

Лекция 5

Нейронные сети

Бручес Елена

15.07.2025

Основные методы в NLP

- Правила/эвристики;
- Статистические методы;
- Методы машинного обучения

Машинное обучение

Машинное обучение – это область компьютерные наук, которая изучает алгоритмы, которые анализируют данные и делают предсказания, основываясь на этих данных.

Приложения машинного обучения:

- NLP;
- Computer Vision;
- Optical Character Recognition (OCR);
- Business analytics;
- etc

Machine Learning

Machine Learning tasks:

- **Supervised learning**: The computer is presented with example inputs and their desired outputs, given by a "teacher", and the goal is to learn a general rule that **maps** inputs to outputs;
- **Unsupervised learning**: No labels are given to the learning algorithm, leaving it on its own to find structure in its input;
- **Reinforcement learning**: A computer program interacts with a dynamic environment in which it must perform a certain goal. The program is provided feedback in terms of rewards and punishments as it navigates its problem space.

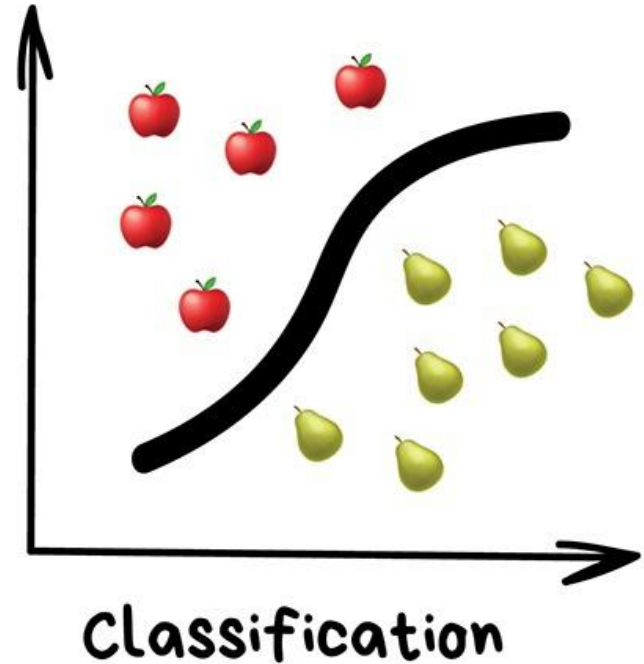
Machine Learning

Machine learning tasks:

- In **classification**, inputs are divided into two or more classes, and the learner must produce a model that assigns unseen inputs to one or more (**multi-label classification**) of these classes;
- In **regression**, also a supervised problem, the outputs are continuous rather than discrete;
- In **clustering**, a set of inputs is to be divided into groups. Unlike in classification, the groups are not known beforehand, making this typically an unsupervised task;
- **Density estimation** finds the **distribution** of inputs in some space;
- **Dimensionality reduction** simplifies inputs by mapping them into a lower-dimensional space.

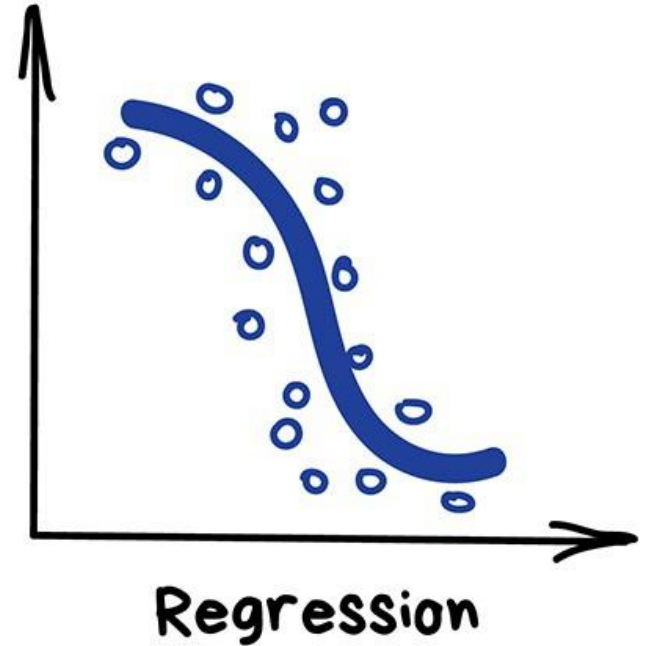
Classification

In **classification**, inputs are divided into two or more classes, and the learner must produce a model that assigns unseen inputs to one or more (**multi-label classification**) of these classes



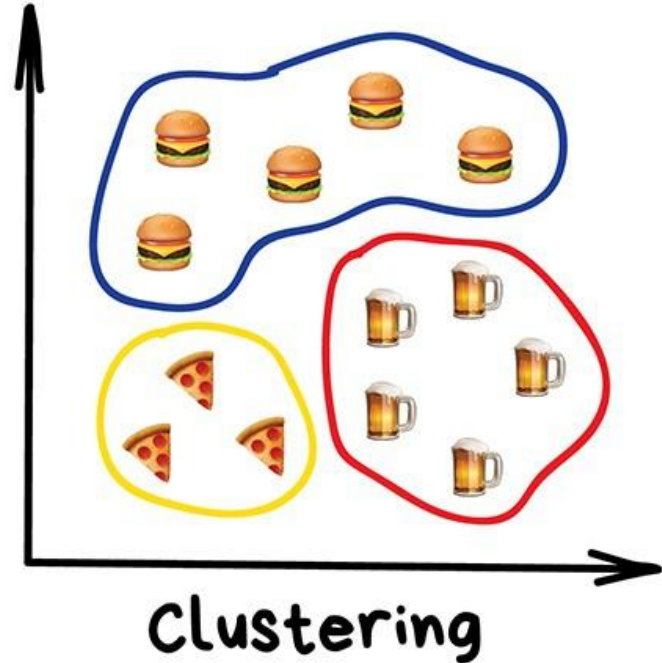
Regression

In **regression**, also a supervised problem, the outputs are continuous rather than discrete



Clustering

In **clustering**, a set of inputs is to be divided into groups. Unlike in classification, the groups are not known beforehand, making this typically an unsupervised task



(Feature, Label) - Sample

(Feature, Label) - Sample

....

(Feature, Label) - Sample



Данные

Терминология ML

Label (метка) - это то, что мы предсказываем. Это может быть цена дома, тип животного на картинке, часть речи слова - в общем, почти всё что угодно.

Feature (признак) - это входная переменная. Самая простая модель может использовать один единственный признак; более сложные модели могут использовать миллионы признаков.

Например, в задаче определения спам/не спам признаками могут быть:

- слова в тексте письма;
- адрес отправителя;
- время дня, в которое было отправлено письмо;
- наличие в письме фразы "Ваш адрес был выбран победителем";
- и др.

Терминология ML

Sample - это конкретный пример данных. Samples бывают двух видов:

1. labeled (размеченные);
2. unlabeled (неразмеченные).

Размеченный sample включает и признак, и label:

labeled examples: $\{features, label\}$: (x, y)

Неразмеченный sample содержит только признак:

unlabeled examples: $\{features, ?\}$: $(x, ?)$

Итак, предположим, что Вы хотите разработать модель, которая будет предсказывать, является ли письмо спамом или нет. В качестве данных у нас есть набор писем, которые были помечены пользователями как "спам" или "не спам". Какие утверждения верны?

1. Слова в теме письма будут хорошими labels.
2. Вы будете использовать неразмеченные данные, чтобы обучить модель.
3. Те письма, которые не помечены "спам" или "не спам", - неразмеченные примеры.
4. Не все labels, которыми помечены письма, могут быть надёжными.

Итак, предположим, что Вы хотите разработать модель, которая будет предсказывать, является ли письмо спамом или нет. В качестве данных у нас есть набор писем, которые были помечены пользователями как "спам" или "не спам". Какие утверждения верны?

