

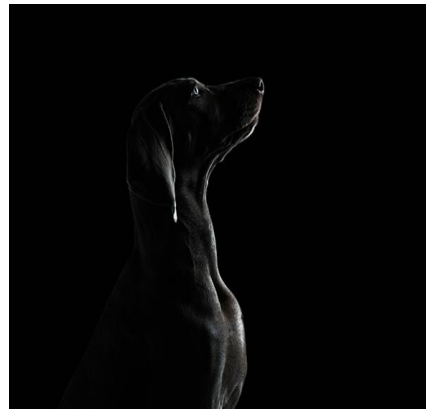
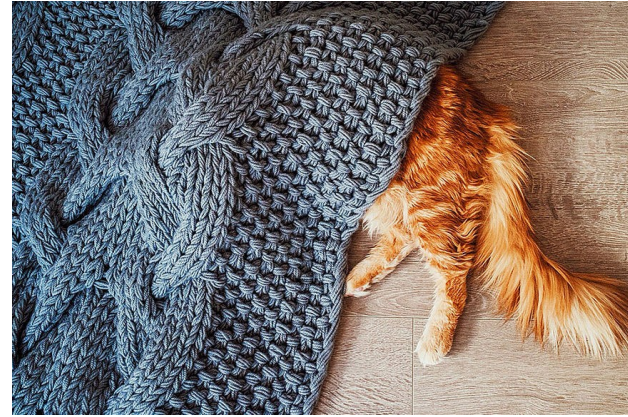
2

Элементы

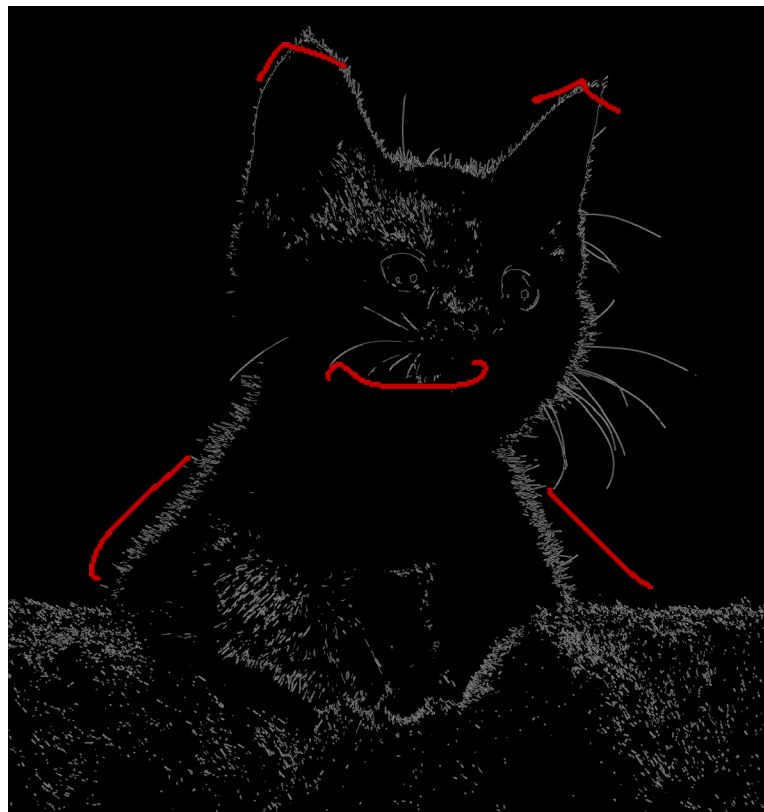
Машинного Обучения



# Кошечки или собачки



# Алгоритм?



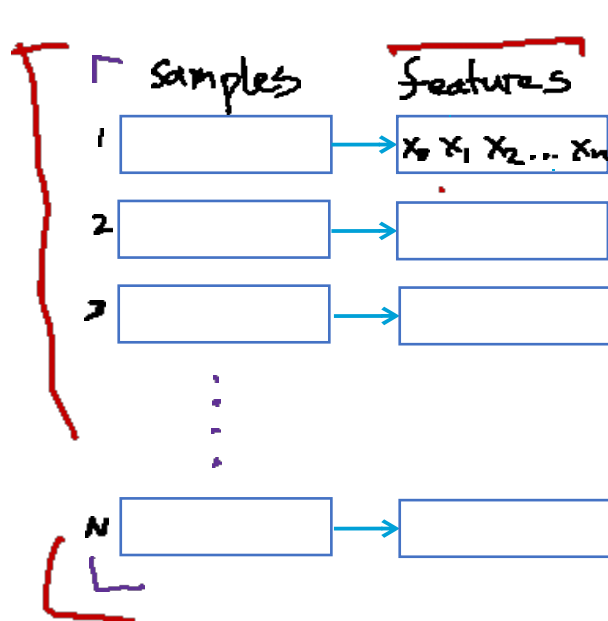
# Машинное обучение Machine Learning

**Machine learning** is a field of [computer science](#) that gives [computers](#) the ability to learn without being explicitly programmed.<sup>[1]</sup>

