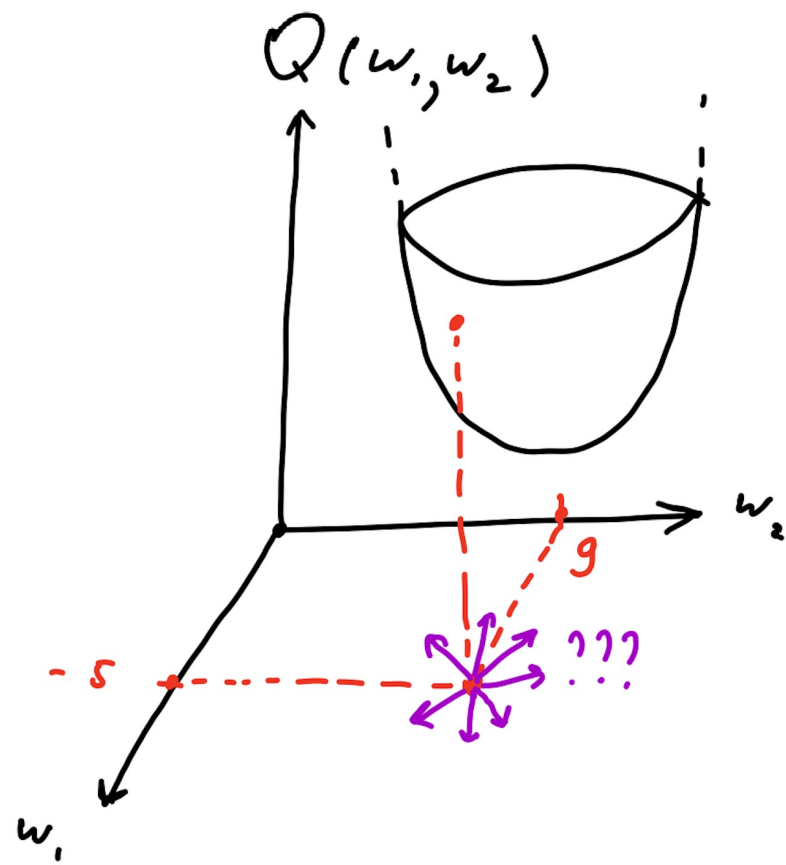


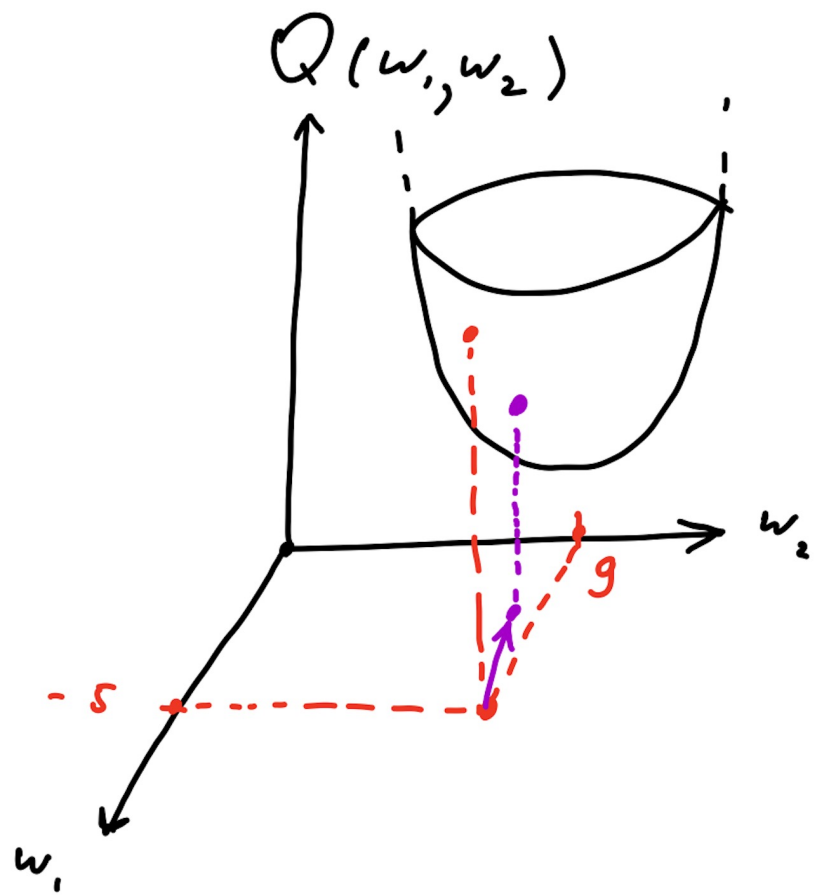


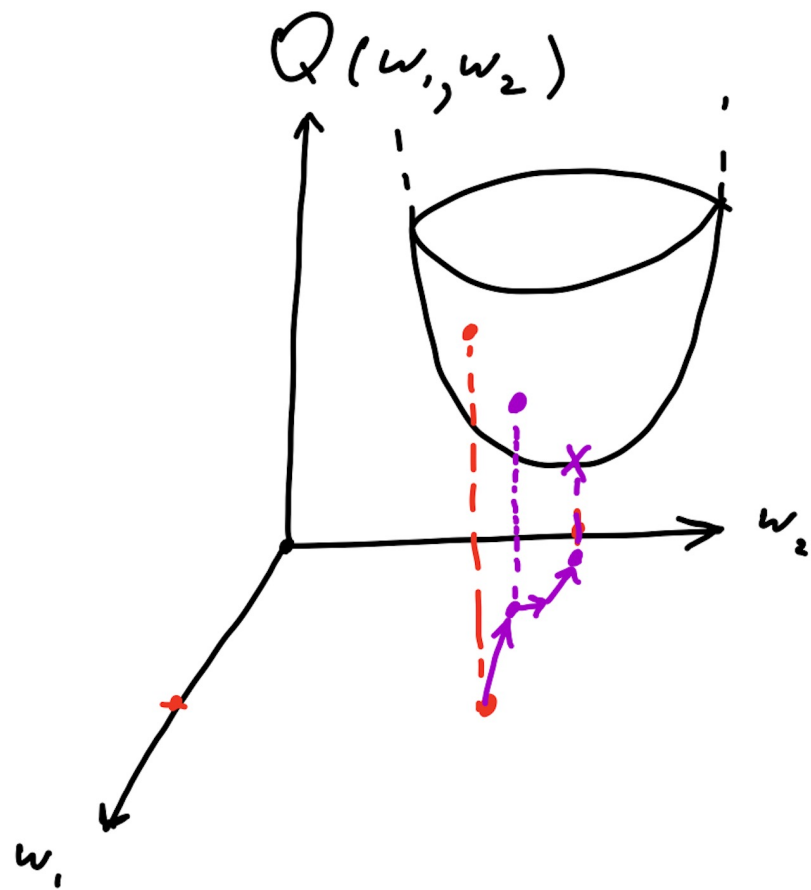










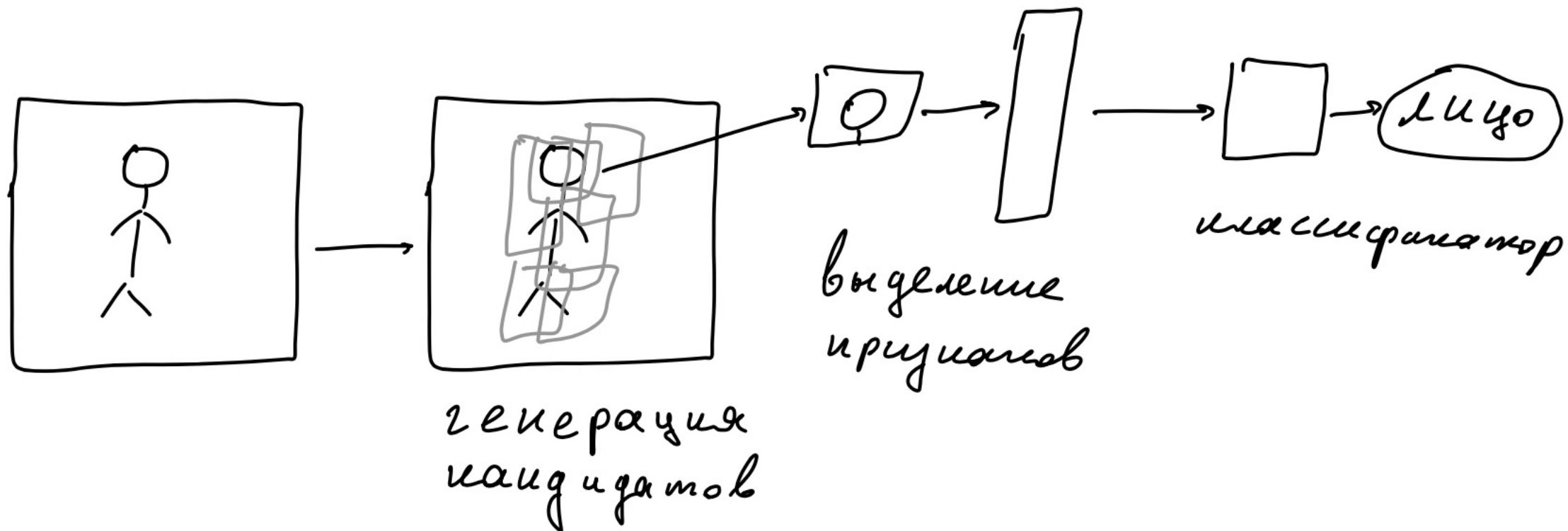




Как работать с  
изображениями?

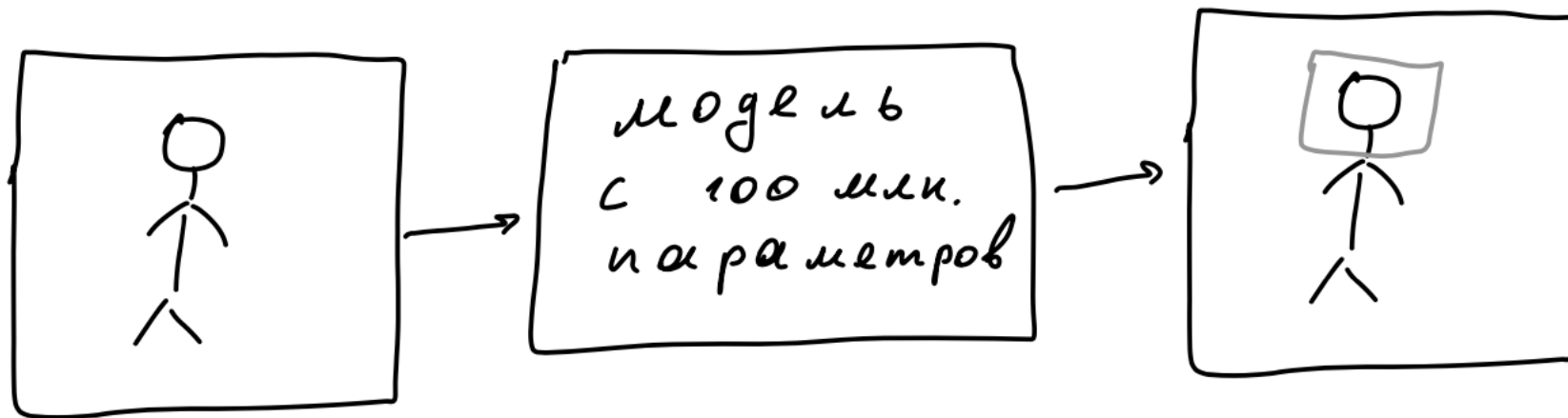
# Принцип 1: end-to-end обучение

Детекция объектов раньше:



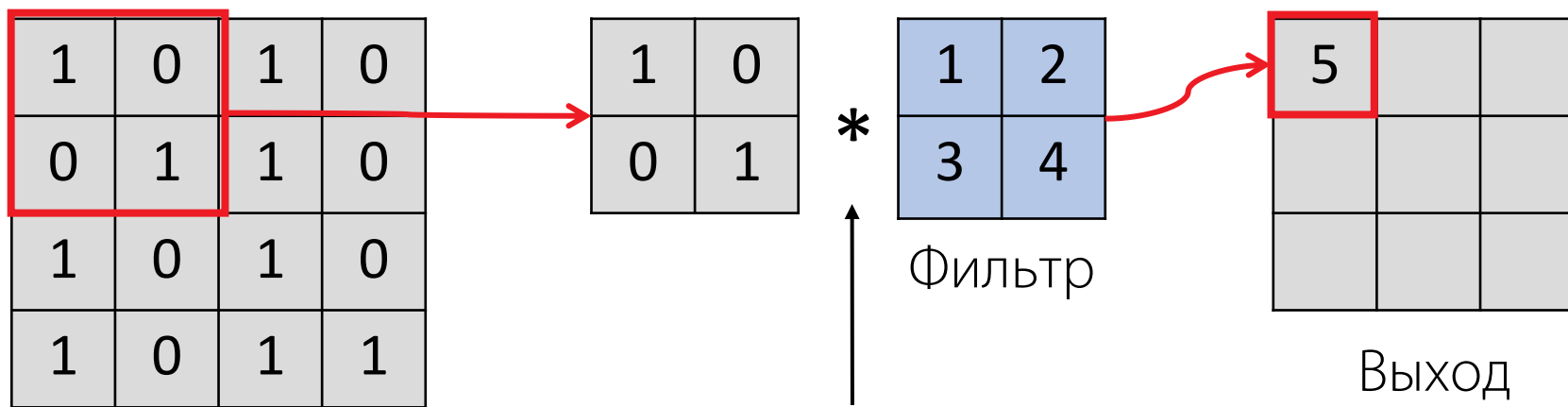
# Принцип 1: end-to-end обучение

Детекция объектов сегодня:

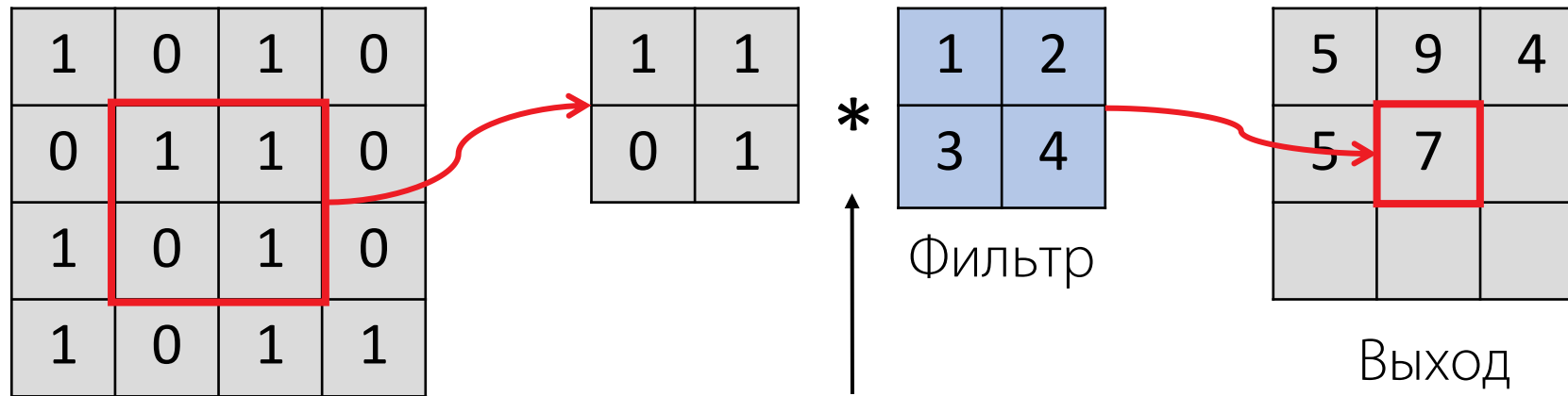


Принцип 2: многократное извлечение признаков

# Свёртка



# Свёртка



*Поэлементное  
умножение,  
затем  
суммирование*

# Свёртка

$$\begin{array}{|c|c|} \hline 1 & 1 \\ \hline 0 & 1 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \boxed{2}$$

$$\begin{array}{|c|c|} \hline 1 & 1 \\ \hline 1 & 1 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \boxed{2}$$

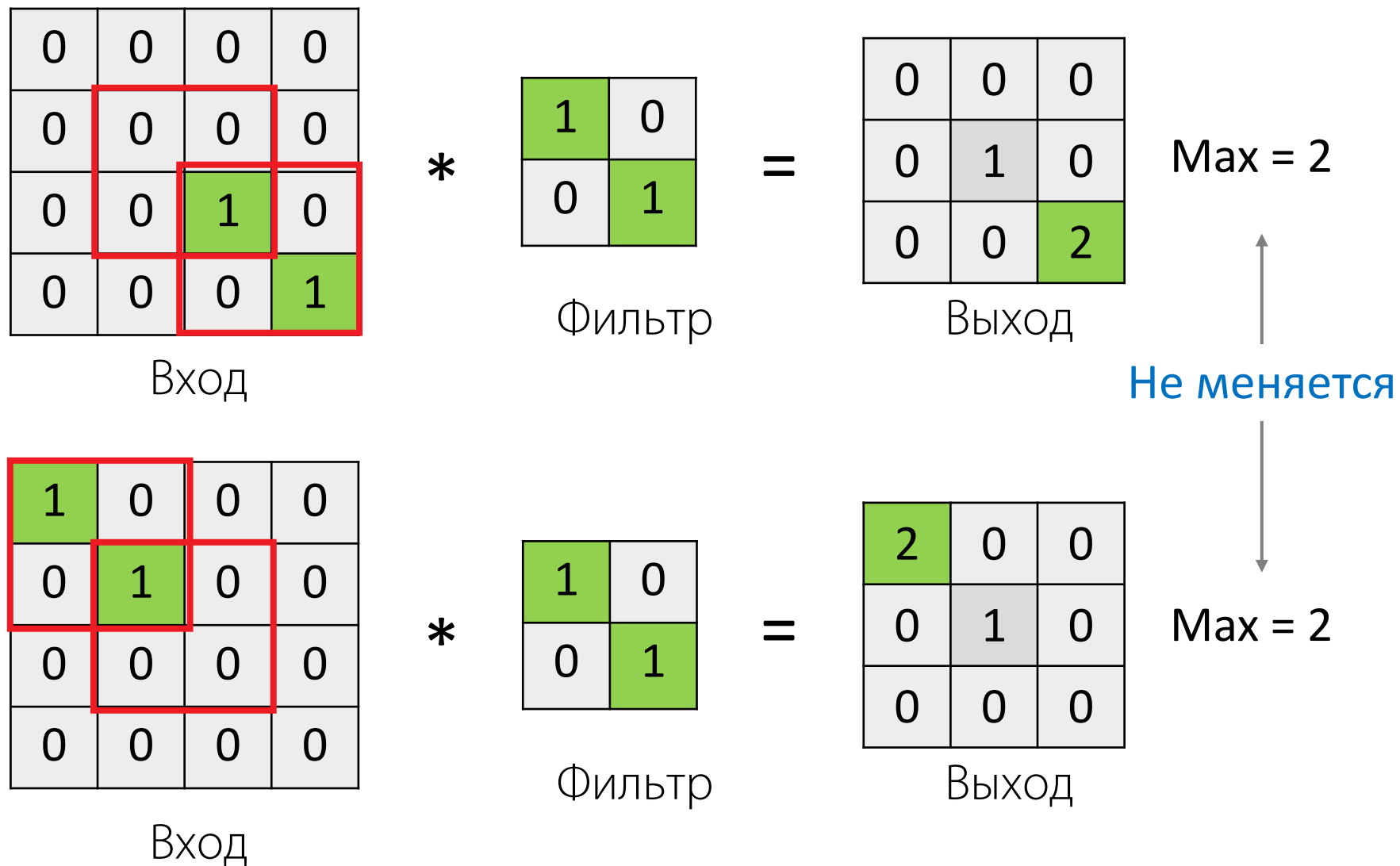
$$\begin{array}{|c|c|} \hline 3 & 0 \\ \hline 0 & 3 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \boxed{6}$$

$$\begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 0 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \boxed{1}$$

$$\begin{array}{|c|c|} \hline 5 & 0 \\ \hline 0 & 5 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \boxed{10}$$

$$\begin{array}{|c|c|} \hline 0 & 2 \\ \hline 3 & 0 \\ \hline \end{array} * \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 1 \\ \hline \end{array} = \boxed{0}$$

# Максимум свёртки инвариантен к сдвигам



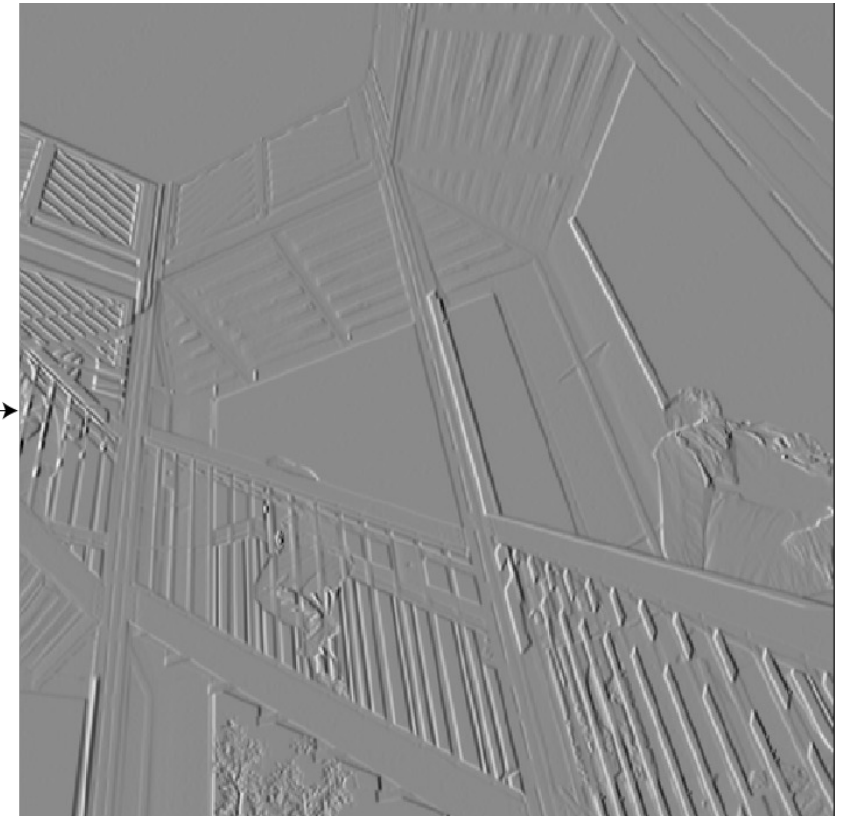


# Свёртки в компьютерном зрении



$$\begin{bmatrix} +1 & 0 & -1 \\ +2 & 0 & -2 \\ +1 & 0 & -1 \end{bmatrix}$$

Horizontal Sobel kernel

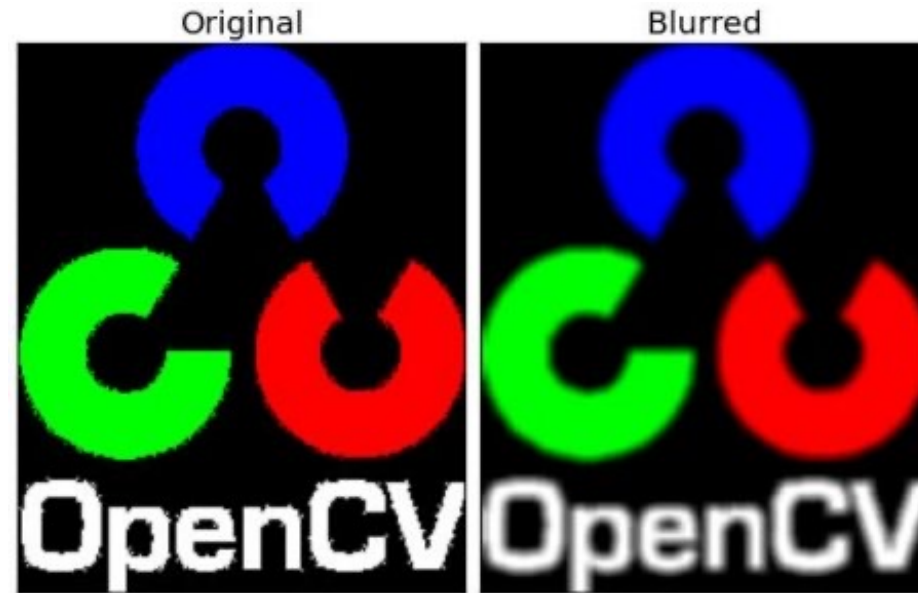


# Свёртки в компьютерном зрении



$$\begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 1 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix} + \begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 1 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix} - \frac{1}{9} \begin{bmatrix} \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \end{bmatrix} = \begin{bmatrix} \bullet 0 & \bullet 0 & \bullet 0 \\ \bullet 0 & \bullet 2 & \bullet 0 \\ \bullet 0 & \bullet 0 & \bullet 0 \end{bmatrix} - \frac{1}{9} \begin{bmatrix} \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \\ \bullet 1 & \bullet 1 & \bullet 1 \end{bmatrix}$$