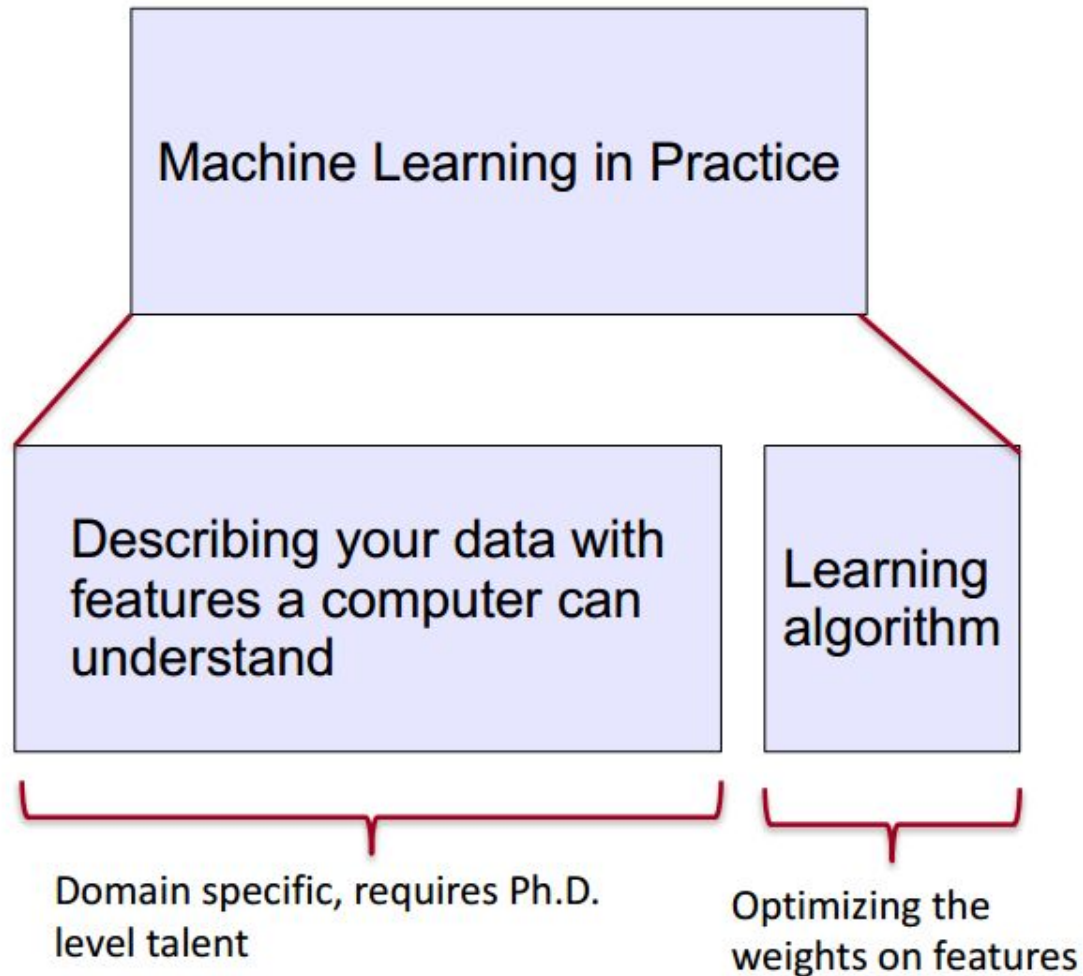
1. Слова в теме письма будут хорошими labels.
2. Вы будете использовать неразмеченные данные, чтобы обучить модель.
3. Те письма, которые не помечены "спам" или "не спам", - неразмеченные примеры.
4. **Не все labels, которыми помечены письма, могут быть надёжными.**

Теперь представим, что обувной он-лайн магазин хочет создать модель, которая будет рекомендовать пользователям обувь. То есть, модель будет рекомендовать определённые пары обуви Кейт и совершенно другие - Джону. Какие утверждения верны?

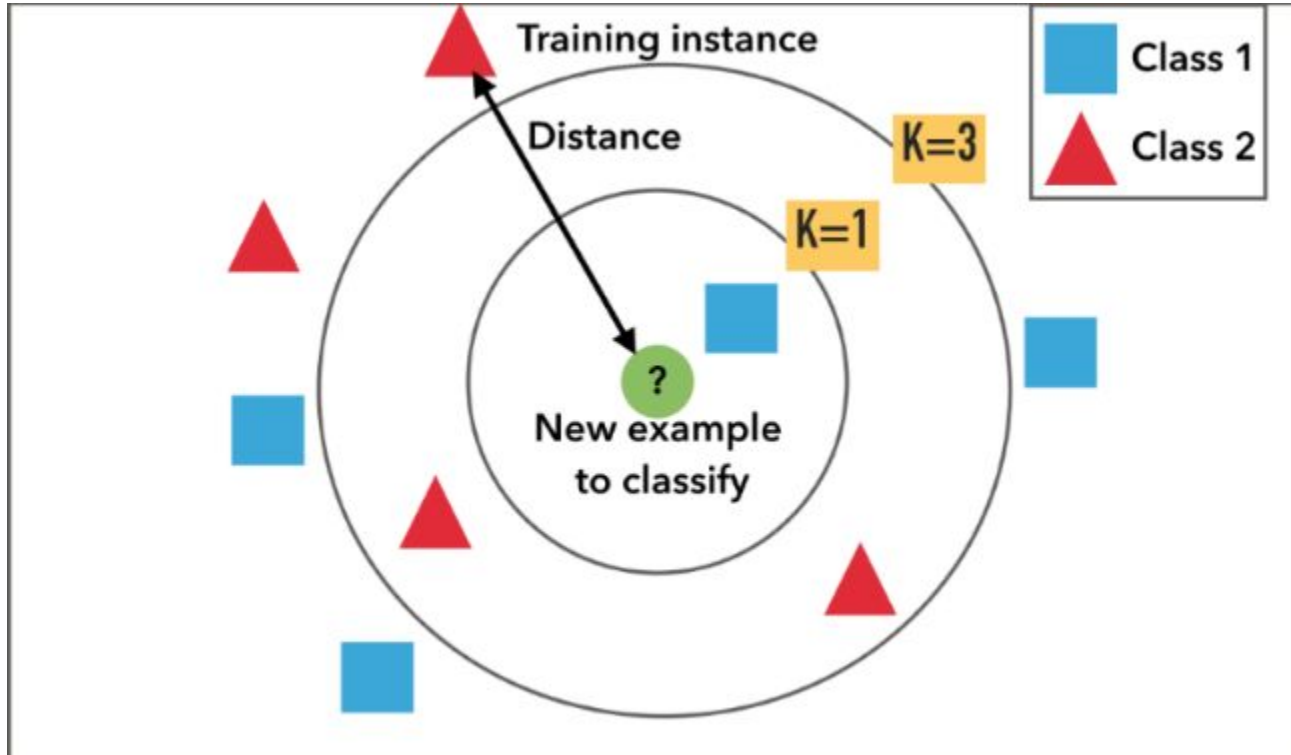
1. "Пользователь кликнул на описание обуви" - хороший label.
2. "Размер обуви" - полезный признак.
3. "Красота обуви" - полезный признак.
4. "Обувь, которую обожает пользователь" - хороший label.

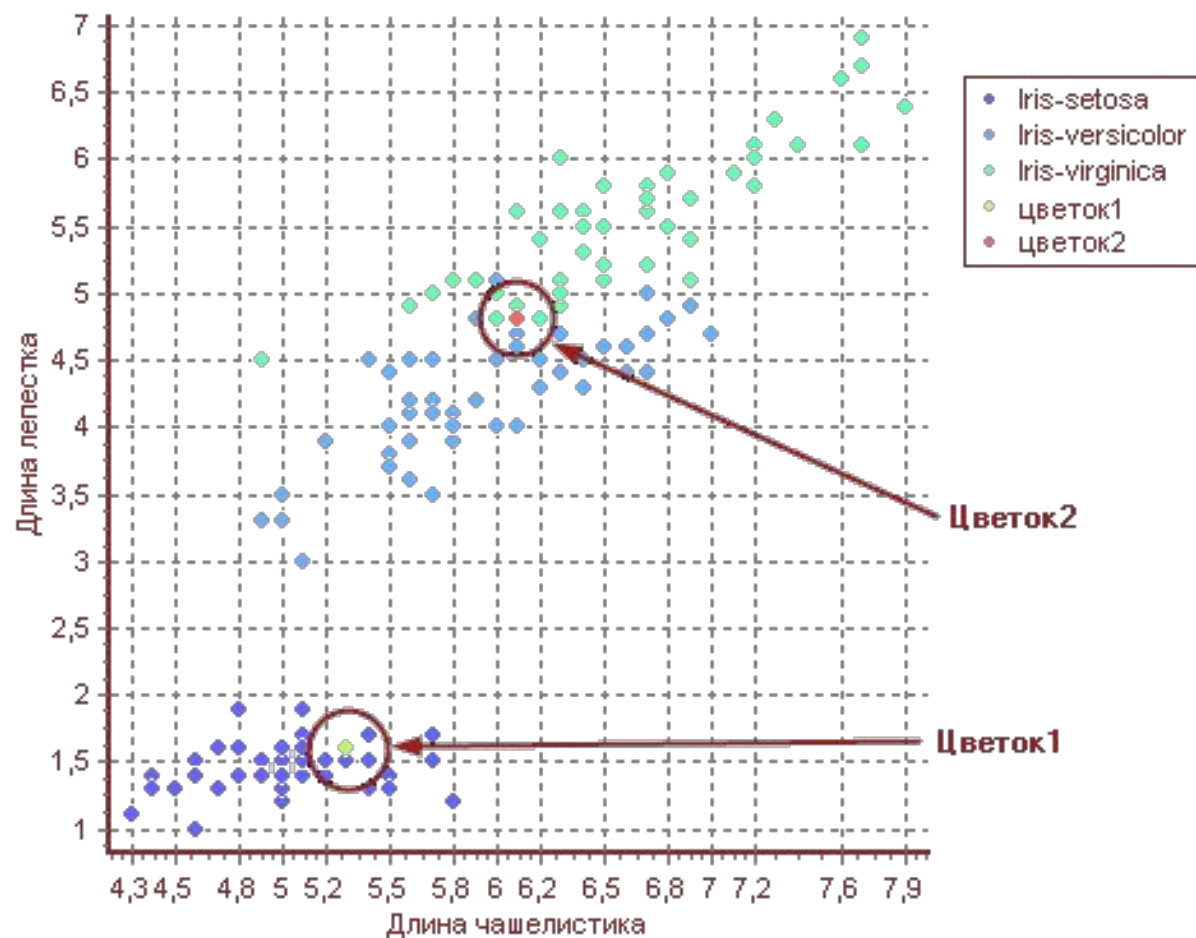
Теперь представим, что обувной он-лайн магазин хочет создать модель, которая будет рекомендовать пользователям обувь. То есть, модель будет рекомендовать определённые пары обуви Кейт и совершенно другие - Джону. Какие утверждения верны?