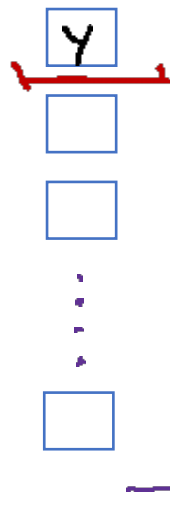
[Wikipedia](#)

## Supervised Learning

data



Labels



ML  
algorithm

- Linear regression
- decision tree
- SVM
- ....

prediction

Model

Контроль  
Соответствия



# Street View House Numbers

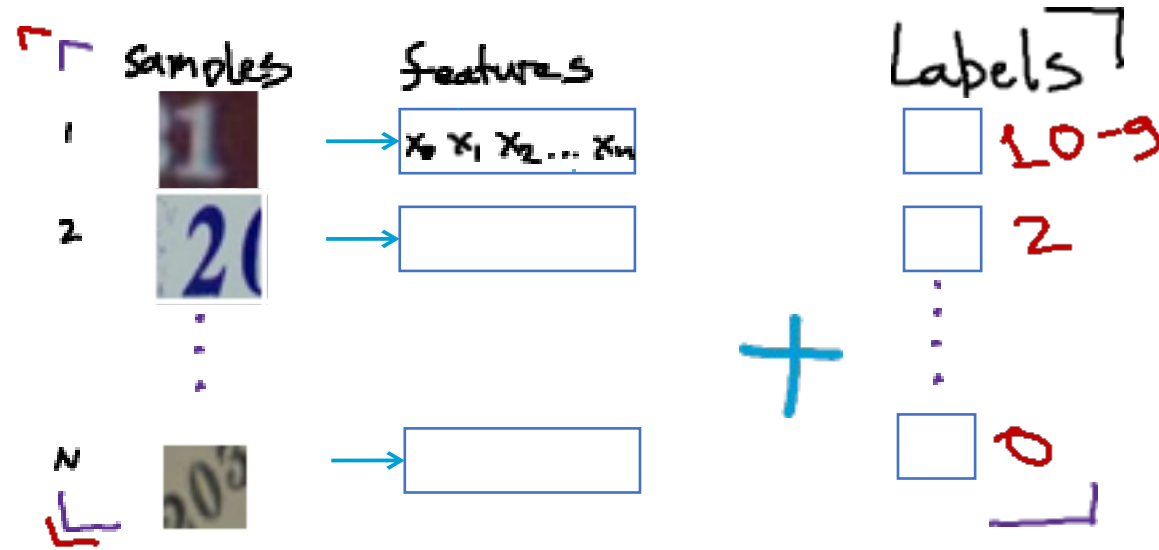


32x32 pixels  
10 classes  
~70000 train  
~25000 test

training data

predict

32-32-3  
32x32 = 1024  
[.....]