# Свёртки в компьютерном зрении

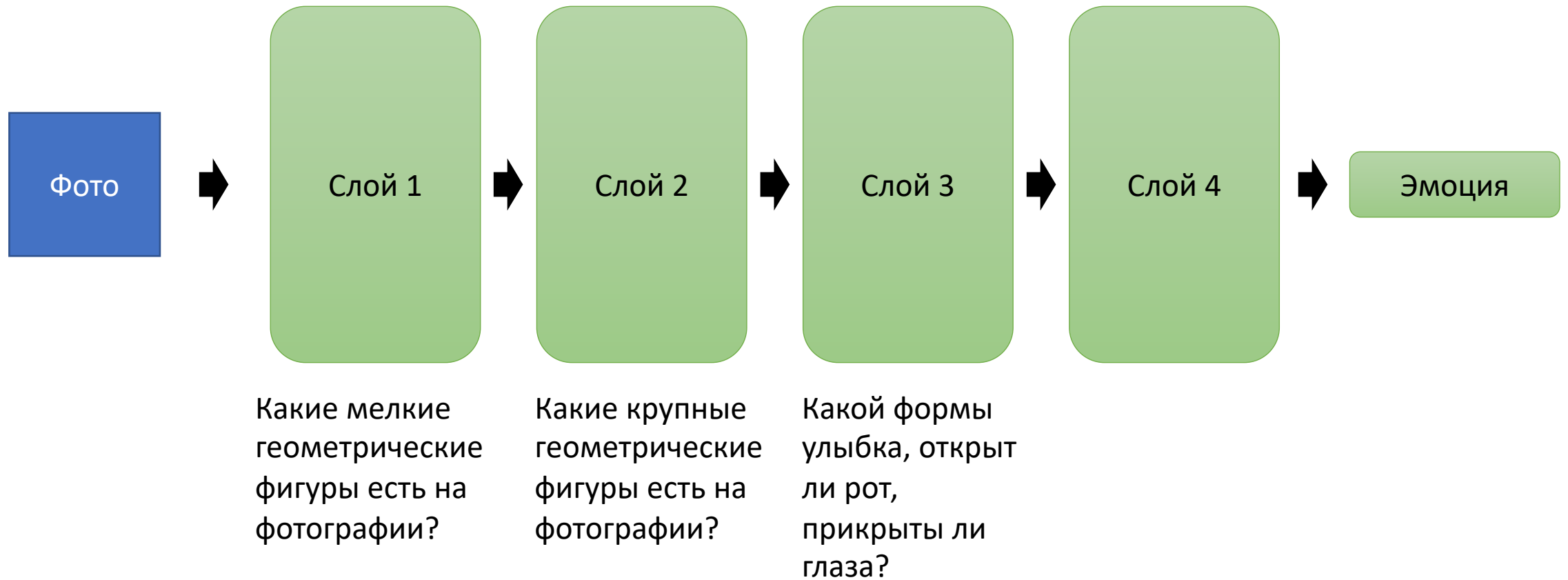


$$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

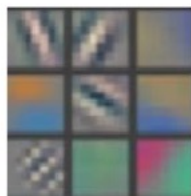
# Свёрточные сети

- Из изображения выделяются всё более верхнеуровневые признаки с помощью свёрток
- Фильтры в свёртках **обучаются**

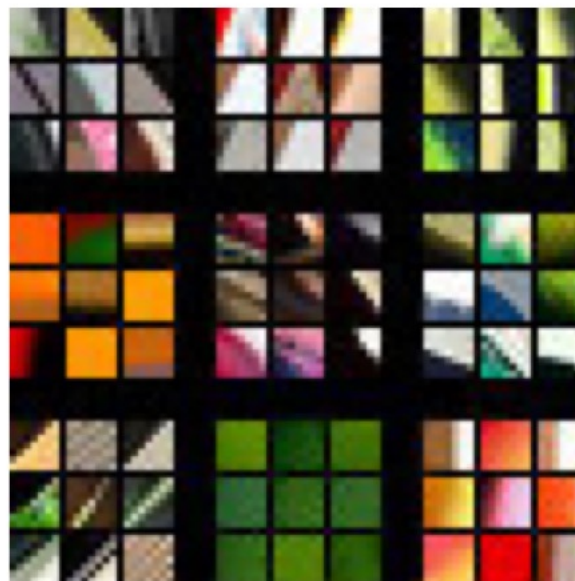
# Свёрточные сети



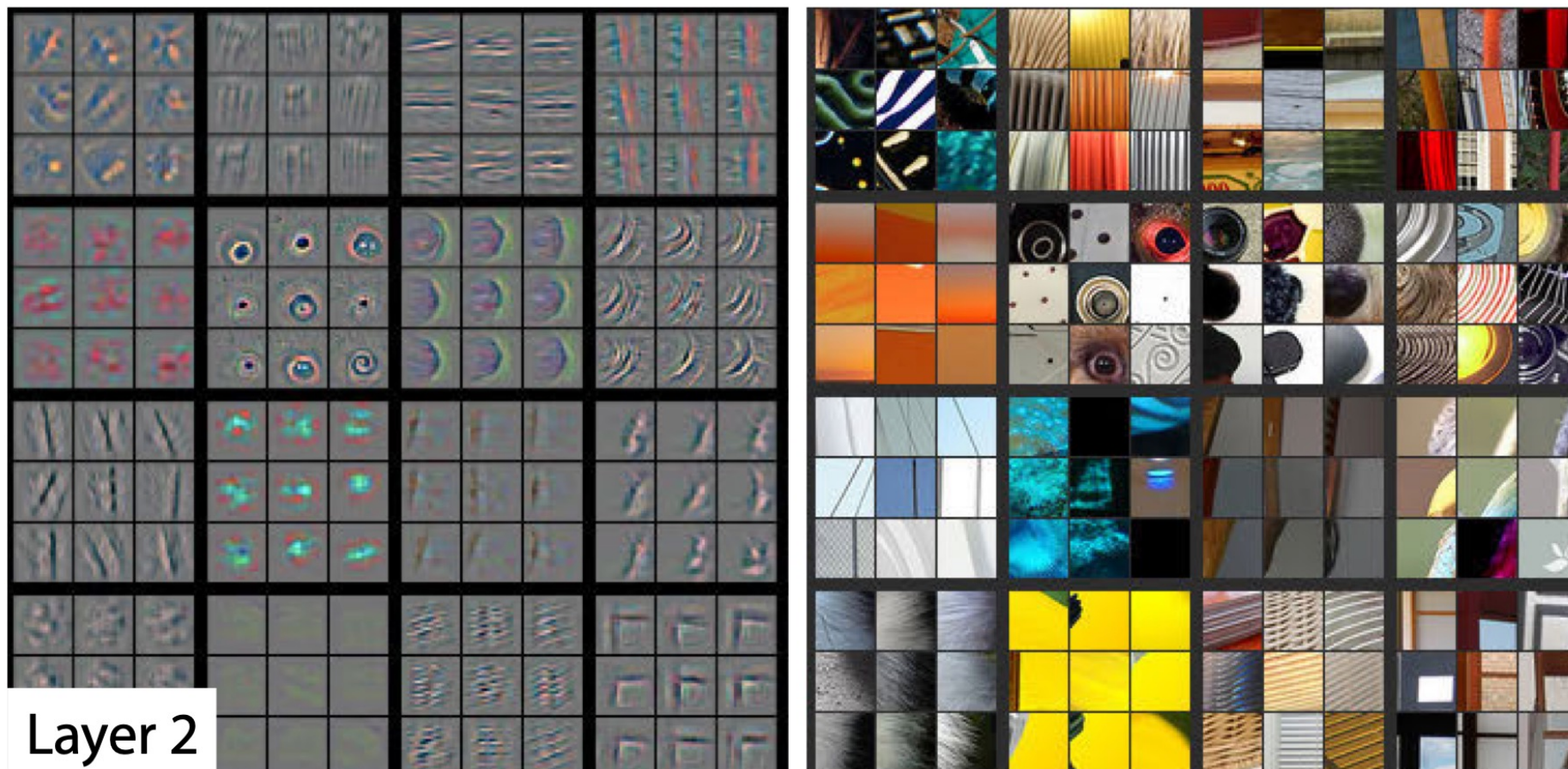
# Что видят фильтры?



