



Central European Institute of Technology
BRNO | CZECH REPUBLIC

Deep learning for biological data

Mgr. Vlastimil Martinek
Mgr. David Cechak





Goals of the Workshop

1. Machine learning **basics**
2. **Demystify** deep learning
3. **Pre-Processing** genomic data
4. Programming **neural nets** for genomics

Out of Scope

1. **Heavy math**
2. Easiest way to do/use neural networks
3. **Elaborate** architecture and training schemes
4. Interpretability
5. Able to deploy **inference server**

Agenda

- 1) Introduction
- 2) Linear regression
- 3) Logistic regression classifier
- 4) Multilayer perceptron (MLP)
- 5) Convolutional neural network (CNN)
- 6) Extras

Workshop Tools



Solving a problem

$$f(X) = y$$

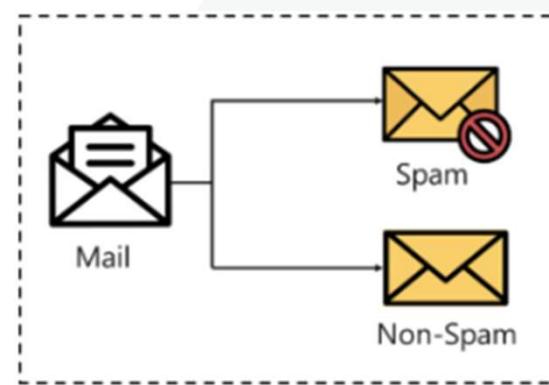
X = biological sequence, microscopy image, molecule, text, map as a graph...

y = contains a gene?, tumor position and volume, binding site, shortest path...

f traditionally designed by hand

What if we let a machine design the function **f**?