1. **"Пользователь кликнул на описание обуви" - хороший label.**
2. **"Размер обуви" - полезный признак.**
3. **"Красота обуви" - полезный признак.**
4. **"Обувь, которую обожает пользователь" - хороший label.**

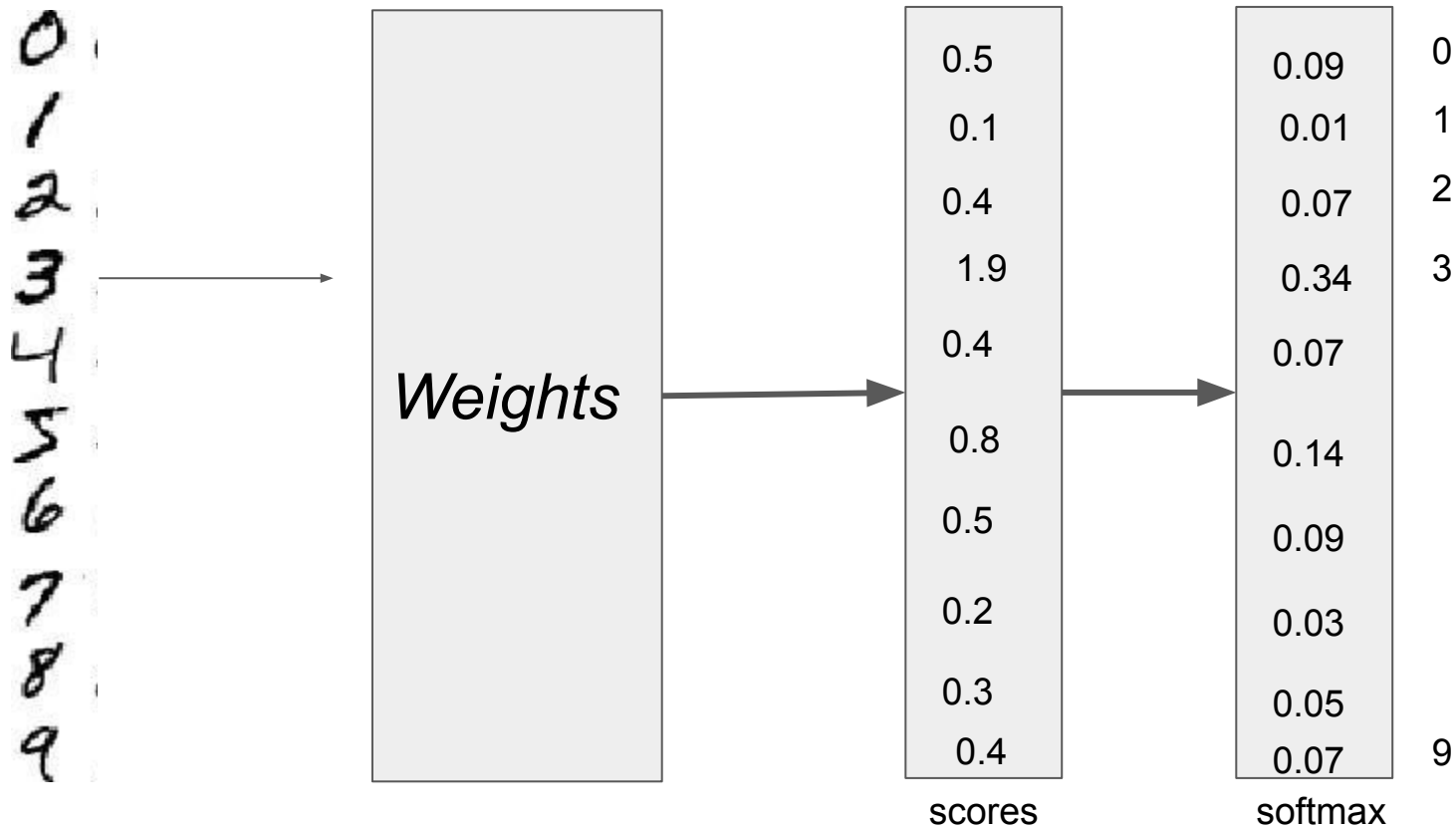


K-Nearest Neighbours

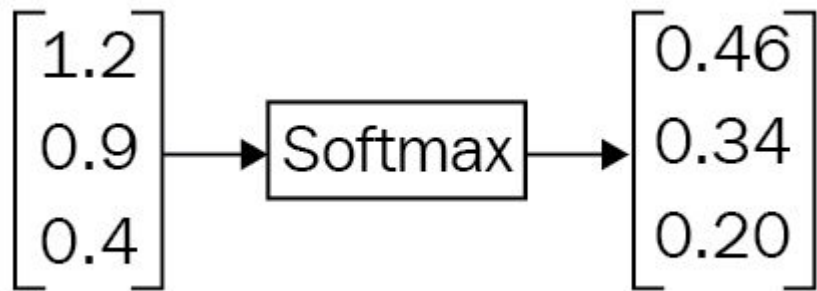




Линейный классификатор



$$\text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$



Обучение модели

Input	Actual output	Desired output
0	0	0