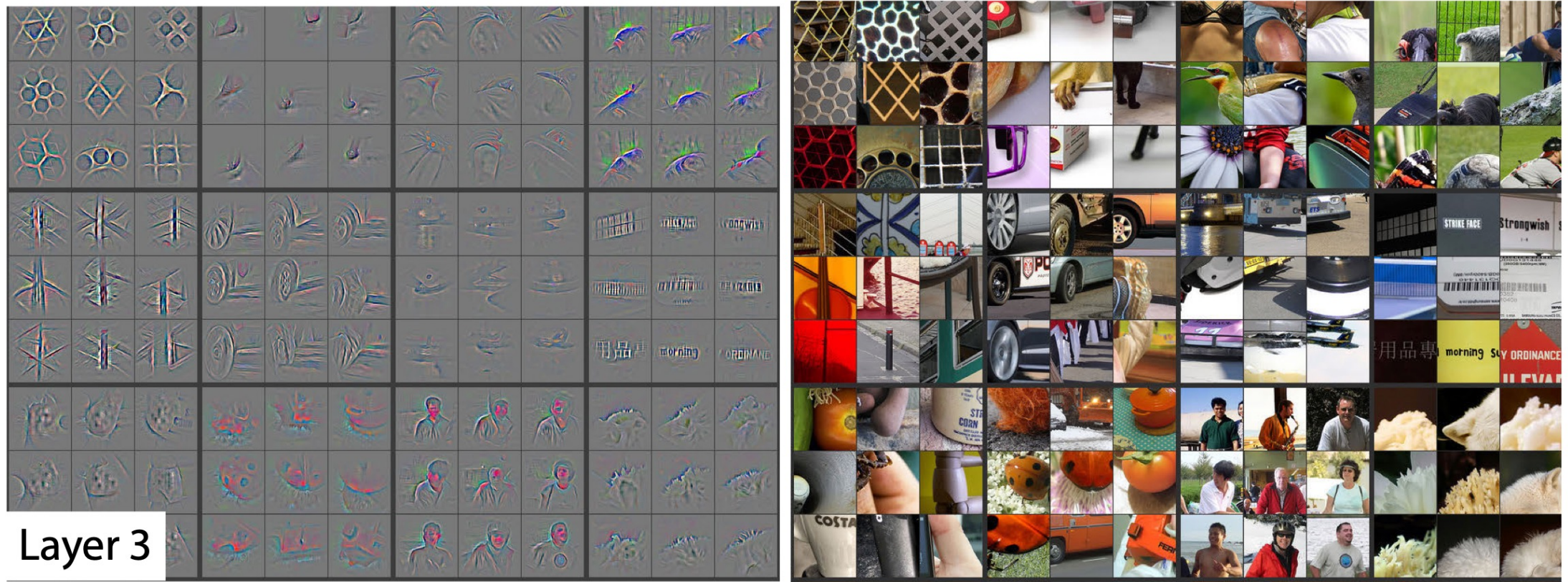
Layer 1



# Что видят фильтры?

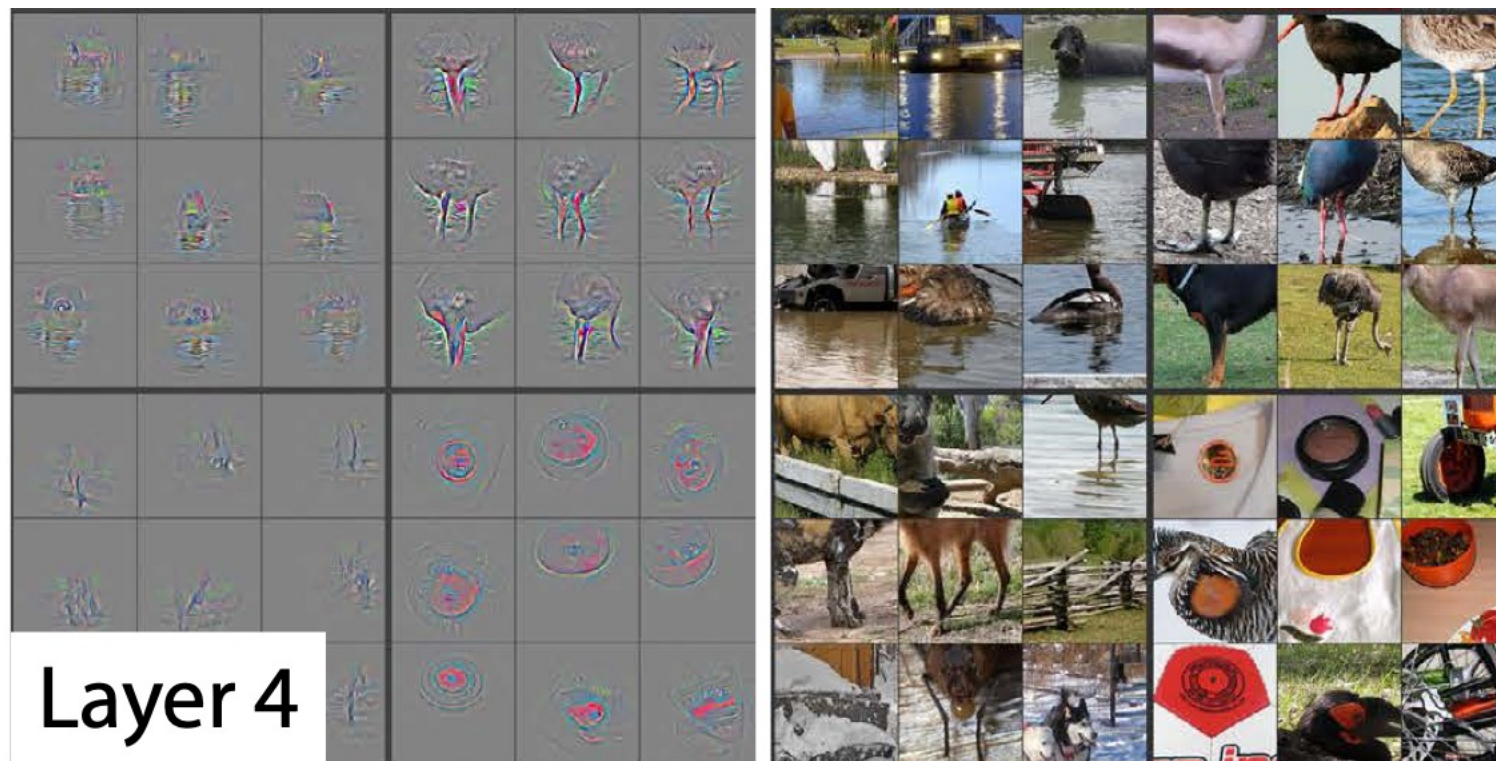


# Что видят фильтры?

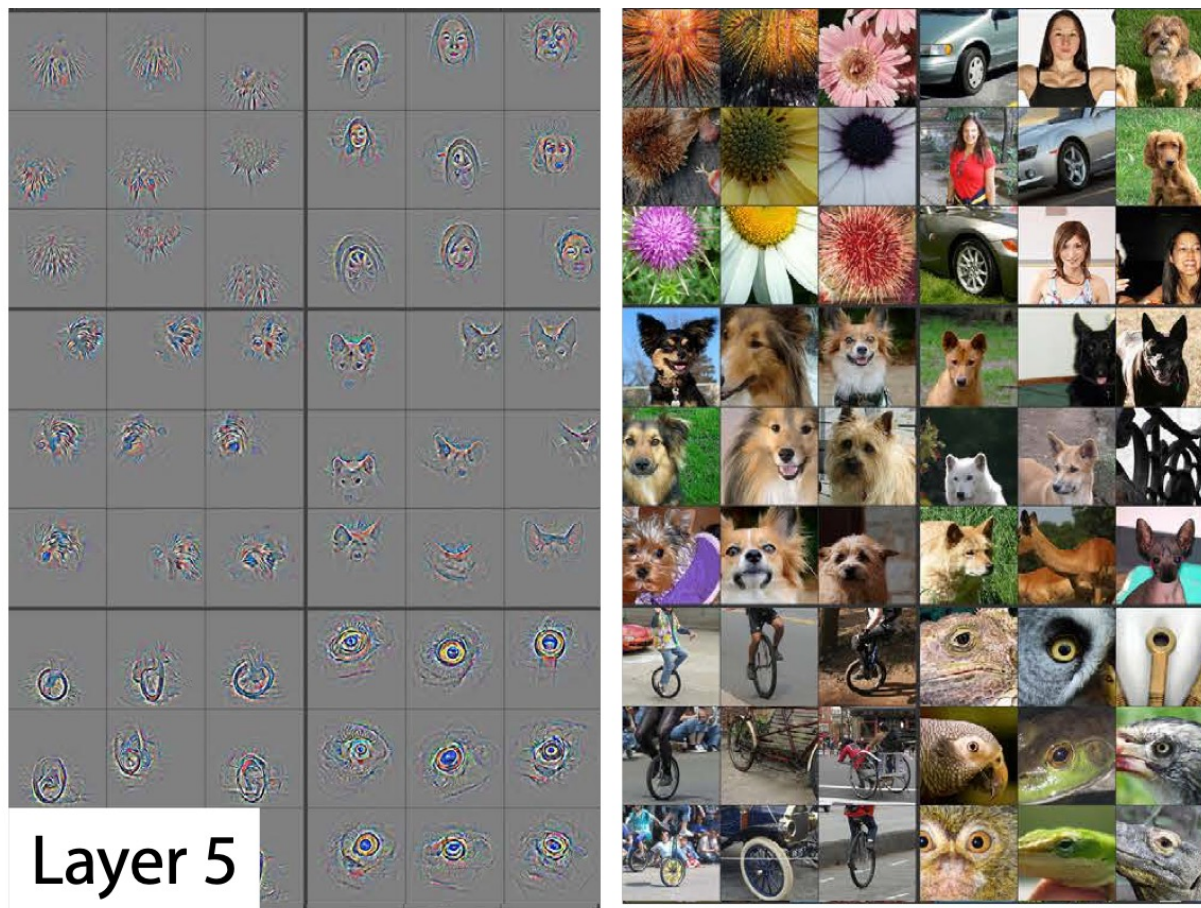




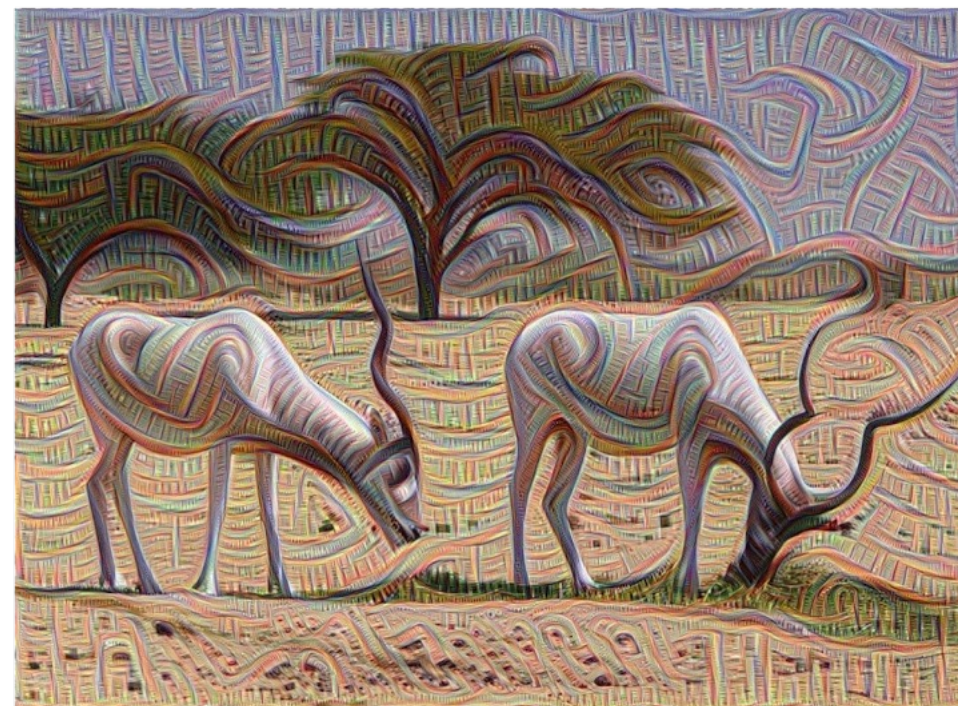
# Что видят фильтры?



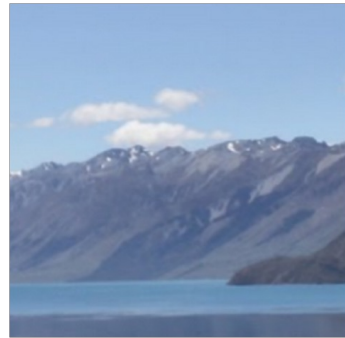
# Что видят фильтры?