ML  
algorithm

test



Model

# Метод ближайших соседей

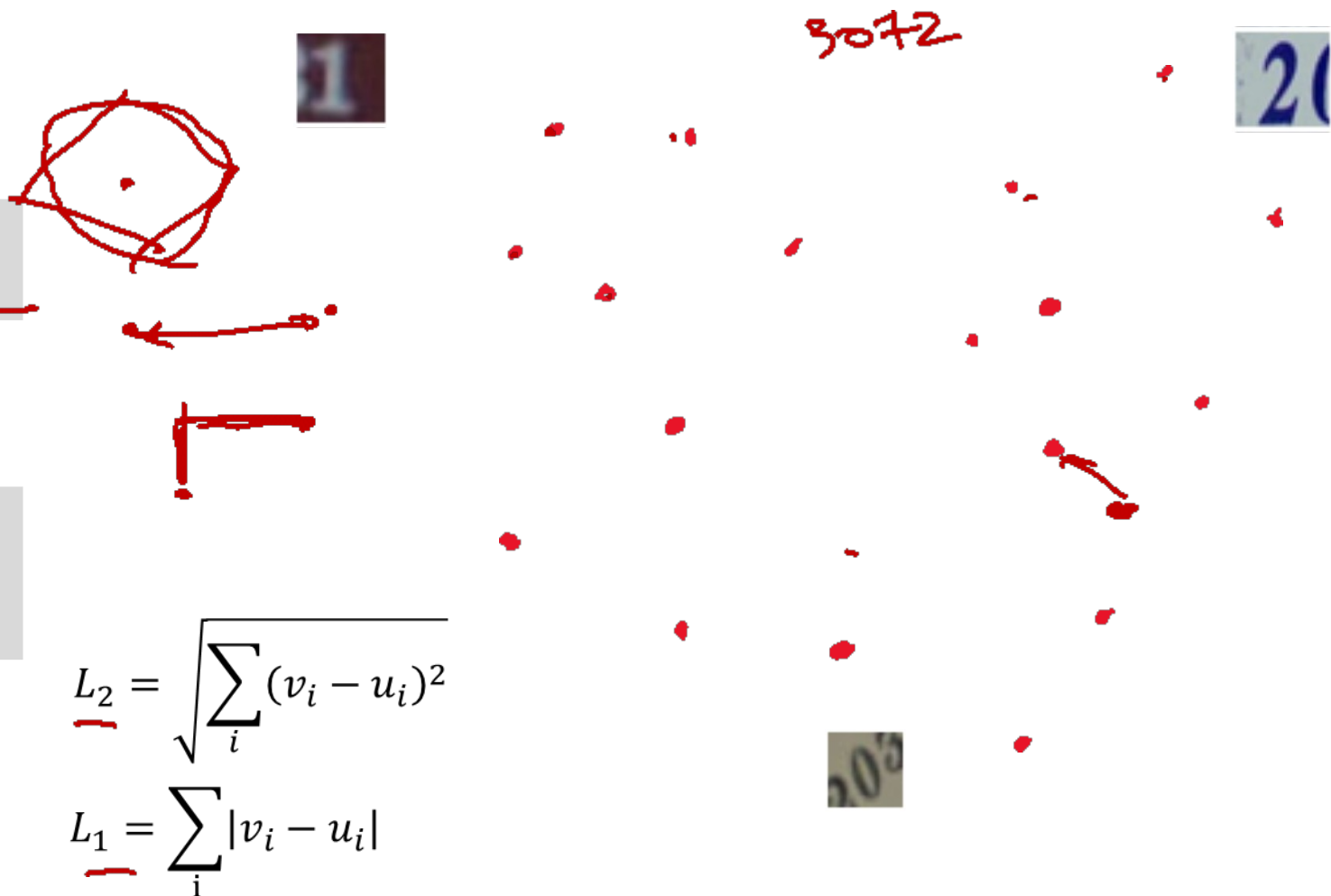
## Nearest neighbor

Train:

просто все запомнить

Predict:

найти ближайший и  
выдать его класс



$$\underline{L_2} = \sqrt{\sum_i (v_i - u_i)^2}$$

$$\underline{L_1} = \sum_i |v_i - u_i|$$

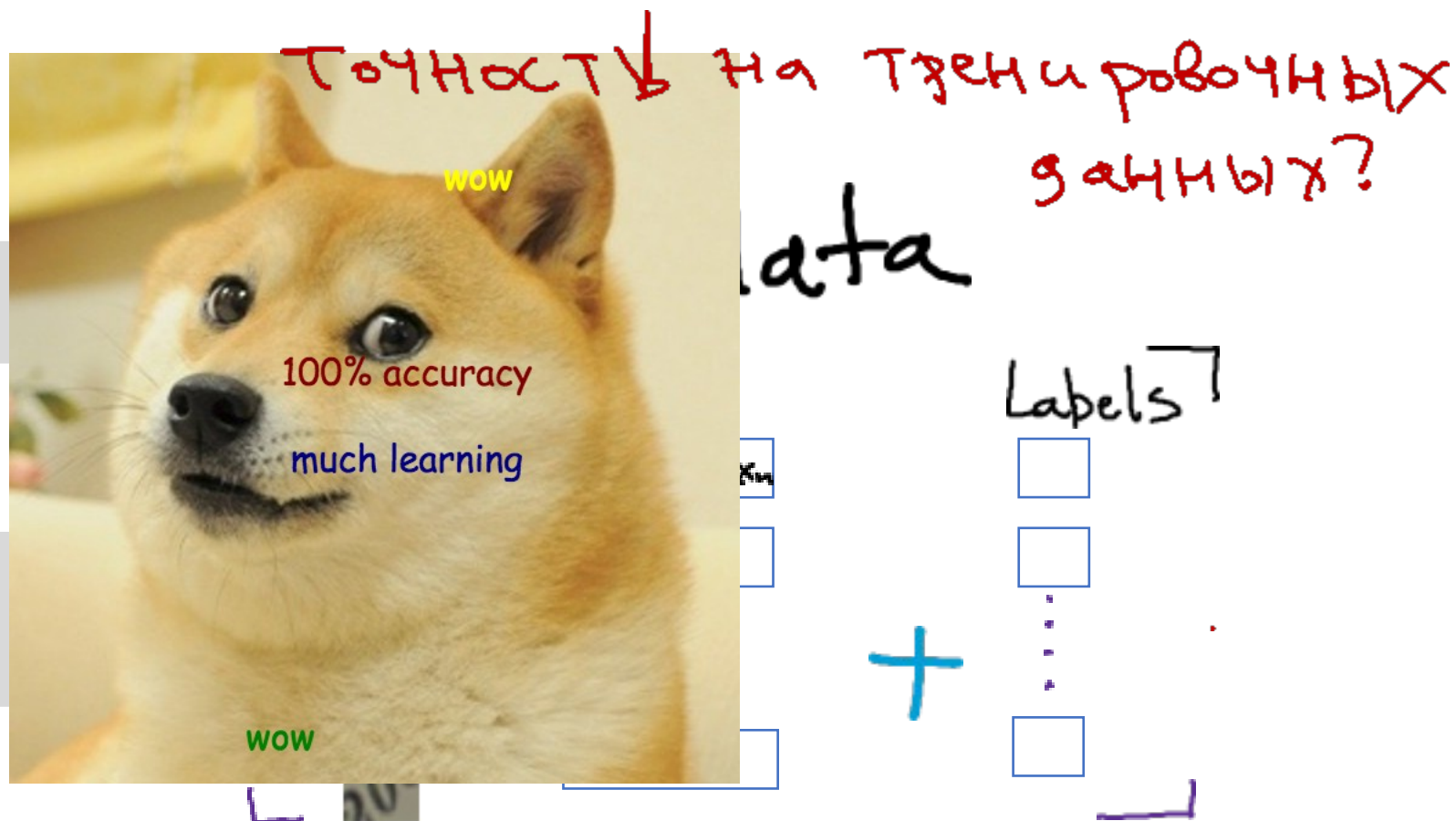
# Точность на тренировочных VS тестовых данных