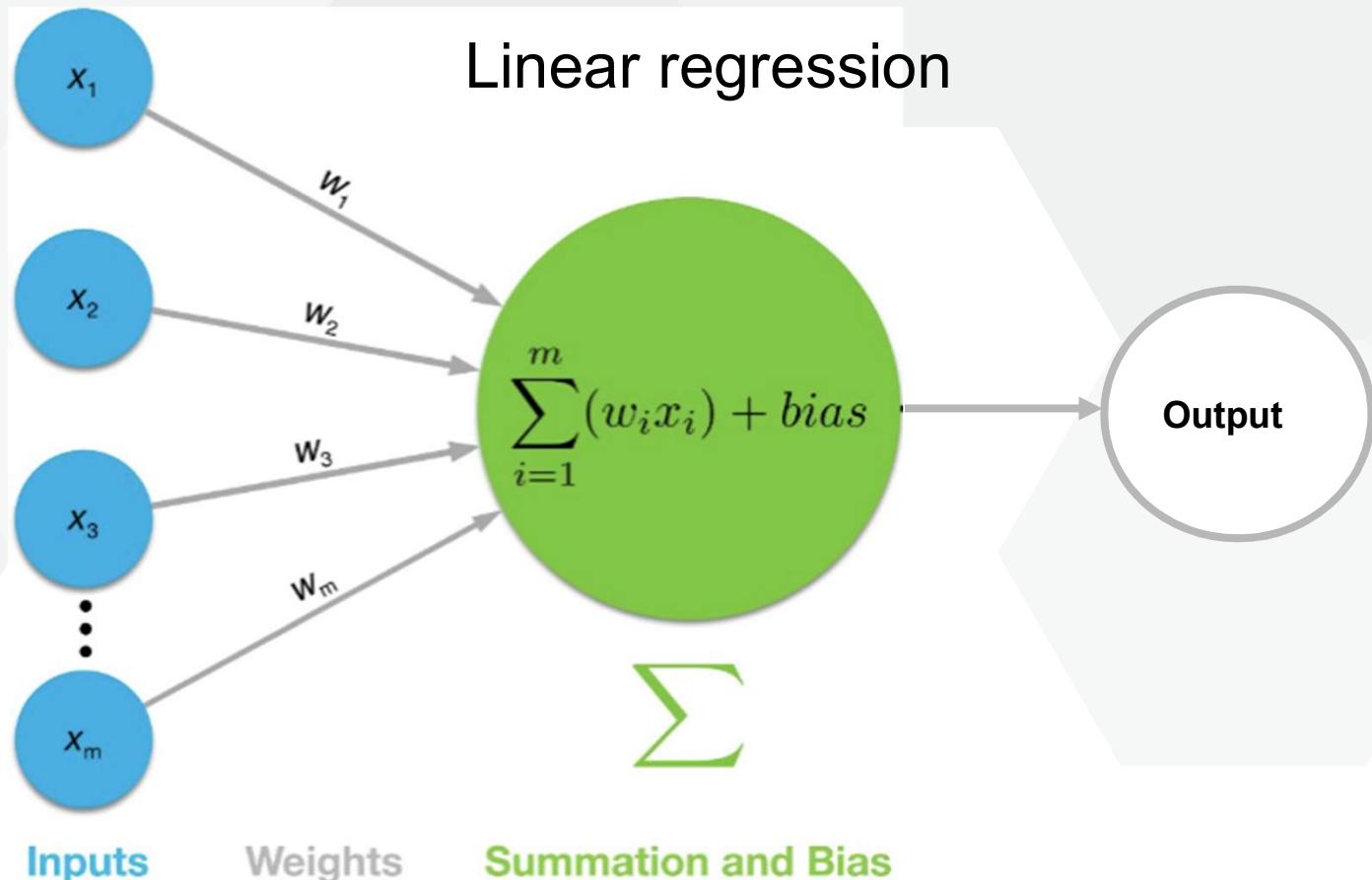
Goal: have a useful predictor



Linear regression

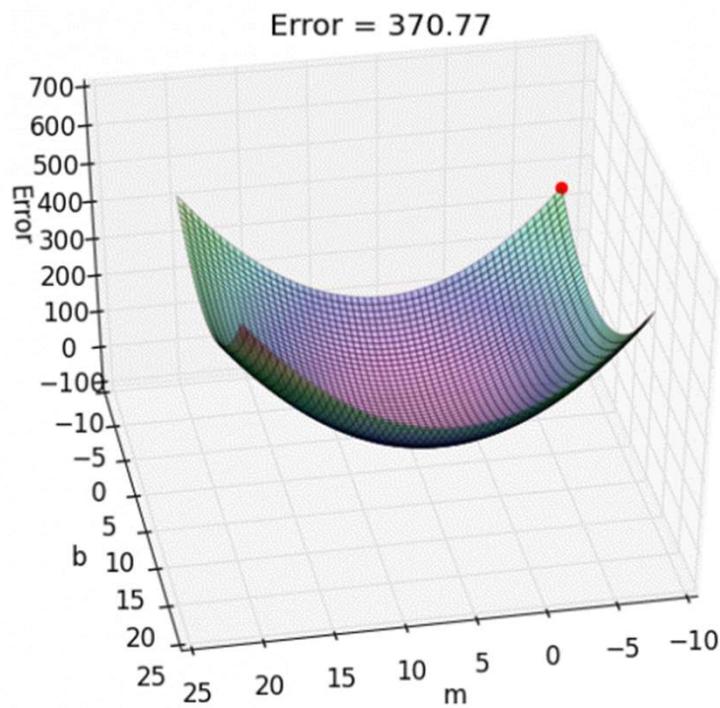


Linear regression



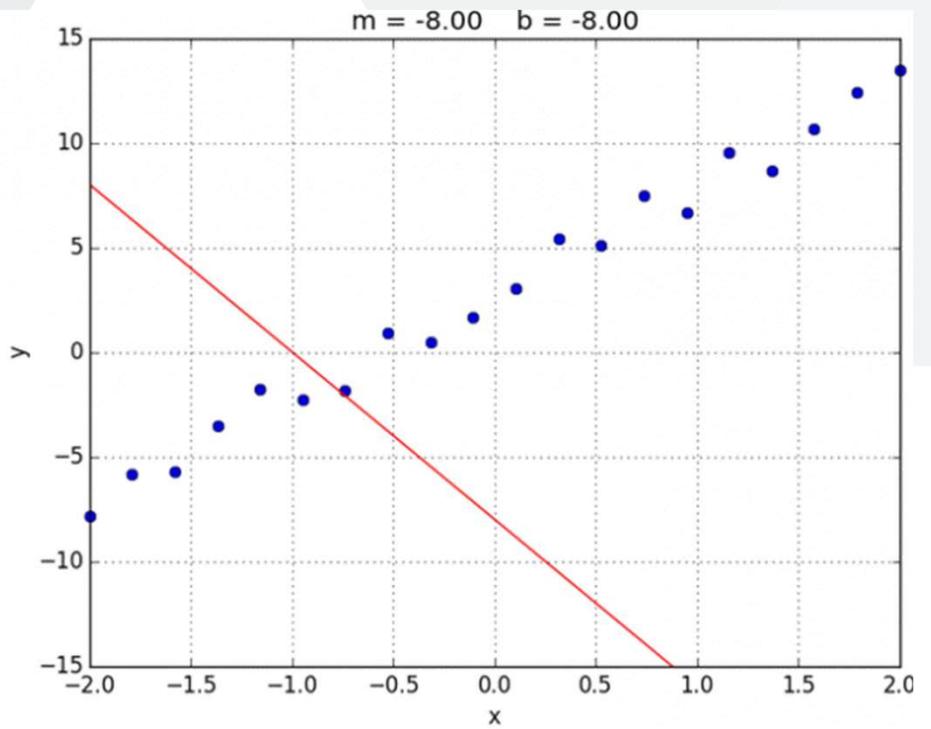