# Максимизация вероятности класса



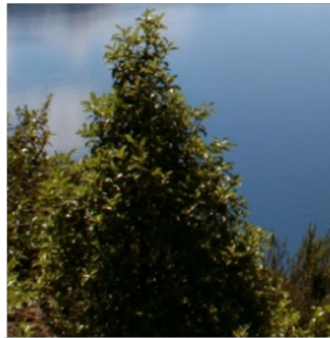
# Максимизация вероятности класса



Horizon



Towers & Pagodas



Trees



Buildings

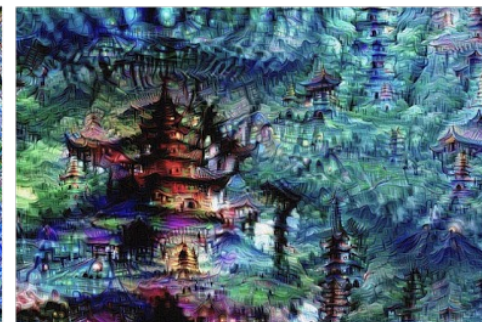
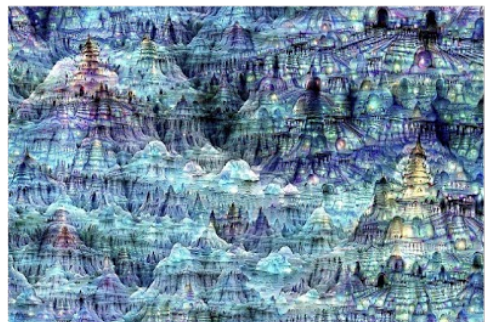
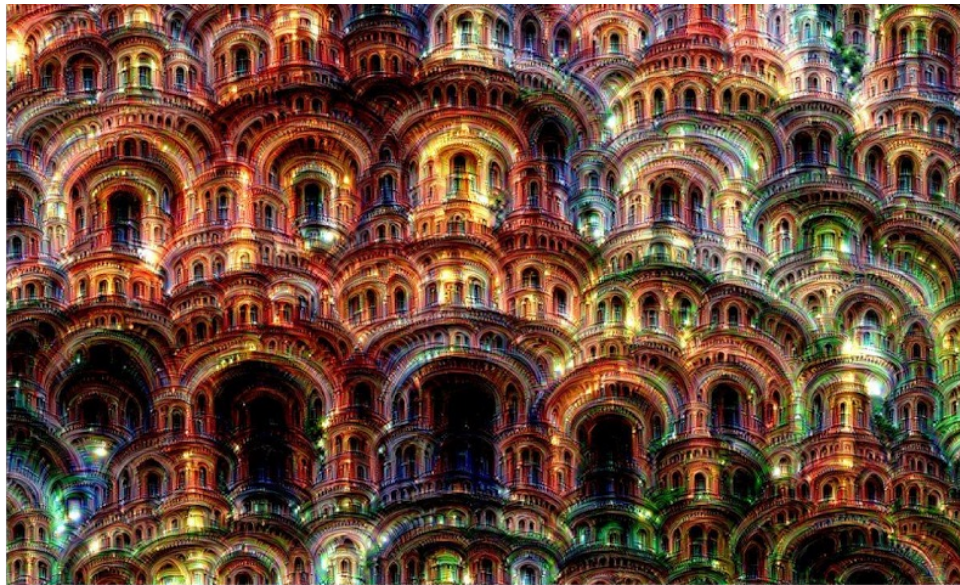
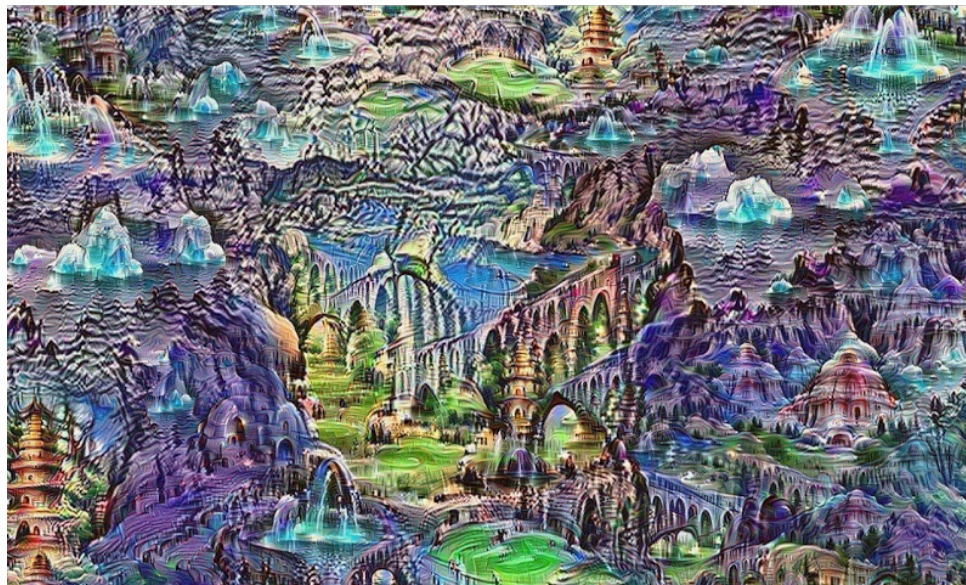


Leaves



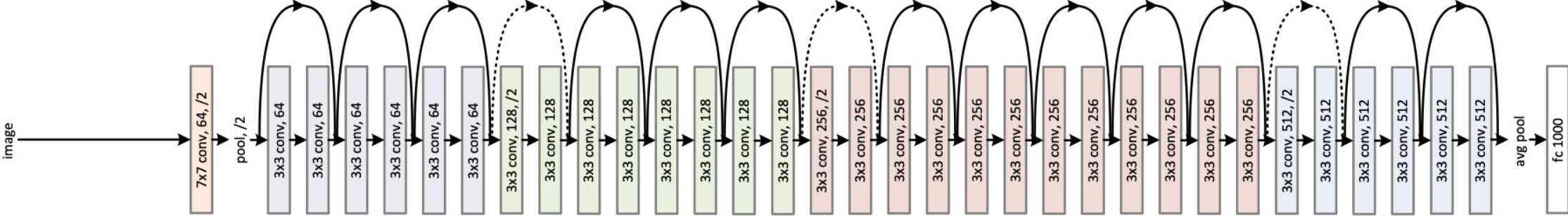
Birds & Insects

# Максимизация вероятности класса

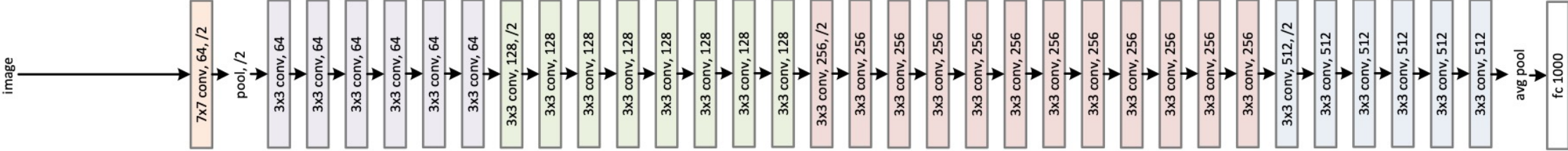


# ResNet (2015)

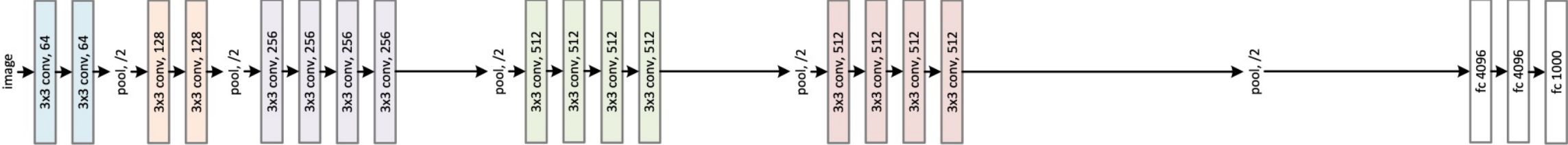
34-layer residual



34-layer plain



VGG-19



Необычные свойства  
функций потерь

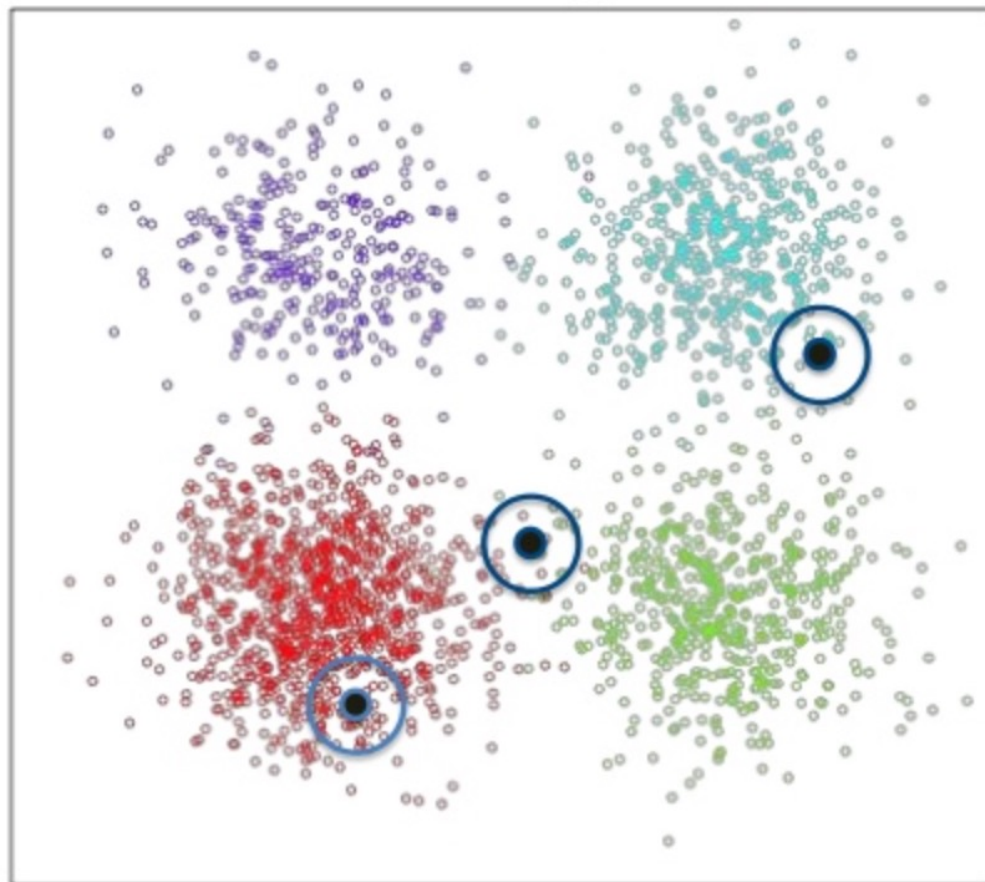
# Функция потерь

$$Q(w) = \sum_{i=1}^{\ell} L(y_i, a(x_i, w)) \rightarrow \min_w$$

- $x_i, y_i$  — объект и правильный ответ
- $a(x, w)$  — модель с параметрами  $w$
- $L(y, z)$  — функция потерь



# Гипотеза компактности



# kNN: обучение

- Дано: обучающая выборка  $X = (x_i, y_i)_{i=1}^{\ell}$
- Задача классификация (ответы из множества  $\mathbb{Y} = \{1, \dots, K\}$ )
  
- Обучение модели:
  - Запоминаем обучающую выборку  $X$

# kNN: применение

Дано: новый объект  $x$

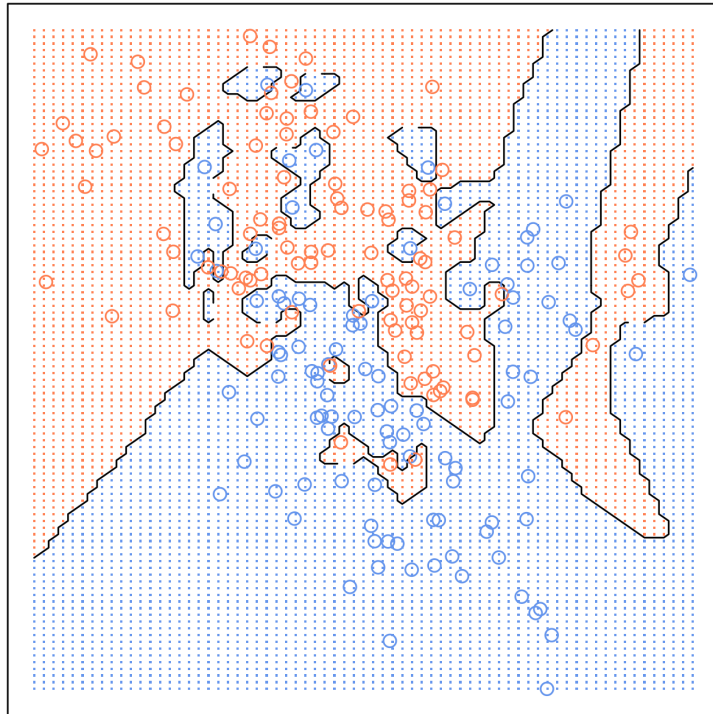
Применение модели:

- Сортируем объекты обучающей выборки по расстоянию до нового объекта:  
 $\rho(x, x_{(1)}) \leq \rho(x, x_{(2)}) \leq \dots \leq \rho(x, x_{(\ell)})$
- Выбираем  $k$  ближайших объектов:  $x_{(1)}, \dots, x_{(k)}$
- Выдаём наиболее популярный среди них класс:

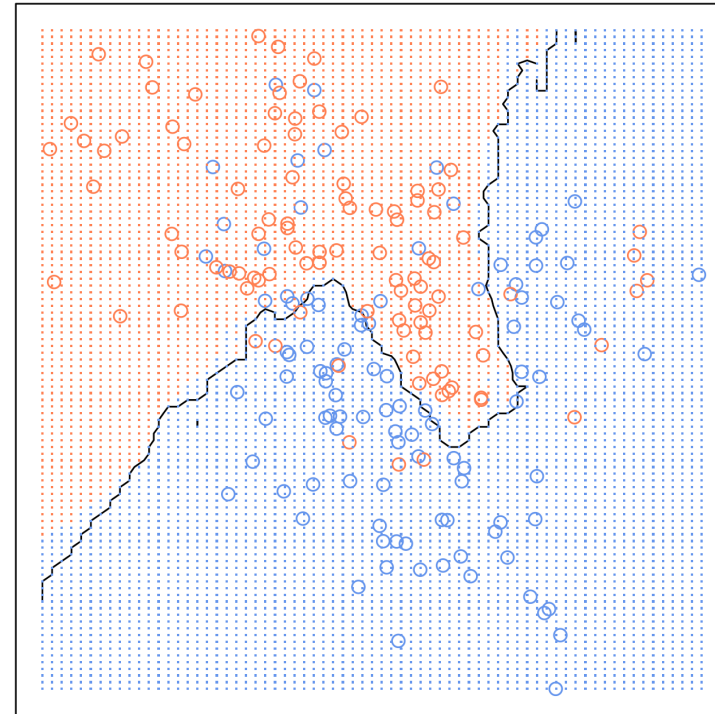
$$a(x) = \arg \max_{y \in \mathbb{Y}} \sum_{i=1}^k [y^{(i)} = y]$$

# Как выбрать $k$ ?

1-nearest neighbours



20-nearest neighbours



<https://kevinzakka.github.io/2016/07/13/k-nearest-neighbor/>

# Обобщающая способность

Как готовиться к экзамену?