Train:

просто все запомнить

Predict:

найти ближайший и  
выдать его класс





# Метод k-ближайших соседей

## K-nearest neighbors

1

20

train / test

80 / 20

70 / 30

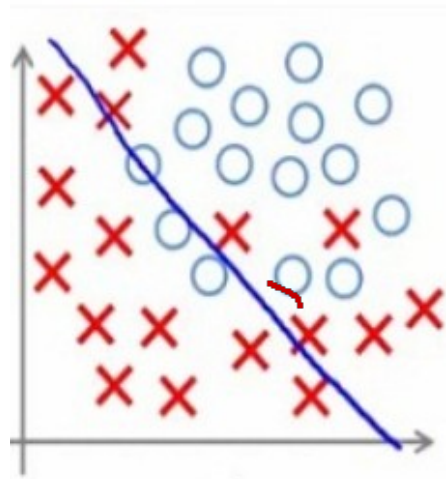


K — нужен  
параметр

Как выбрать K?

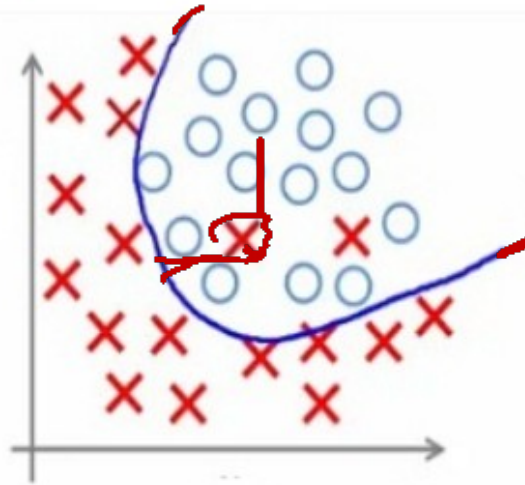
# Переобучение и недообучение

## Overfitting vs underfitting

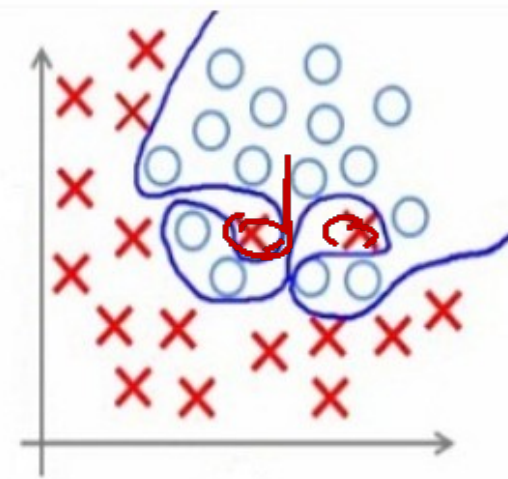


Under-fitting

(too simple to explain the variance)