1	3	2
2	6	4
3	9	6
4	12	8

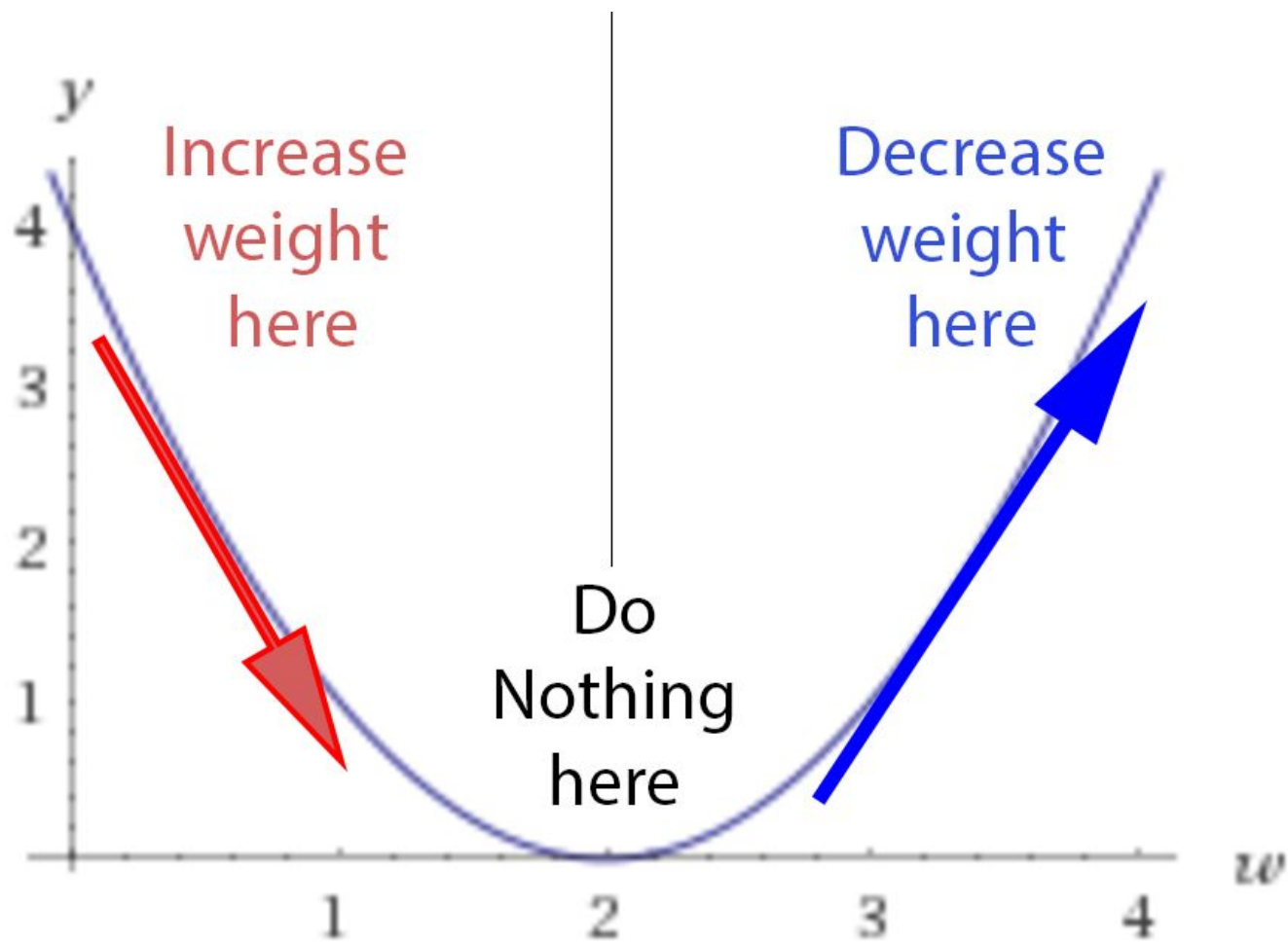
Обучение модели

Input	actual	Desired	Absolute Error	Square Error
0	0	0	0	0
1	3	2	1	1
2	6	4	2	4
3	9	6	3	9
4	12	8	4	16
Total:	-	-	10	30

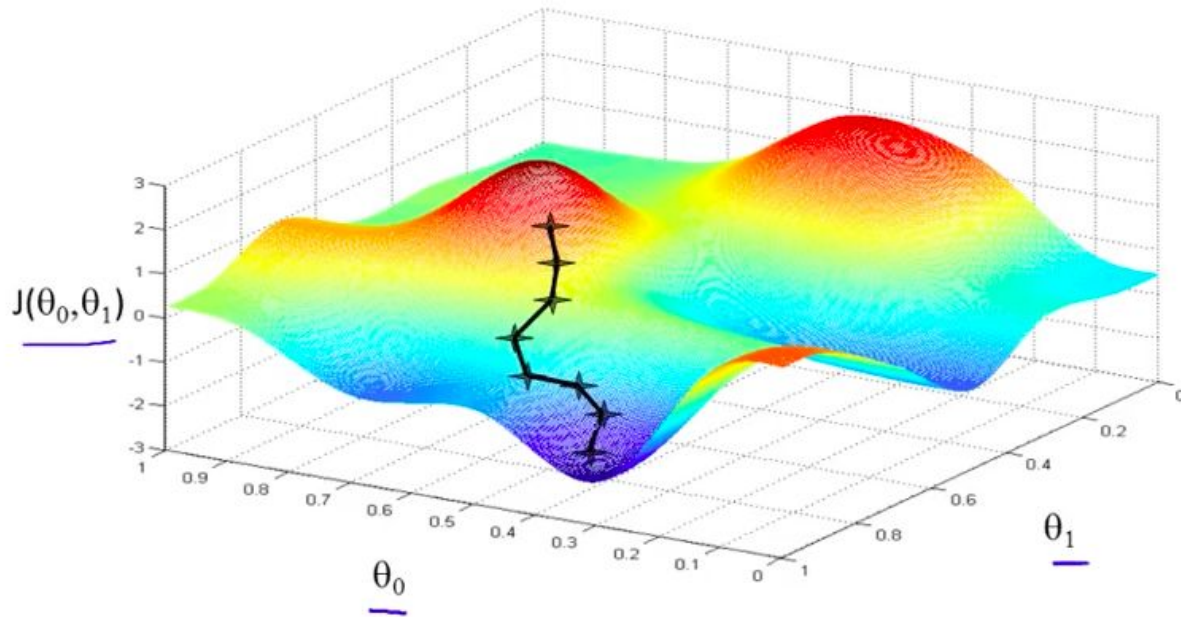
Обучение модели

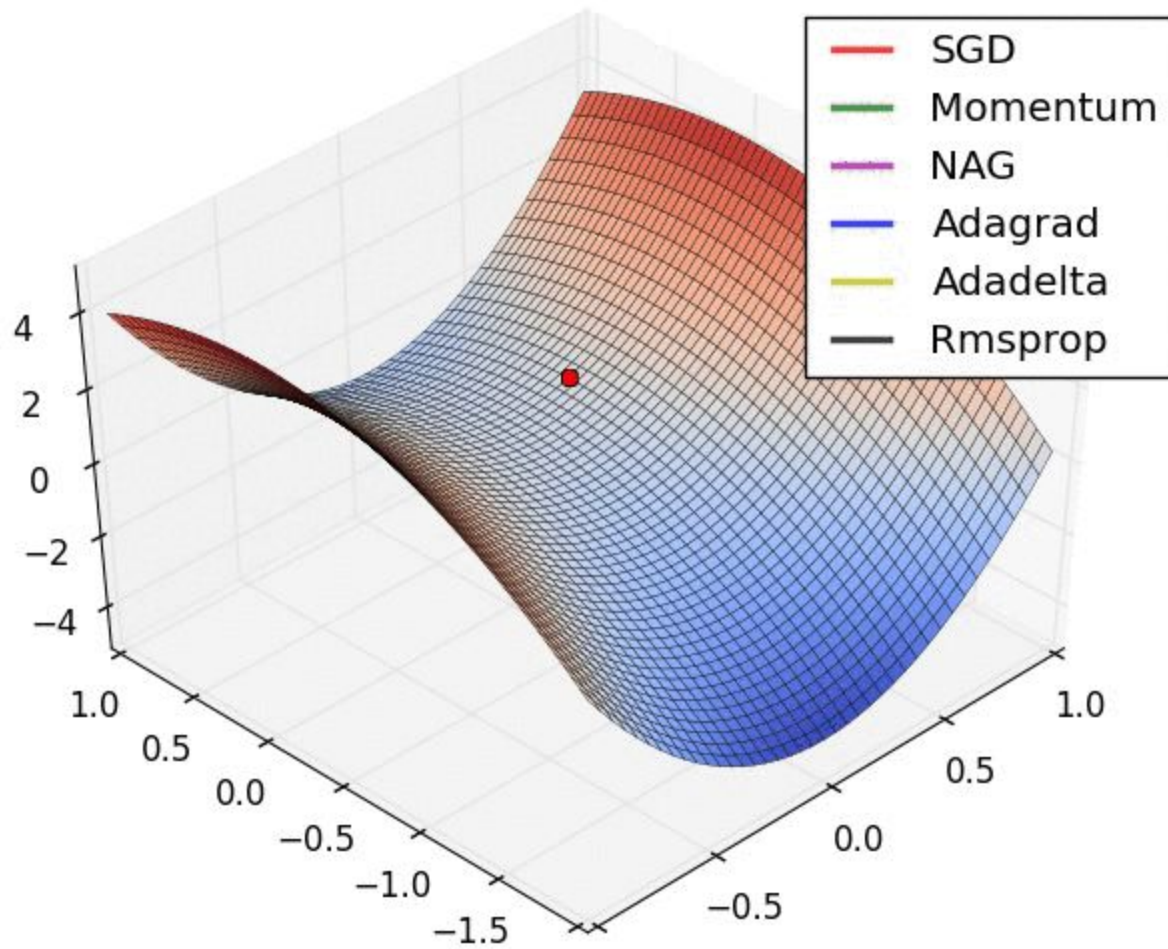
Input	Output	W=3	rmse(3)	W=3.0001	rmse
0	0	0	0	0	0
1	2	3	1	3.0001	1.0002
2	4	6	4	6.0002	4.0008
3	6	9	9	9.0003	9.0018
4	8	12	16	12.0004	16.0032
Total:	-	-	30	-	30.006