Заучить все примеры с  
занятий

Разобраться в предмете и  
усвоить алгоритмы решения  
задач

# Обобщающая способность

Как готовиться к экзамену?

Заучить все примеры с  
занятий

Переобучение (overfitting)

Разобраться в предмете и  
усвоить алгоритмы решения  
задач

Обобщение (generalization)

# Обобщающая способность

Как готовиться к экзамену?

Заучить все примеры с  
занятий

Переобучение (overfitting)

Хорошее качество на обучении  
Низкое качество на новых данных

Разобраться в предмете и  
усвоить алгоритмы решения  
задач

Обобщение (generalization)

Хорошее качество на обучении  
Хорошее качество на новых  
данных

# Отложенная выборка



Обучение



Тест



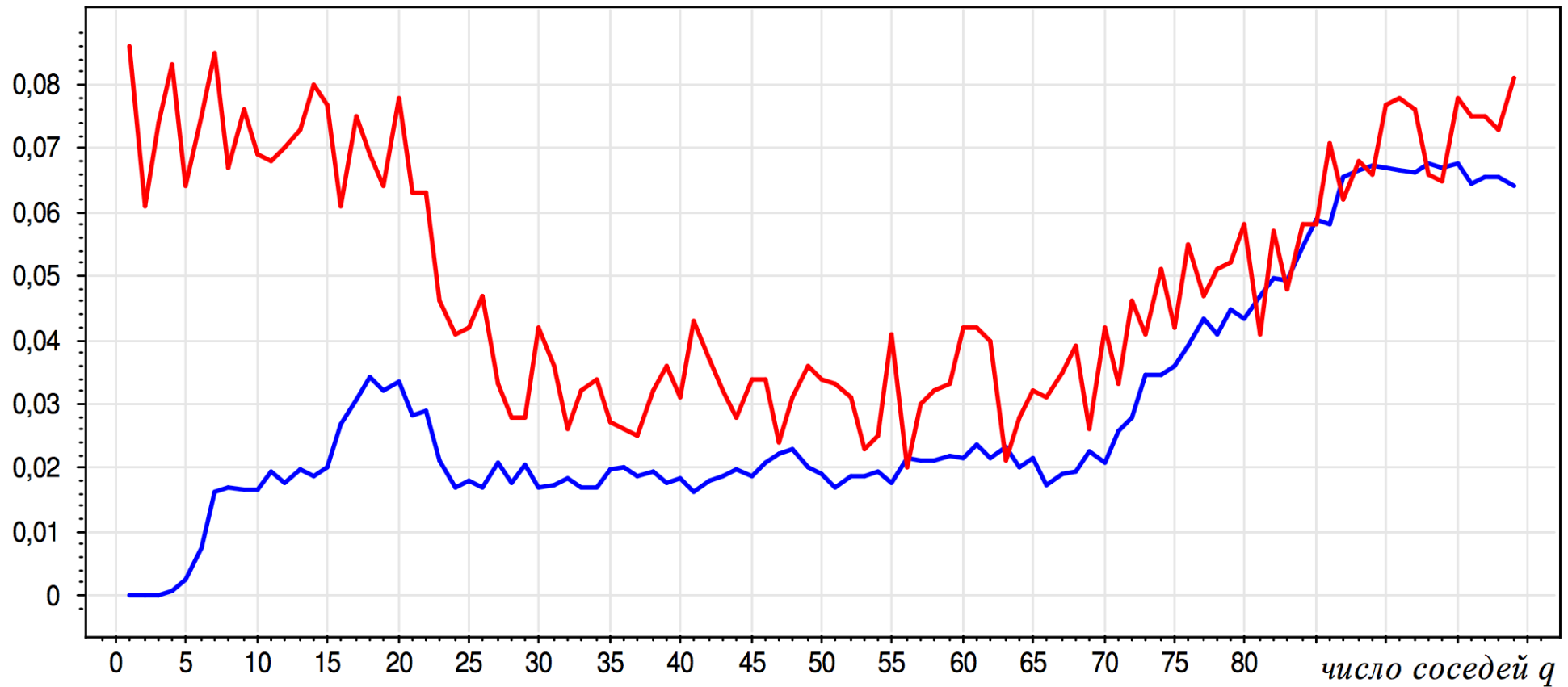
# Отложенная выборка



- Слишком большое обучение — тестовая выборка нерепрезентативна
- Слишком большой тест — модель не сможет обучиться
- Обычно: 70/30, 80/20

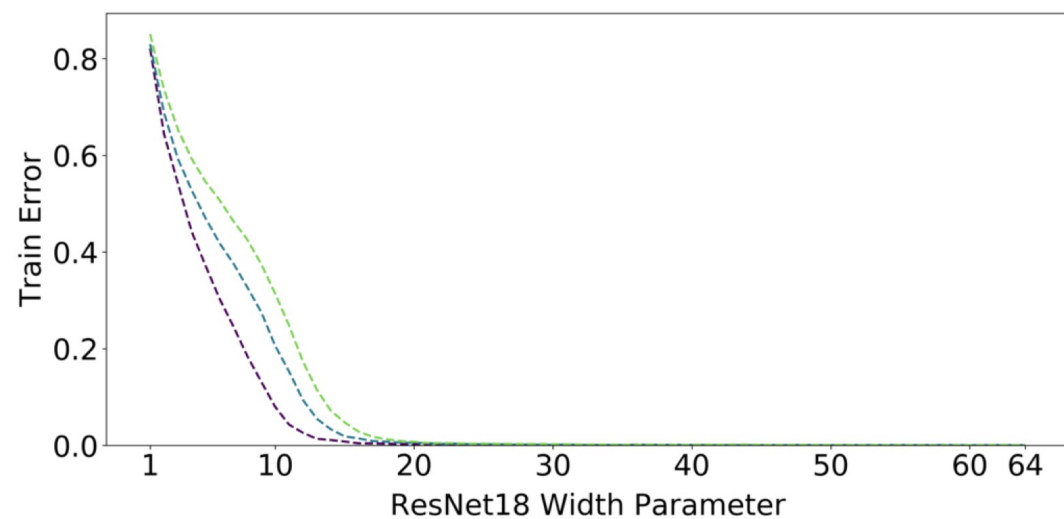
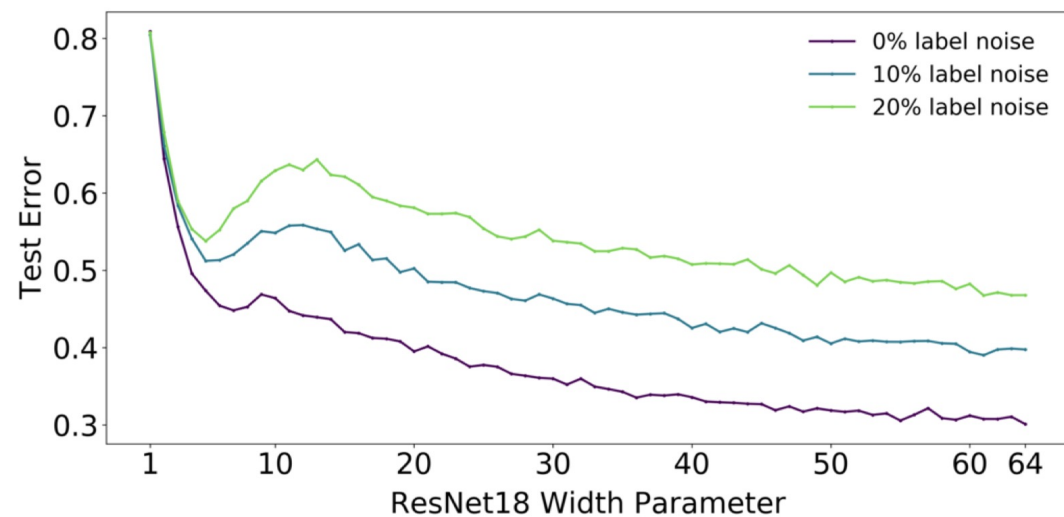
# Подбор числа соседей