Appropriate-fitting

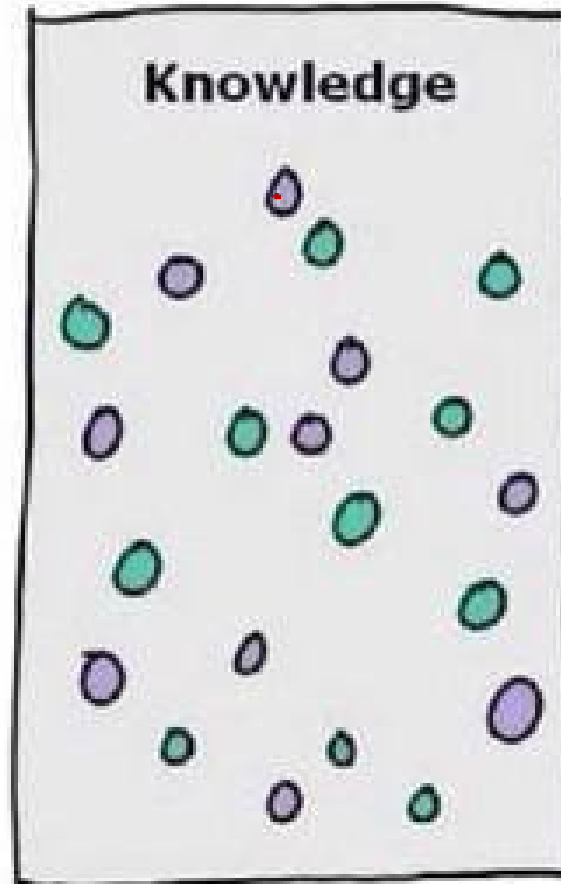


Over-fitting

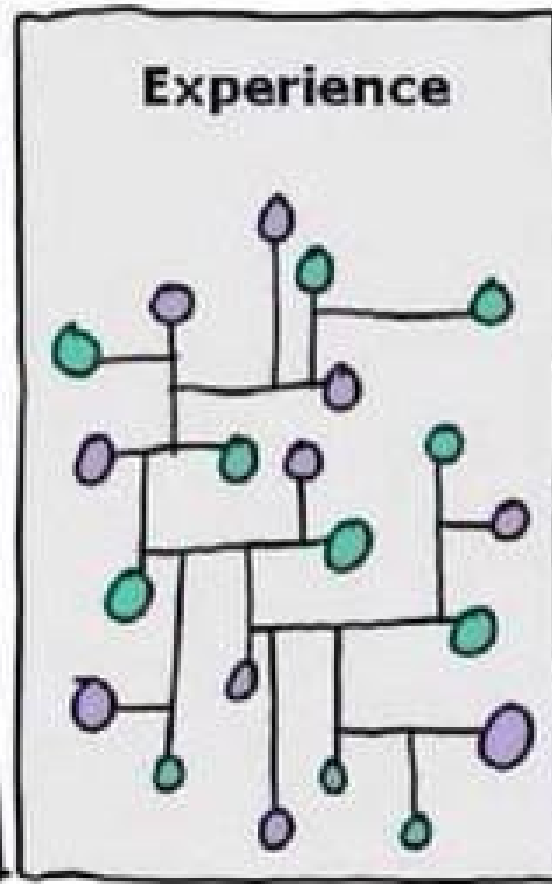
(forcefitting -- too good to be true)