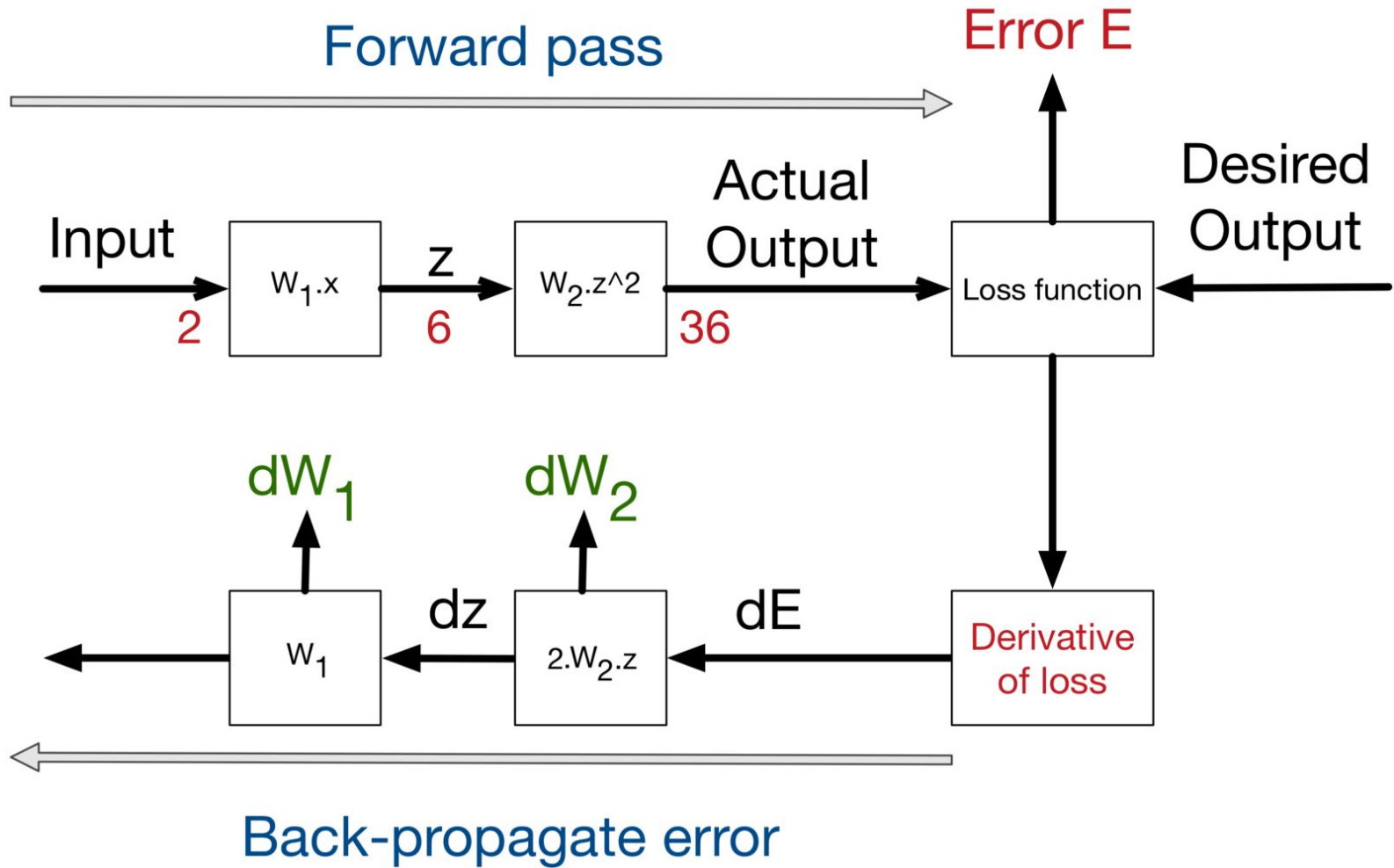
$W = 2.9999$, rmse - ?

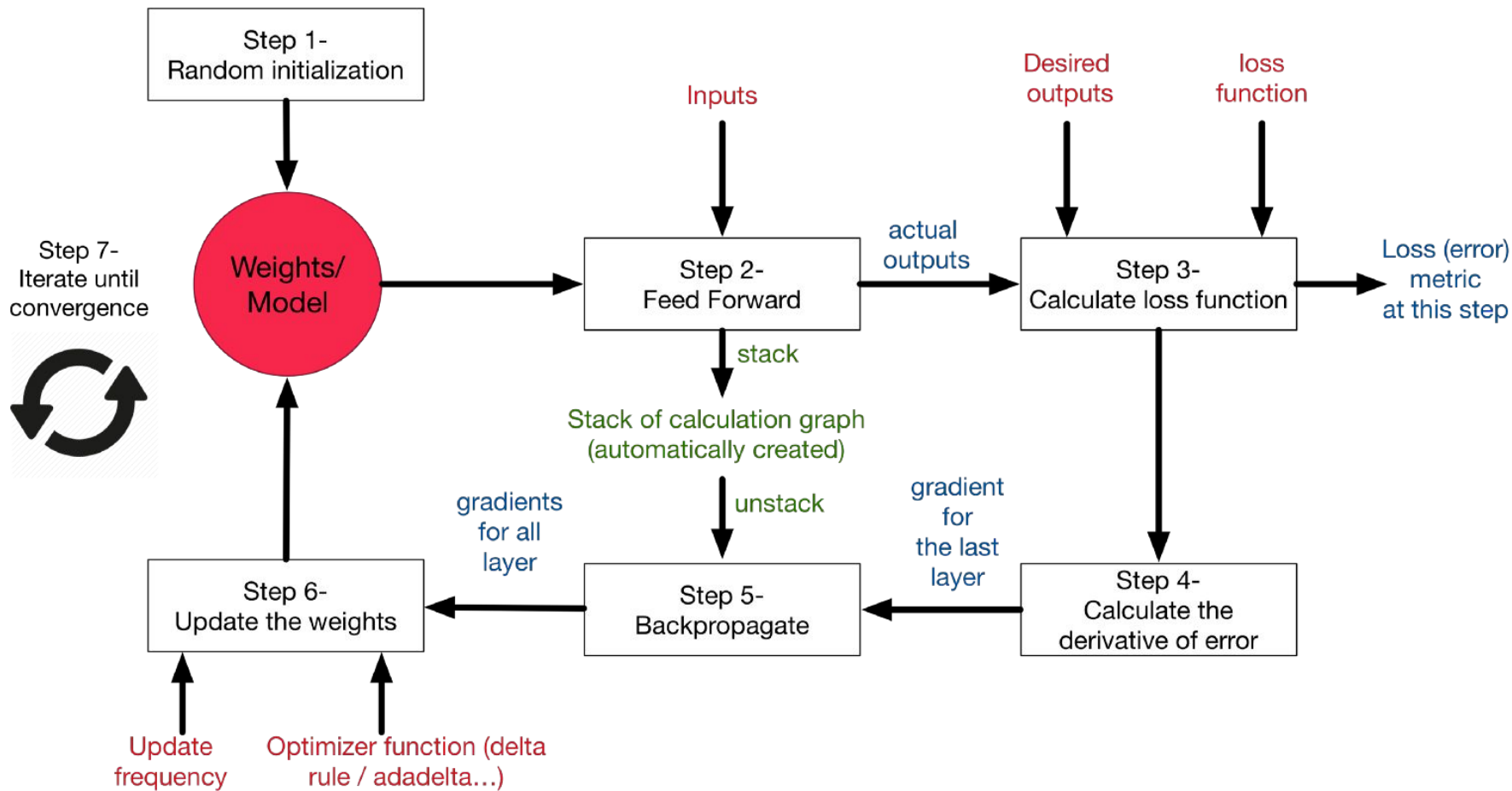


Стохастический градиентный спуск

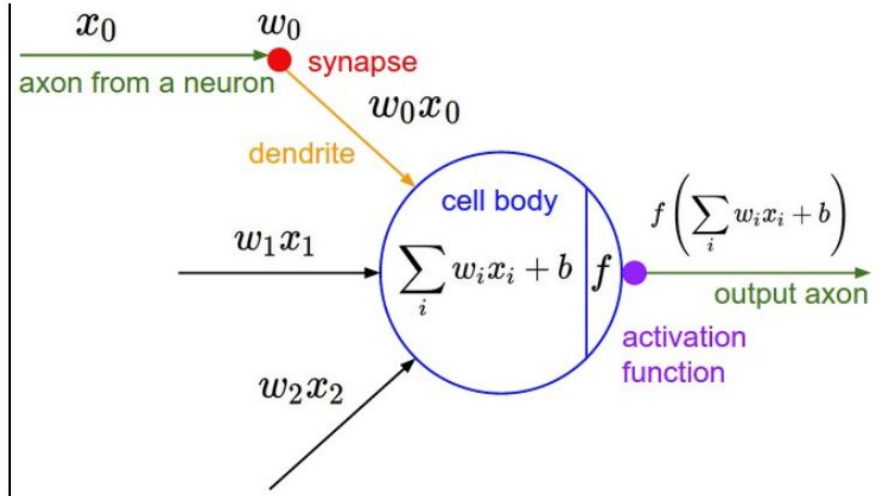
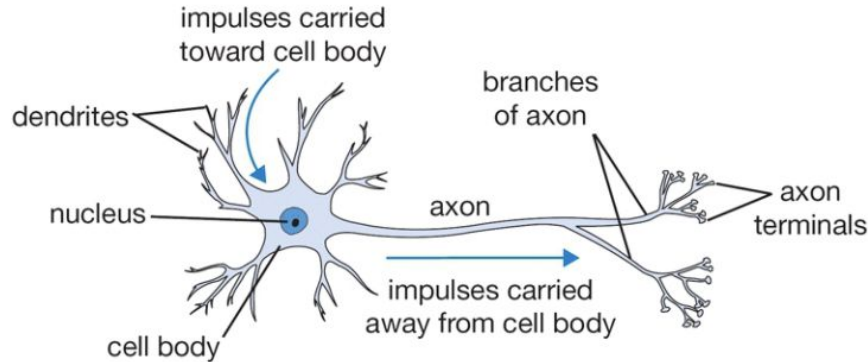