*частота ошибок*

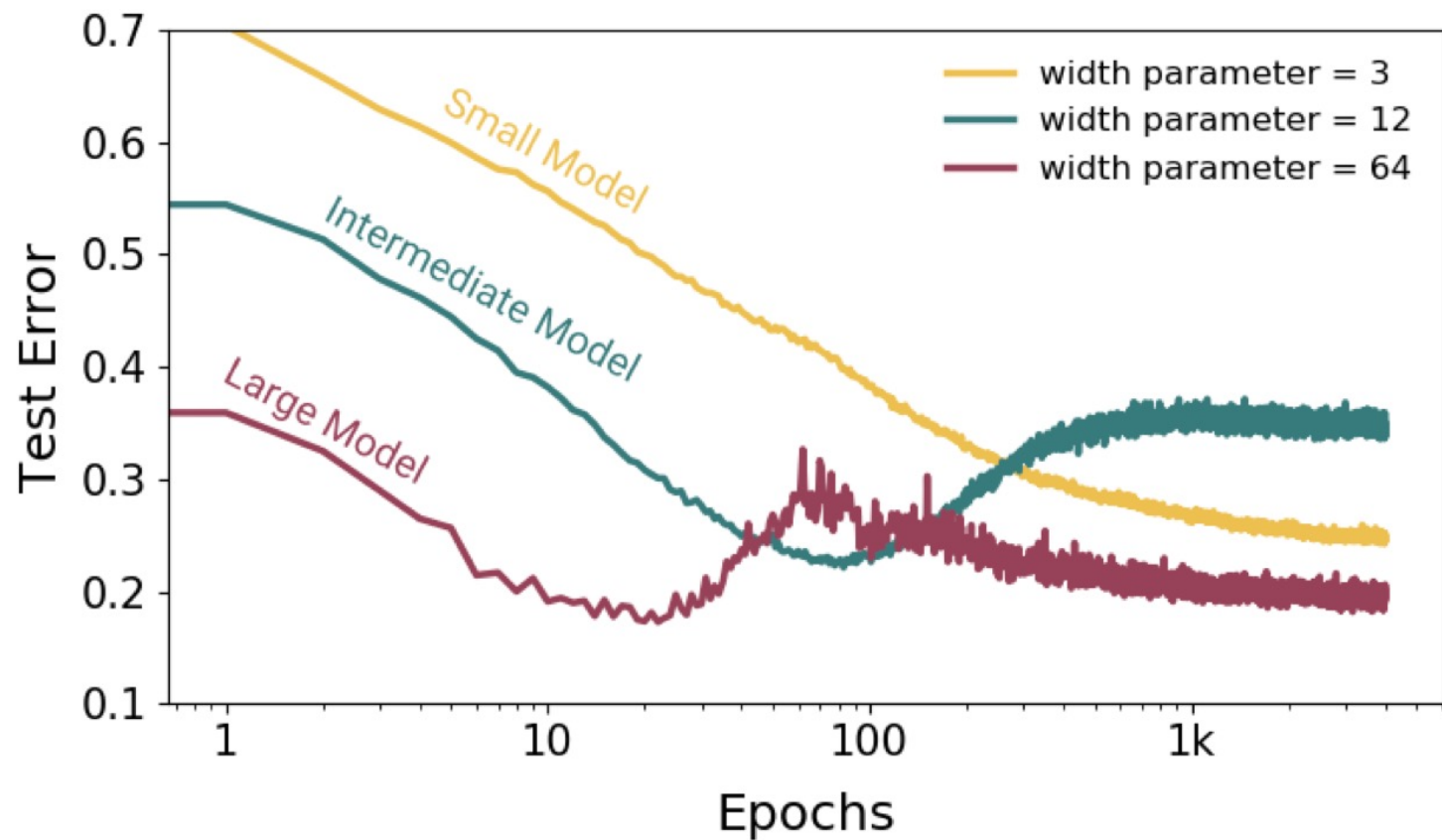


<http://www.machinelearning.ru/wiki/index.php?title=MO>

# Что если увеличивать число параметров?



# Что если продолжать градиентный спуск?

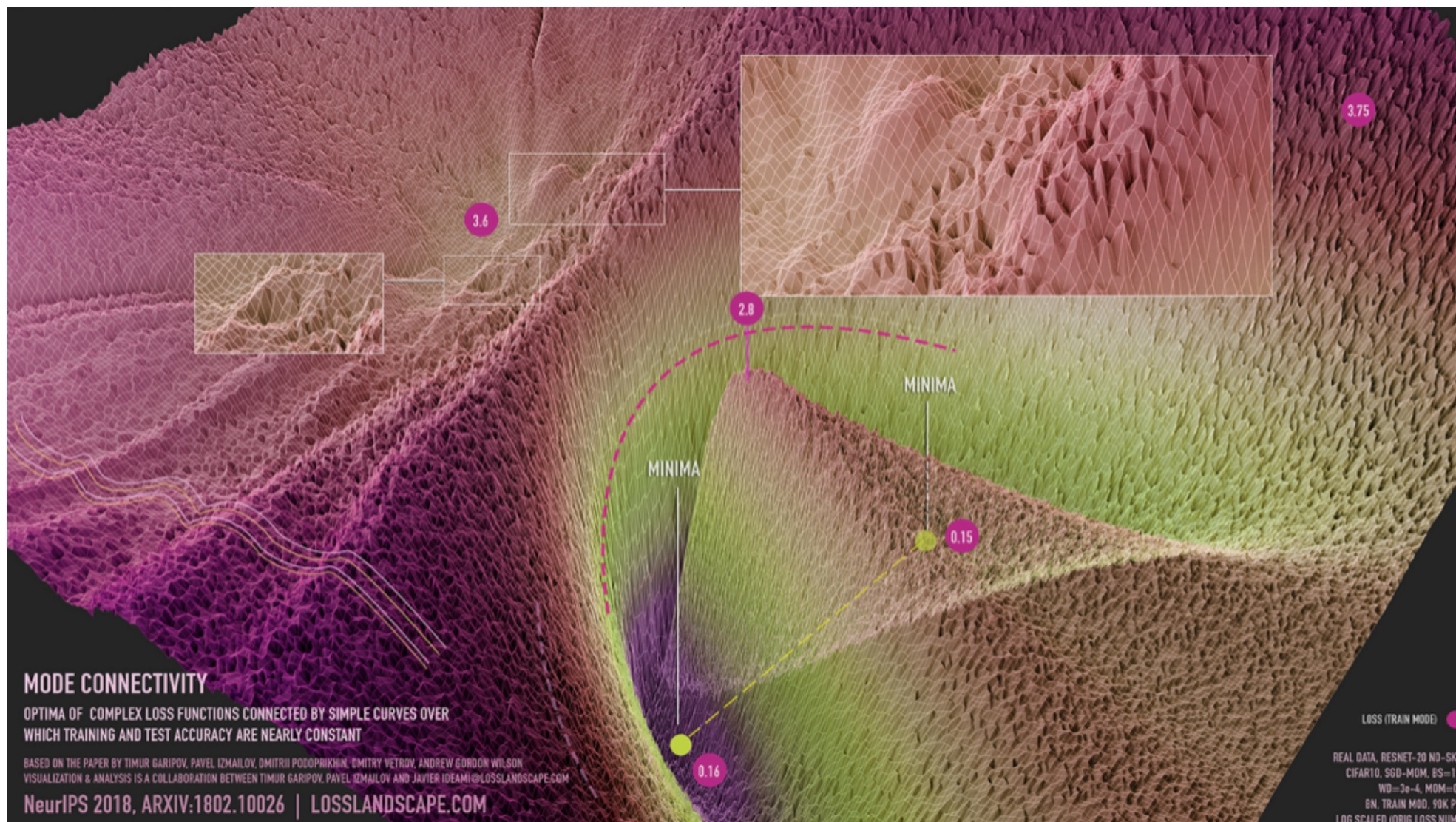


# Эффект двойного спуска

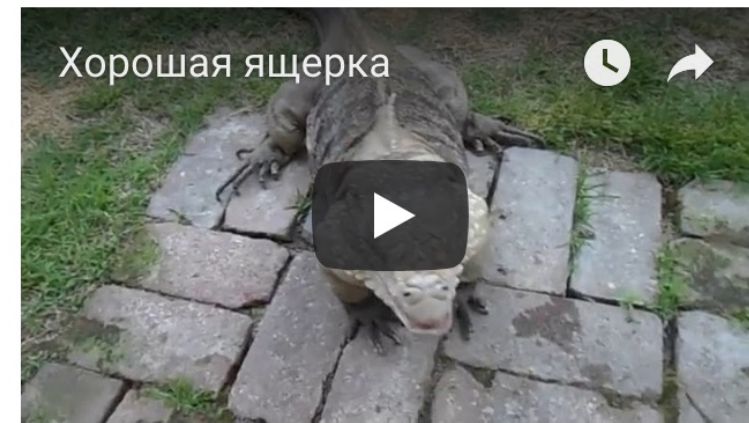
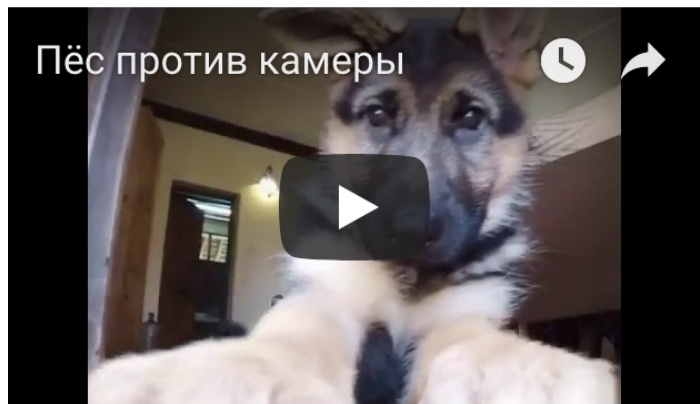
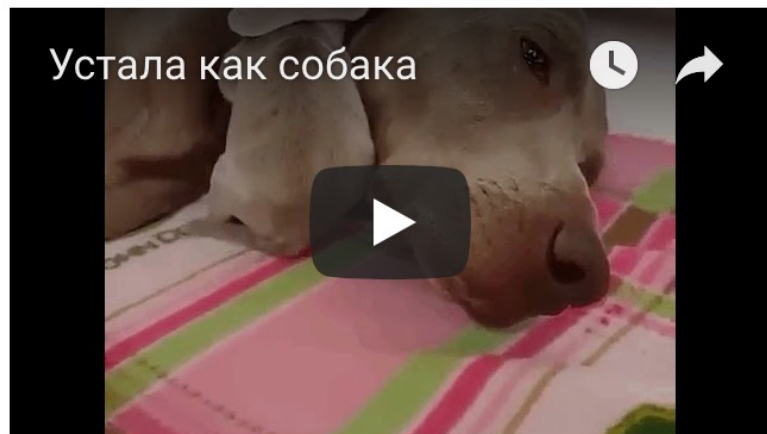
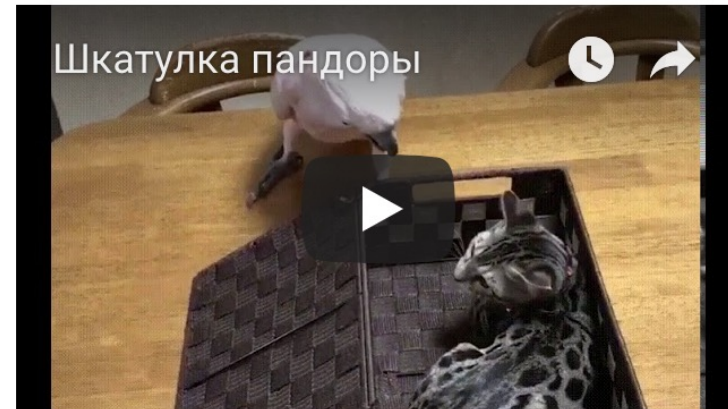
Три режима работы моделей:

- Недопараметризованный режим
- Стандартный режим
- Перепараметризованный режим

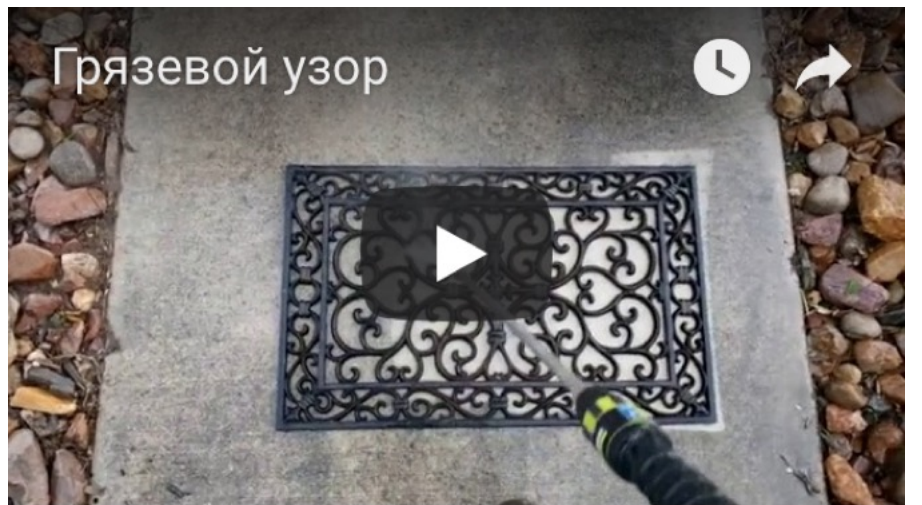
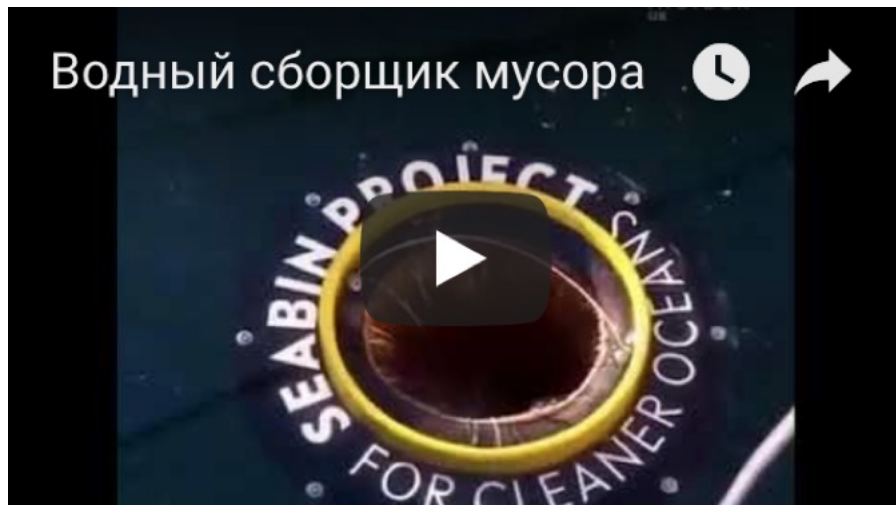
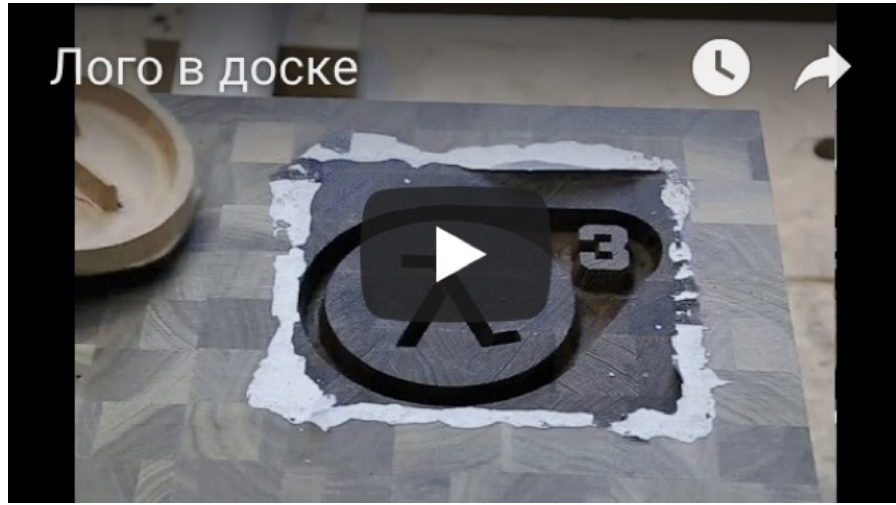
# Рельеф функции потерь