$x = y$

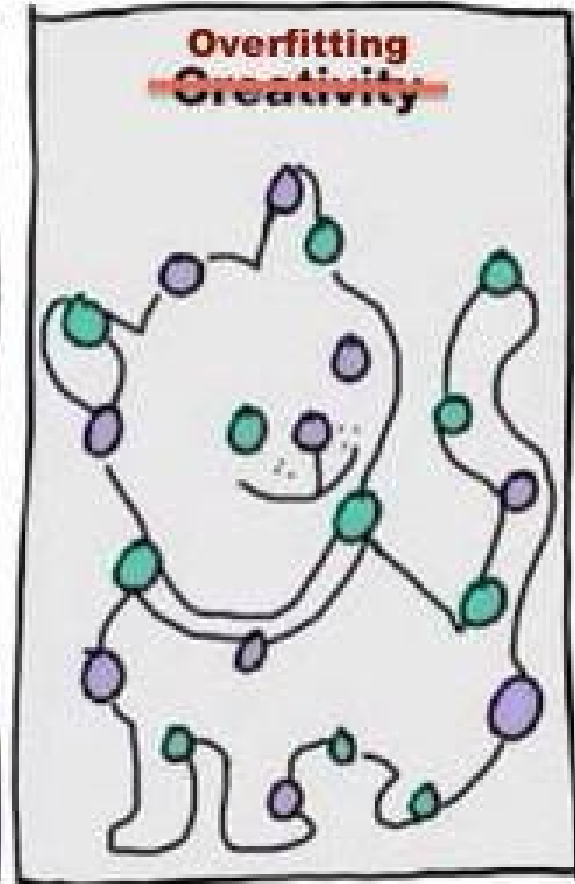
**Knowledge**



**Experience**



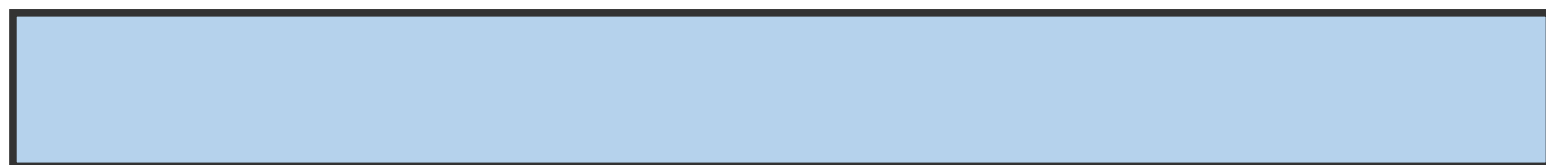
**Overfitting  
~~Creativity~~**





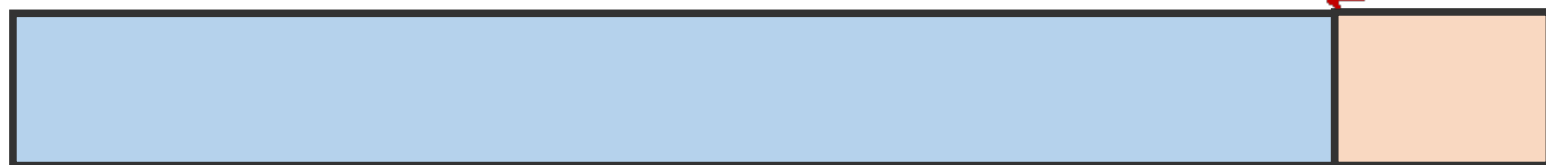
$k=1$

train



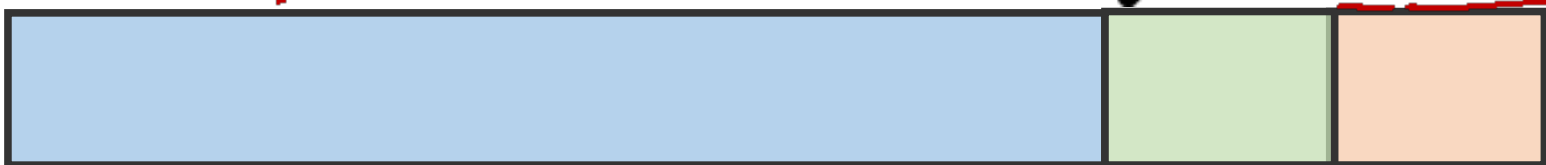
train

test



train

val test



$\frac{1}{k}$

# Кросс-валидация

## Cross-validation

train

Val

test

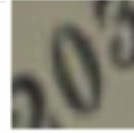


# Как сравнивать?

Бинарная классификация  
Binary classification



vs



1000

990

10'

99% - 99.9%



пропуск

ложное срабатывание

Точность

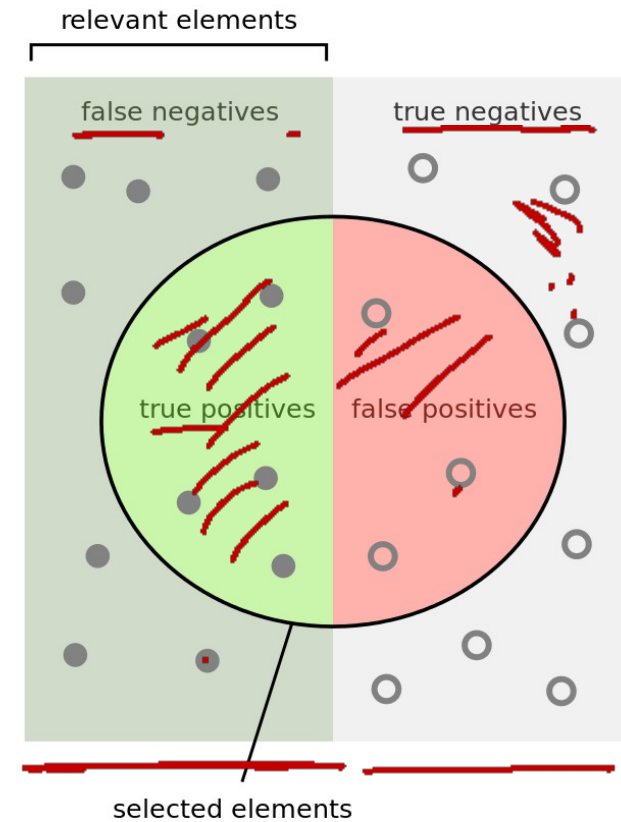
$$\text{Accuracy} \equiv \frac{\text{correct}}{\text{total}} \quad \underline{70\% \quad 70\%}$$



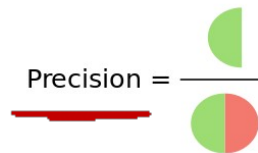
$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$F1 = \frac{2}{\frac{1}{\text{precision}} + \frac{1}{\text{recall}}} = \frac{2 * (\text{precision} * \text{recall})}{\text{precision} + \text{recall}} .8$$

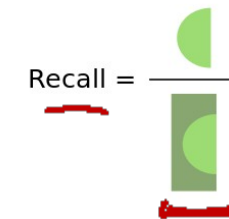


How many selected items are relevant?



Precision =

How many relevant items are selected?



Recall =

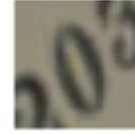
.5

# Как сравнивать?

Бинарная классификация  
Binary classification



vs



Точность

$$\underline{Accuracy} \equiv \frac{correct}{total}$$

$$\underline{Precision} = \frac{TP}{TP + FP}$$

$$\underline{Recall} = \frac{TP}{TP + FN}$$

$$\underline{F_1} = 2 * \frac{precision * recall}{precision + recall}$$

# Как сравнивать?

Многоклассовая классификация  
Multi-class classification

7



VS



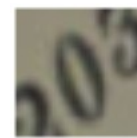
VS

...