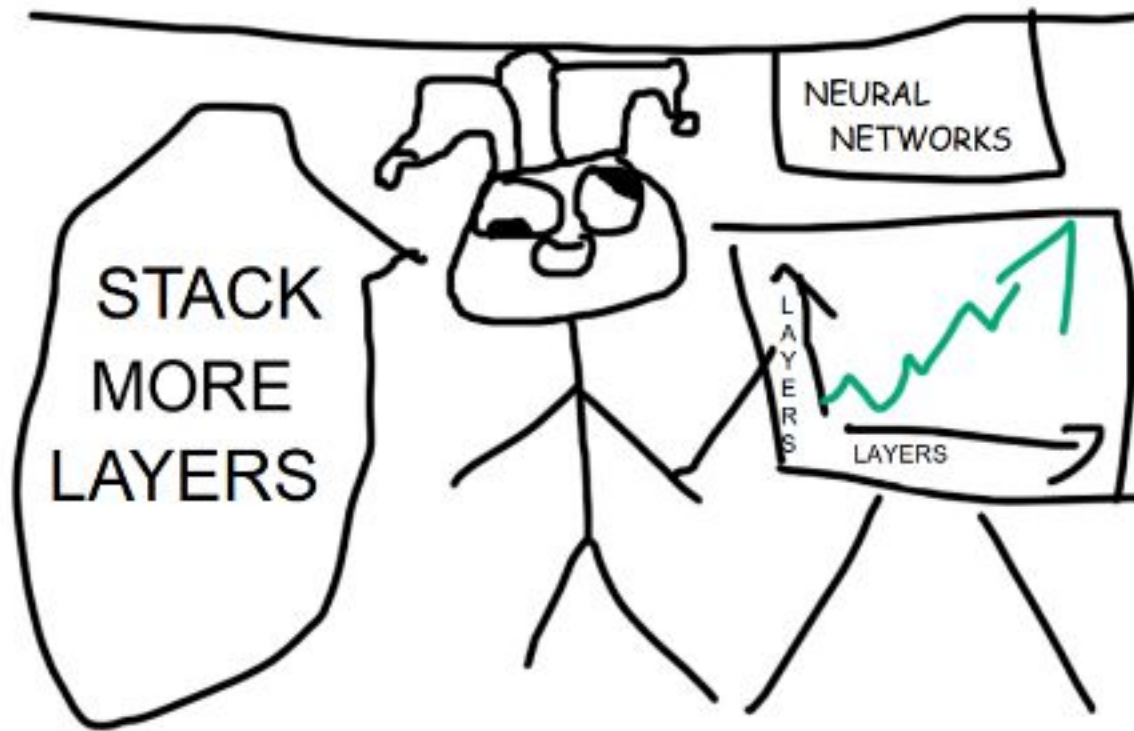




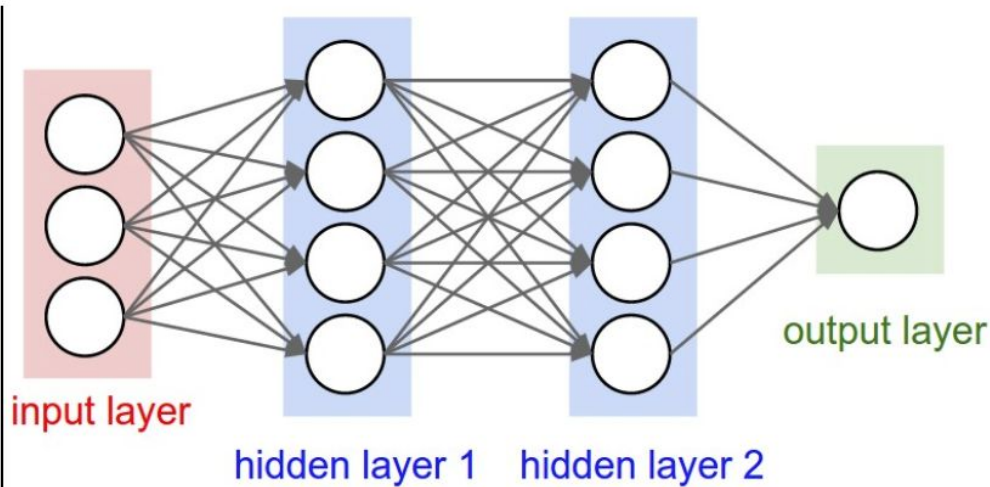
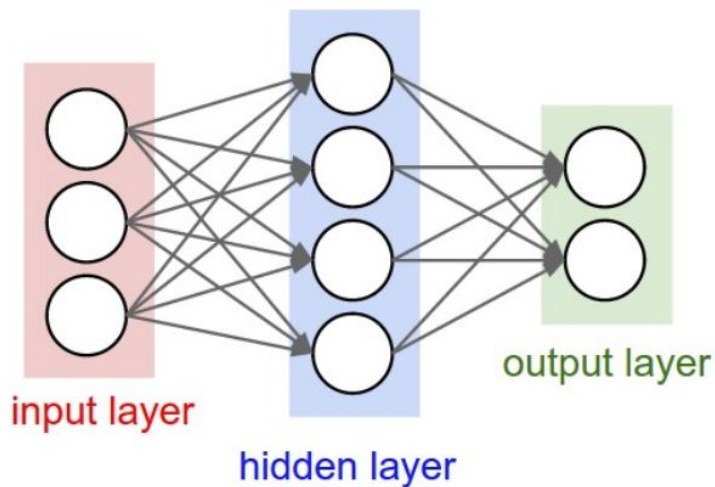
Neural Networks



A cartoon drawing of a biological neuron (left) and its mathematical model (right).



Neural Networks



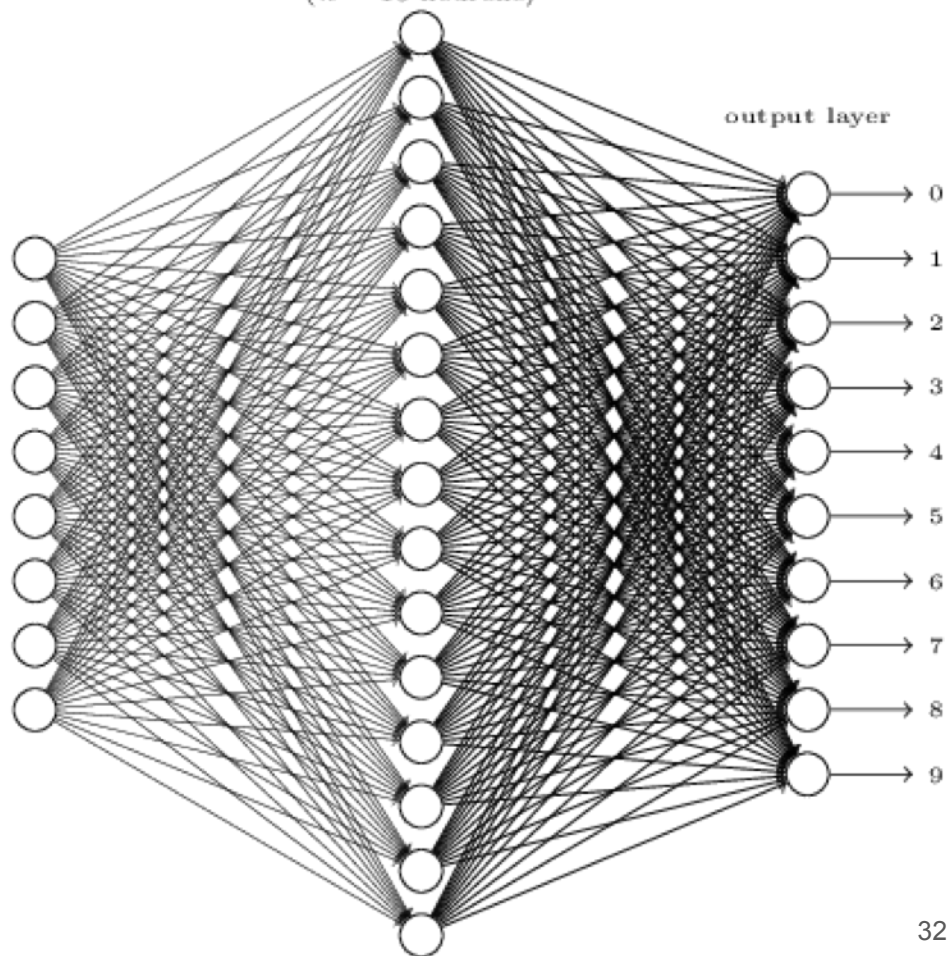
Left: A 2-layer Neural Network (one hidden layer of 4 neurons (or units) and one output layer with 2 neurons), and three inputs.
Right: A 3-layer neural network with three inputs, two hidden layers of 4 neurons each and one output layer. Notice that in both cases there are connections (synapses) between neurons across layers, but not within a layer.