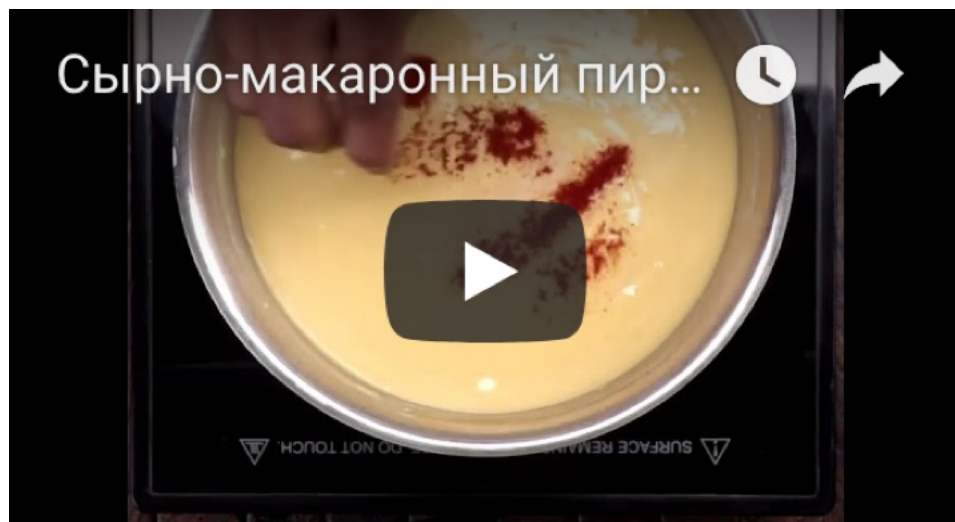
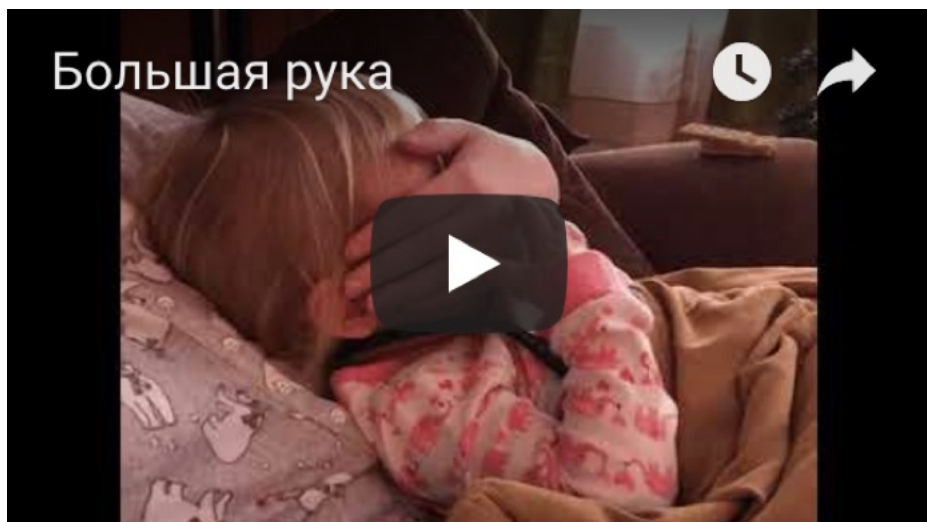
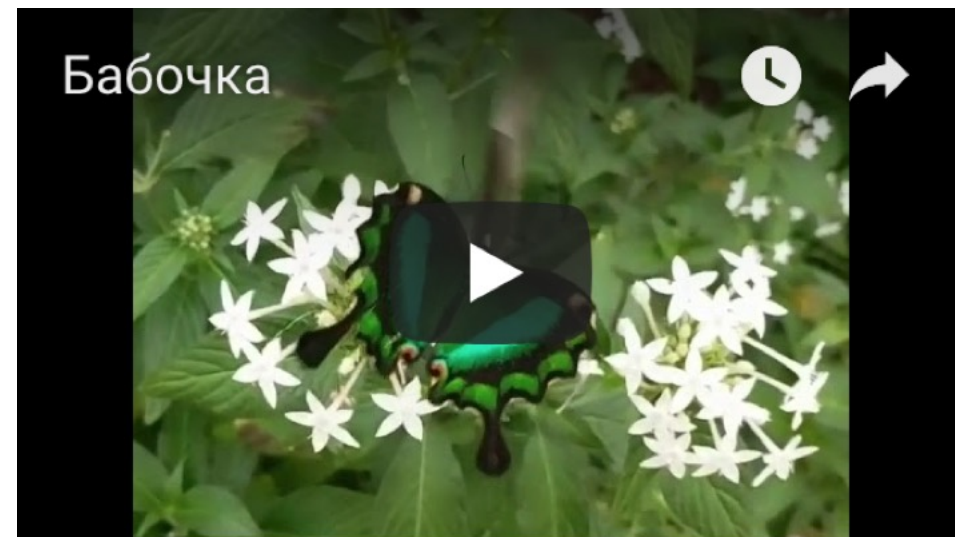
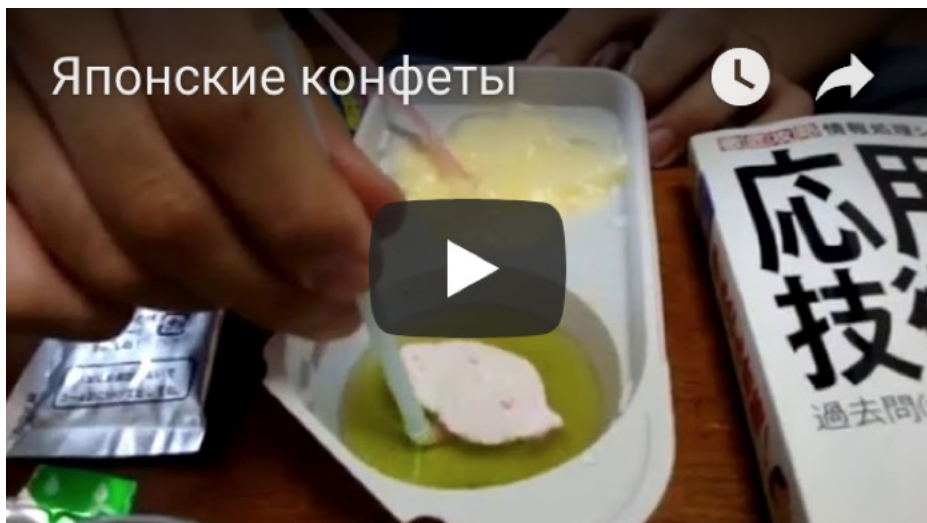


# Рекомендательные системы









# User-based collaborative filtering

Видео

Пользователи

1	1	0		1	
0	1	1			1
			1	1	0
	1	1		0	
	1				1

# User-based collaborative filtering

Видео

Пользователи

1	1	0		1	
0	1	1			1
			1	1	0
	1	1		0	
	1				1

# User-based collaborative filtering

Видео

	1	1	0		1	
	0	1	1			1
Пользователи				1	1	0
		1	1		0	
	1					1

Похожие пользователи

# User-based collaborative filtering

Видео

1	1	0		1	
0	1	1			1
			1	1	0
	1	1		0	
	1				1

Пользователи

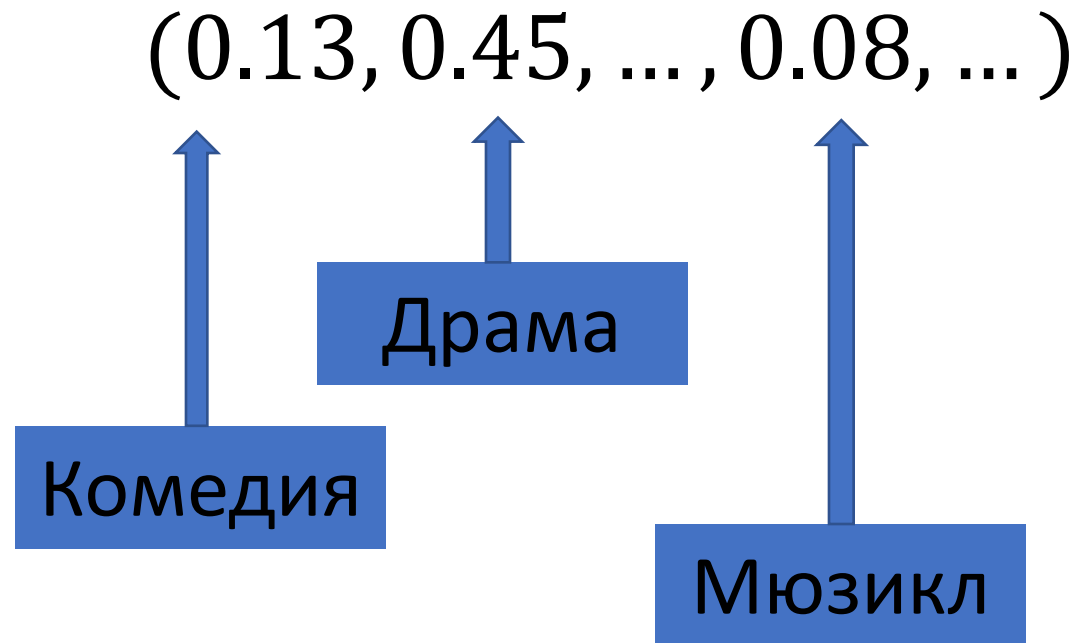
Похожие пользователи

# User-based collaborative filtering

- Решение чисто эвристическое
- Можно ли лучше?

# Векторы интересов

- Для пользователя — насколько он интересуется каждым жанром
- Для фильма — насколько он относится к каждому жанру





# Рейтинг

- Предположение: заинтересованность определяется как скалярное произведение векторов пользователя и фильма

$$(0.1, 0.5, 0.01, 0.92) \times (0, 0, 0.1, 0.95) = 0.875$$

$$(0.1, 0.5, 0.01, 0.92) \times (0.9, 0, 0, 0.1) = 0.182$$

Пользователь

Фильм

# Модели со скрытыми переменными

- Обучим вектор  $p_u$  для каждого пользователя  $u$
- Обучим вектор  $q_i$  для каждого товара  $i$
- Оценка приближается их скалярным произведением:

$$r_{ui} \approx \langle p_u, q_i \rangle$$

- Находим векторы только по известным оценкам
- После этого можем предсказать оценку для любой пары «пользователь-товар»

# Модели со скрытыми переменными

- Оптимизационная задача:

$$\sum_{(u,i) \in R} (r_{ui} - \bar{r}_u - \bar{r}_i - \langle p_u, q_i \rangle)^2 \rightarrow \min_{P, Q}$$

- Решение: градиентный спуск, Alternating Least Squares (ALS) и другие методы

# Модели со скрытыми переменными

2	5	
5		4
	1	
	2	5

# Модели со скрытыми переменными

	(0.9, 0.05)	(0.02, 1.1)	(1.05, 0.01)
(2.1, 5)	2	5	
(4.6, 0)	5		4
(0, 1)		1	
(4.9, 0.9)		1	5

# Коллаборативная фильтрация

- Десятки миллионов пользователей и сотни тысяч видео
- Нужно уметь строить векторы очень быстро и распределённо

# Резюме

- Машинное обучение — подбор алгоритма под данные
- Много интересных подходов для текстов, картинок, звука
- В основе лежит много математики и алгоритмов

tg: @esokolov  
[esokolov@hse.ru](mailto:esokolov@hse.ru)