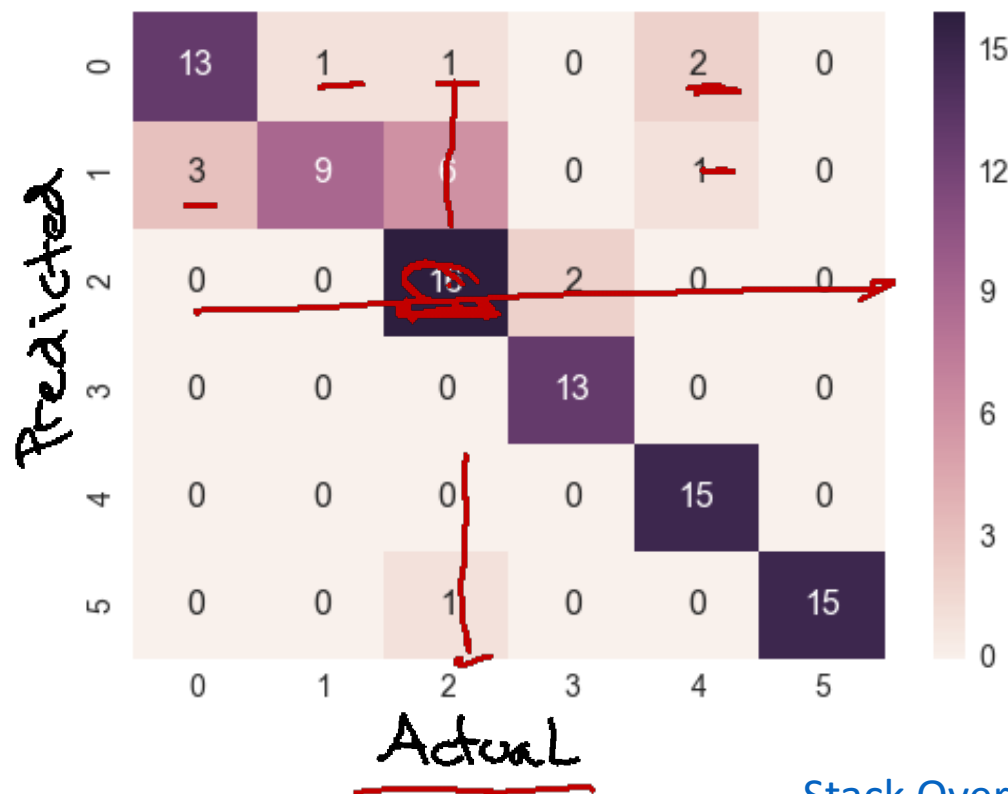
VS



VS



confusion matrix



Точность

$$\text{Accuracy} \equiv \frac{\text{correct}}{\text{total}}$$

$$\text{Precision}_c = \frac{A_{c,c}}{\sum_{i=1}^n A_{c,i}}$$

$$\text{Recall}_c = \frac{A_{c,c}}{\sum_{i=1}^n A_{i,c}}$$

$$\text{Precision} = \frac{\sum_{c=1}^n P_c}{n}$$

$$\text{Recall} = \frac{\sum_{c=1}^n R_c}{n}$$



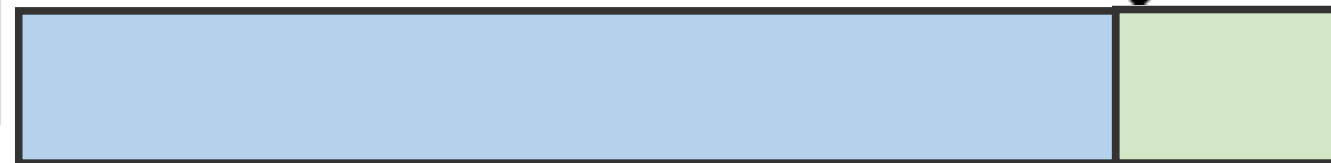


# Machine Learning Flow

train

val

test



Кстати вот он про это подробнее:

[Nuts and Bolts of Applying Deep Learning](#)