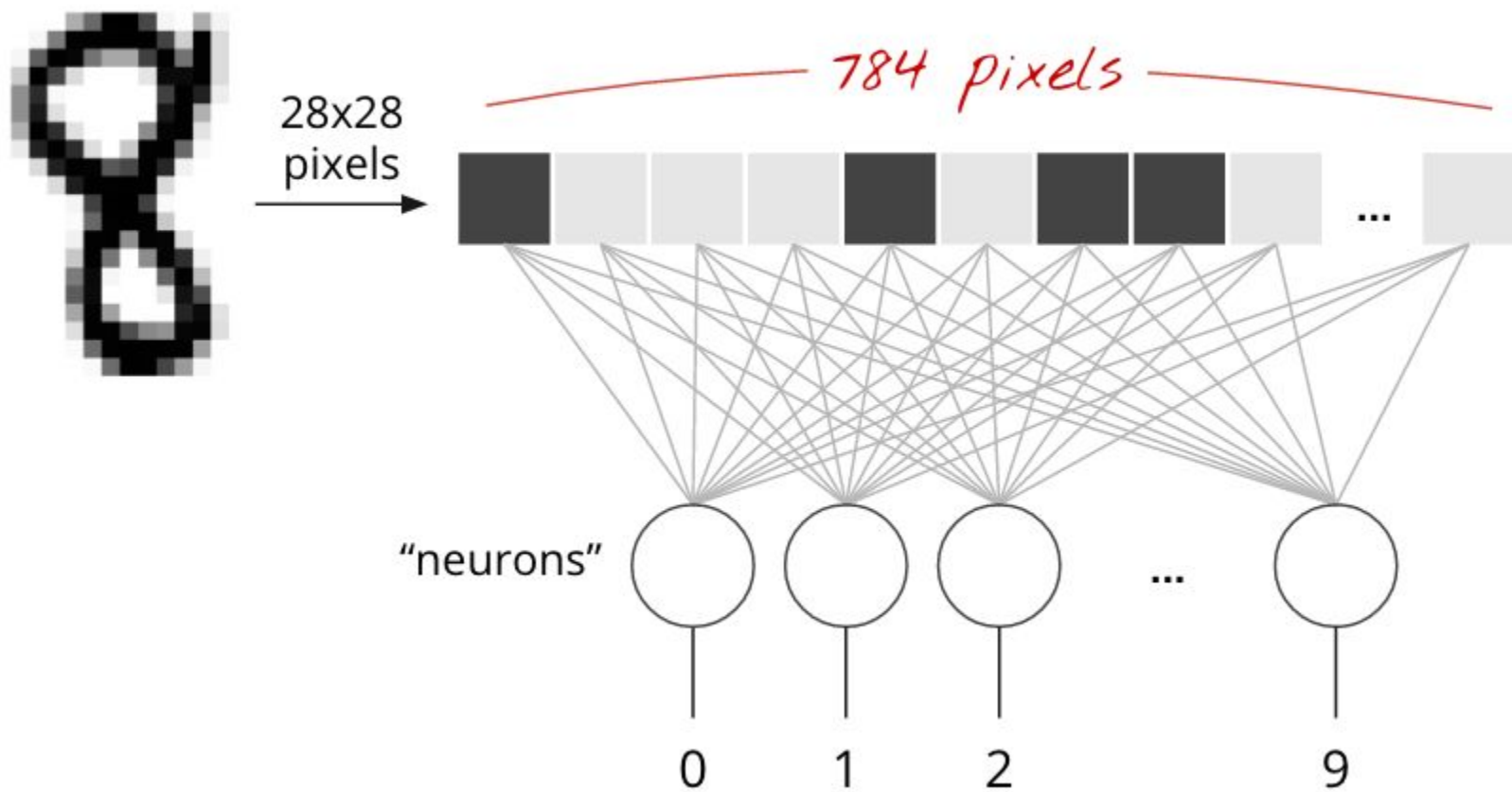
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9

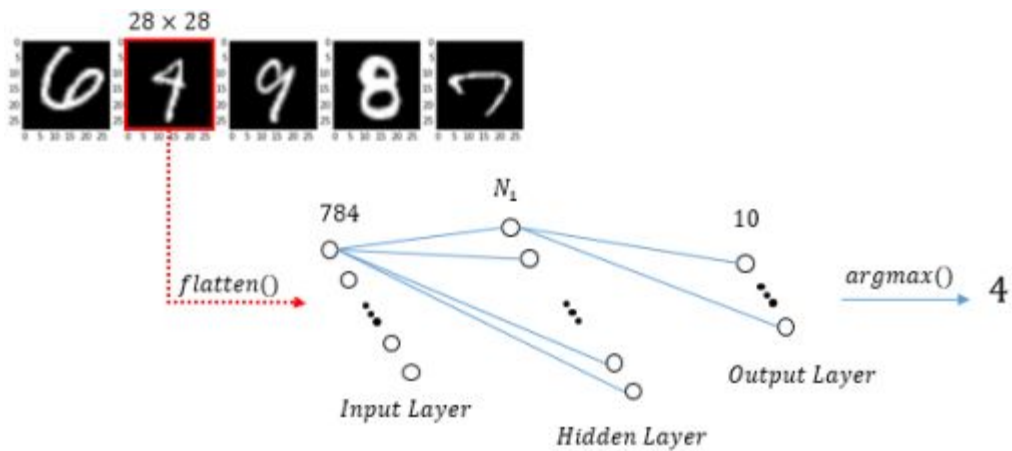
input layer
(784 neurons)

hidden layer
($n = 15$ neurons)

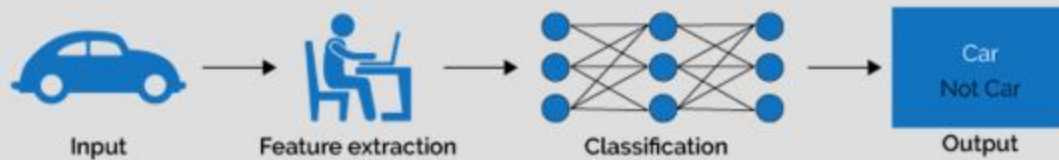
output layer



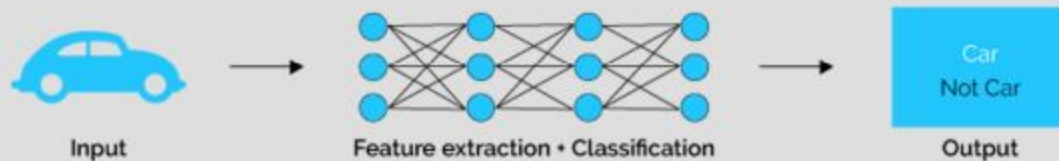




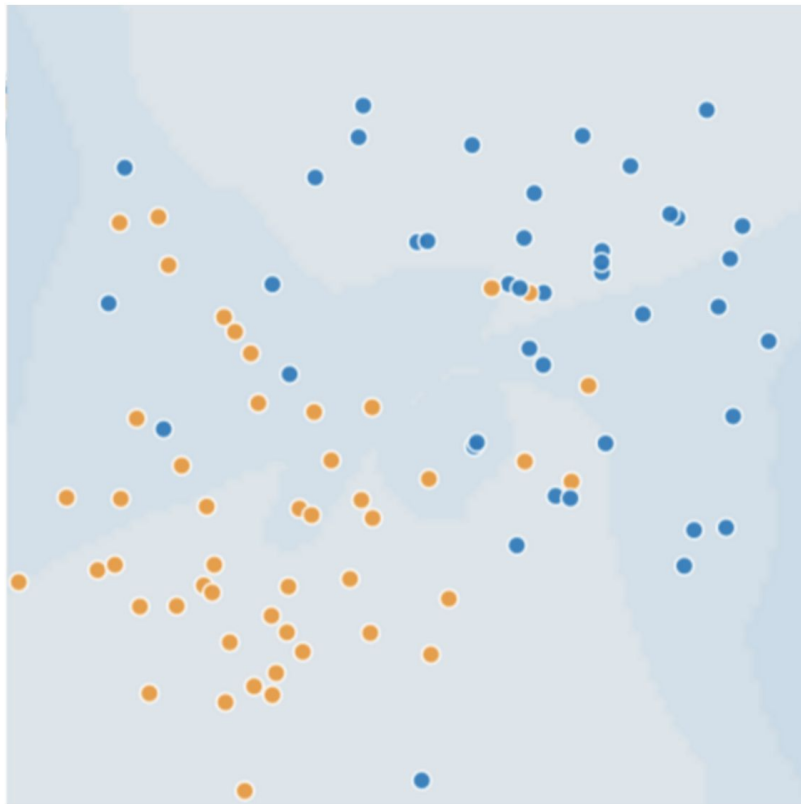
Machine Learning



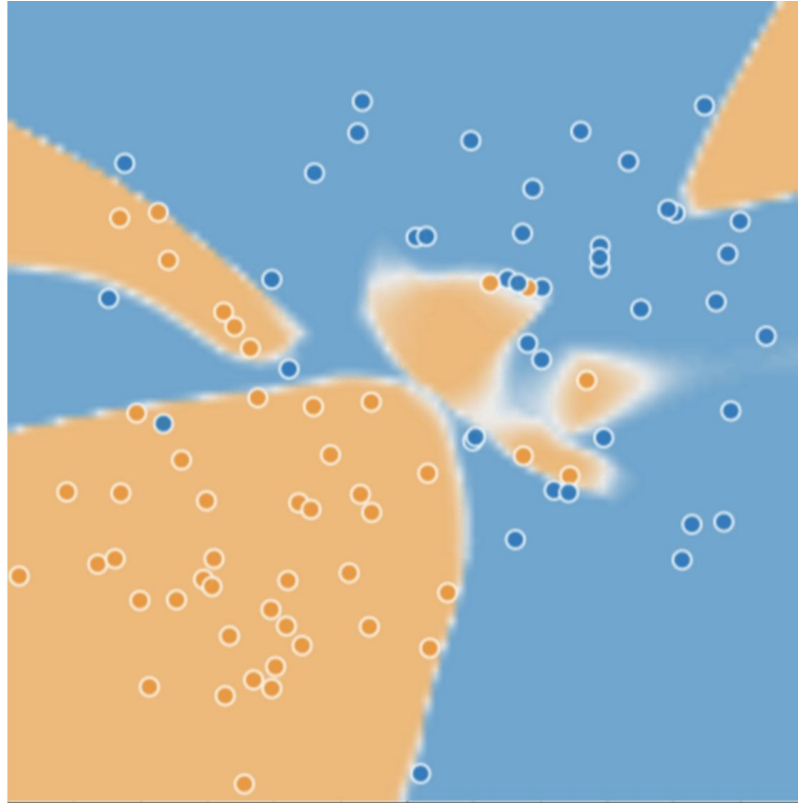
Deep Learning



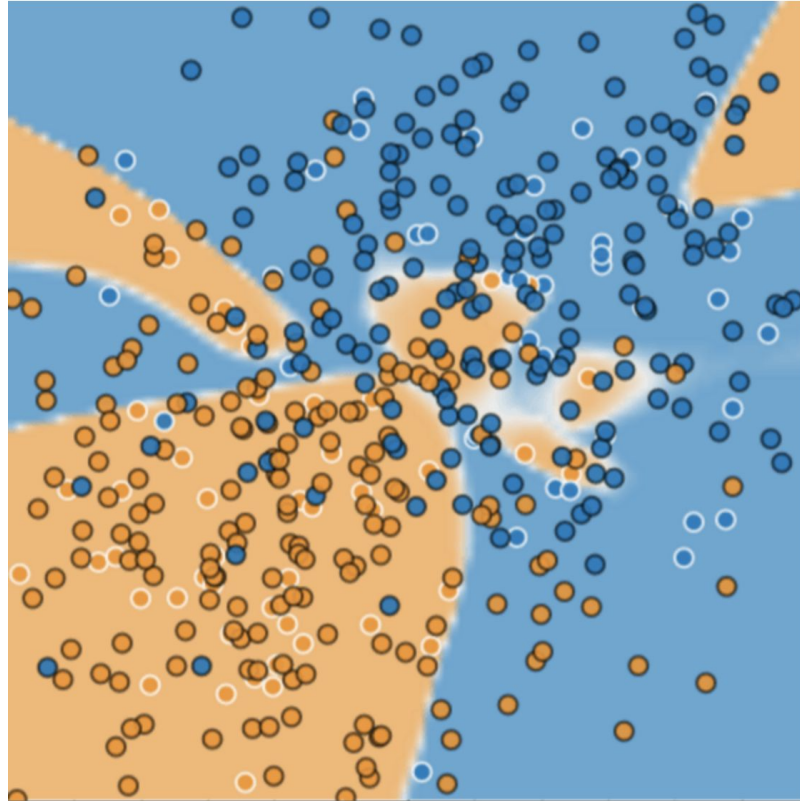
Overfitting



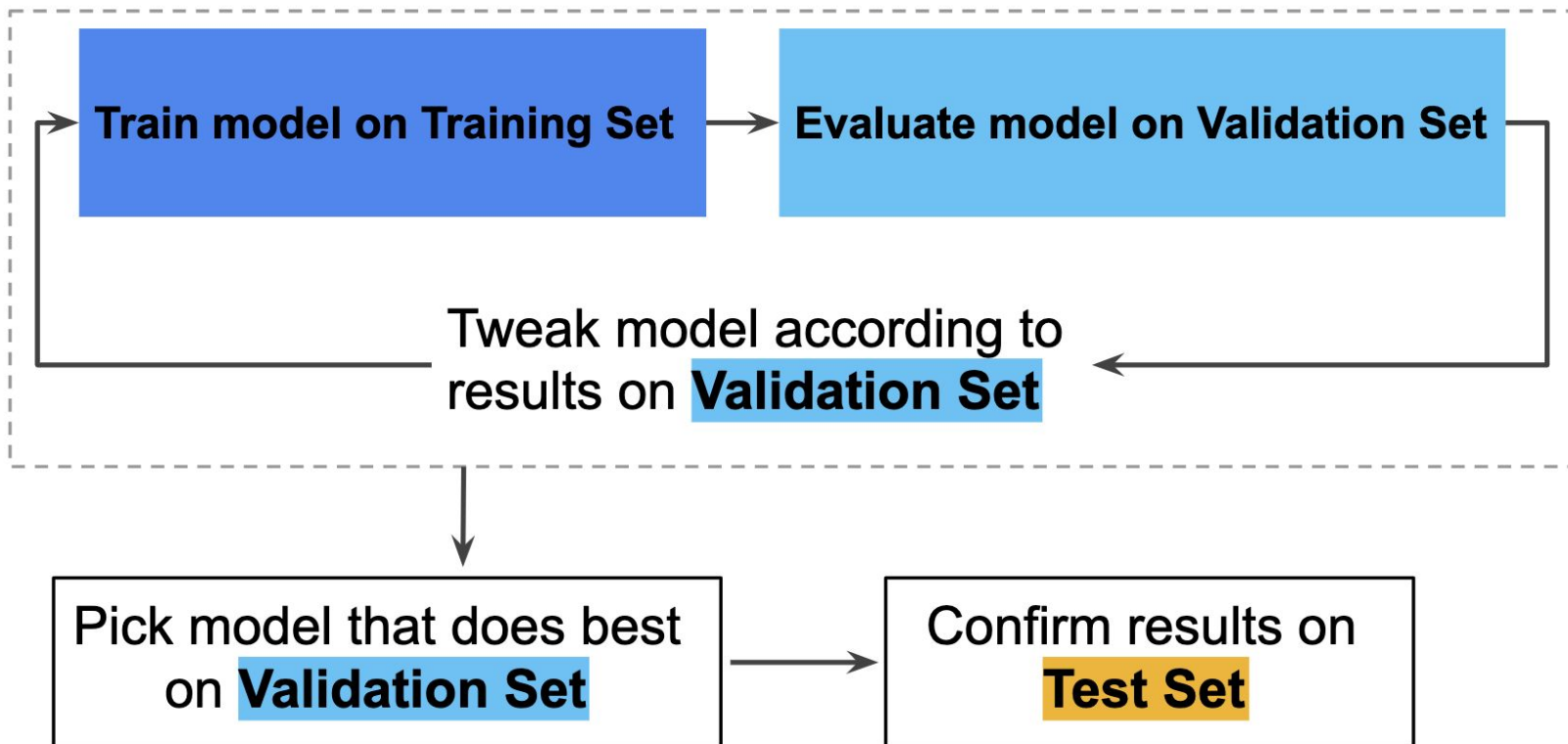
Overfitting



Overfitting



Train/val/test splitting



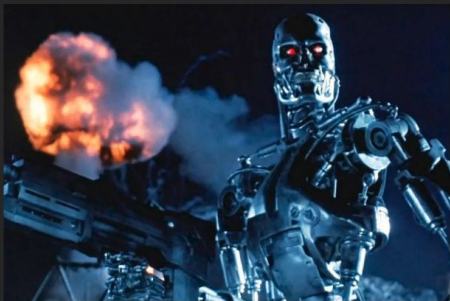
Cross-validation

Split 1	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 1
Split 2	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 2
Split 3	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 3
Split 4	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 4
Split 5	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 5

Training data

Test data

Deep Learning



What society thinks I do



What my friends think I do



What other computer scientists think I do



What mathematicians think I do



What I think I do

```
from keras import *
```

What I actually do

Neural Networks: demo

- ConvNetJS - <https://cs.stanford.edu/people/karpathy/convnetjs/>
- A Neural Network Playground - <https://playground.tensorflow.org/>

Спасибо
за просмотр!