Ошибка на train

большая

*underfitting*

- Более мощную модель
- Больше ресурсов для тренировки
- Другой подход

маленькая ↓

Ошибка на val

большая

*overfitting*

- Больше данных
- Больше регуляризации
- Другой подход

маленькая ↓

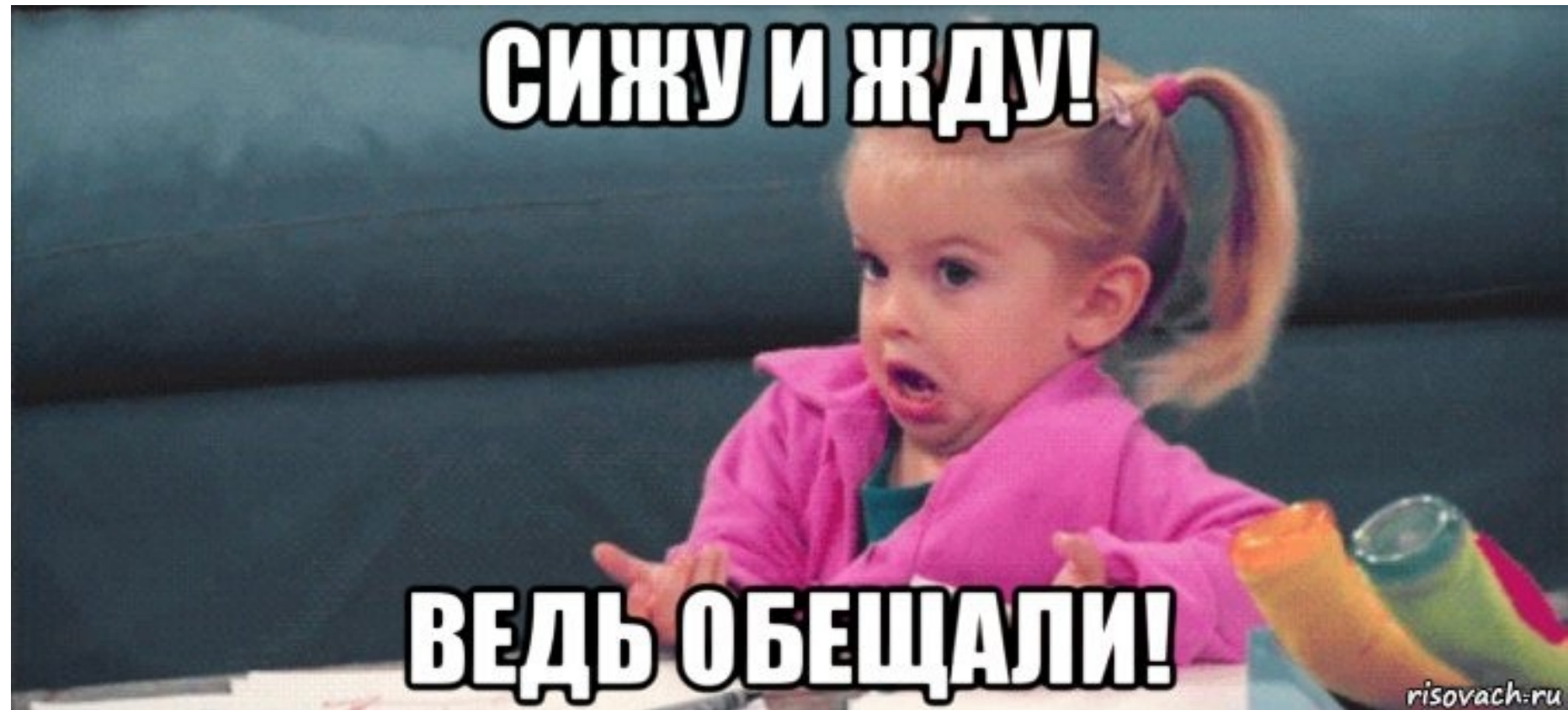
Ошибка на test

большая

- Отличаются train и test
- Больше данных, таких как test

← маленькая

В следующий раз уже будет про  
нейросети





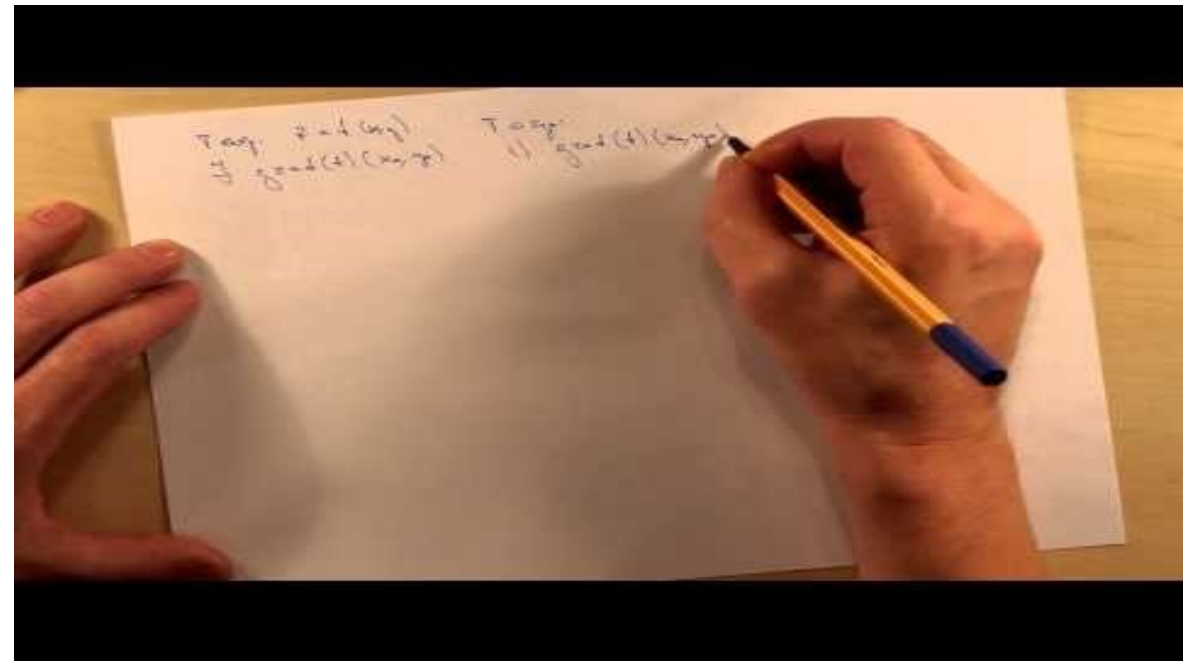
# Домашнее задание!

Повторить производную сложной функции (chain rule)

и Градиент



[Link](#)



[Link](#)