Optimization process - fitting our data

Error = MSE



CEITEC

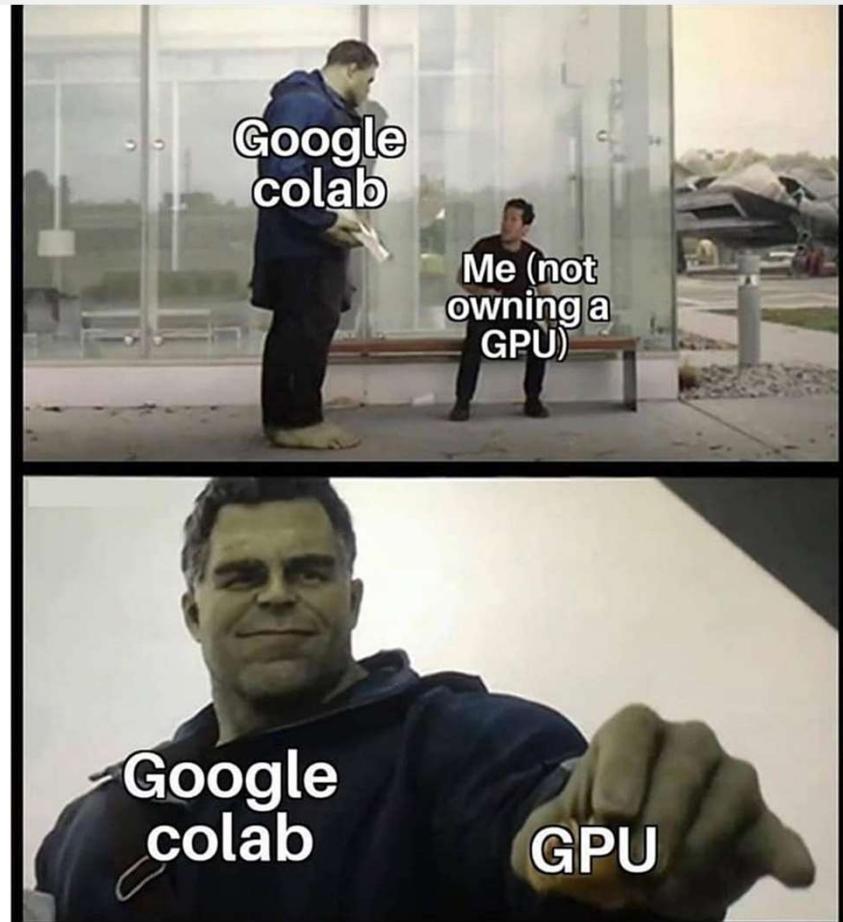
$Y = mX+b$



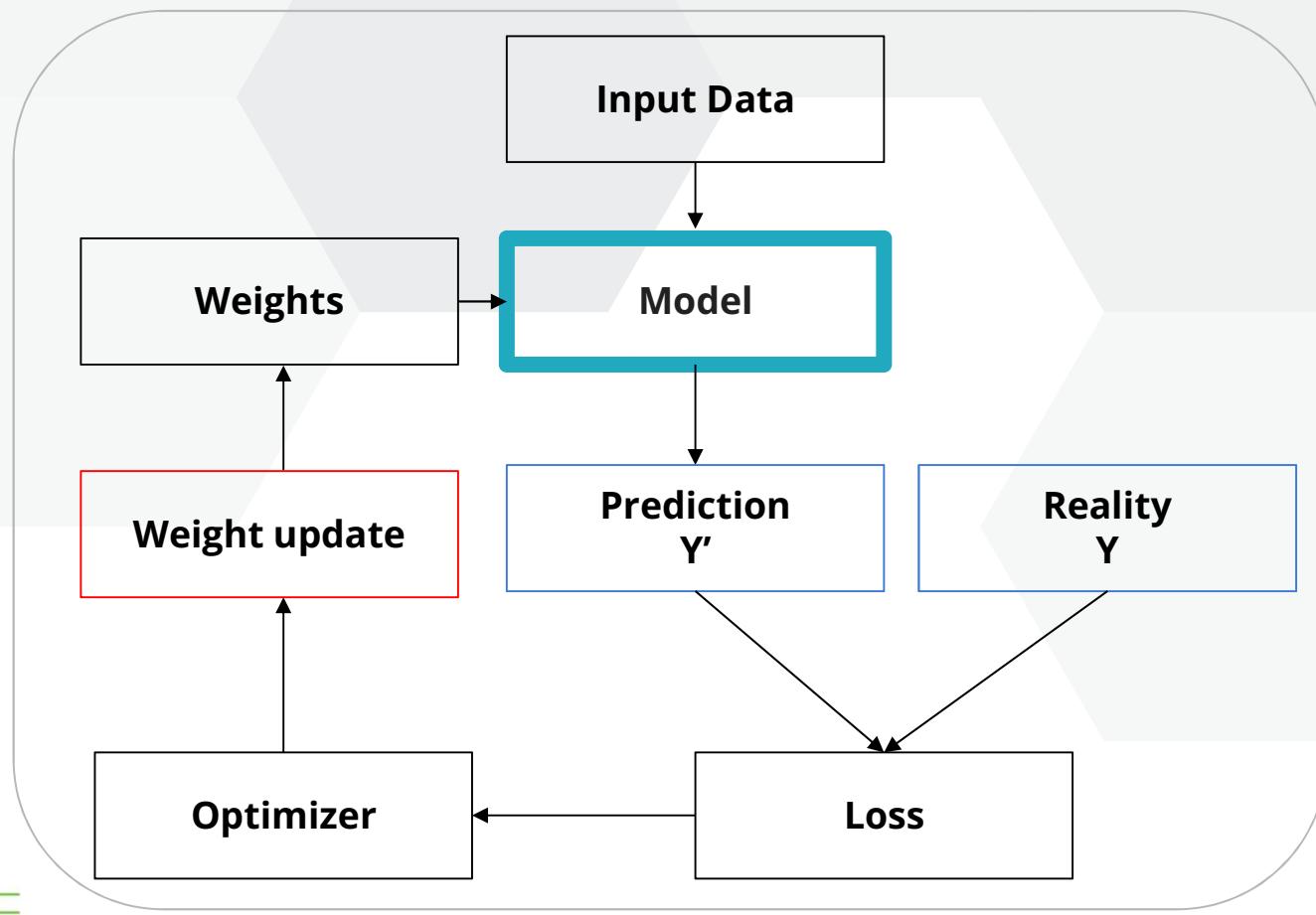
Practical part

[https://github.com/BioGeMT/
MALTAomics-Summer-
School/](https://github.com/BioGeMT/MALTAomics-Summer-School/)

Day 2 Workshop
MALTA_01 Notebook



Task: Learn Optimal Weights to Minimize Loss



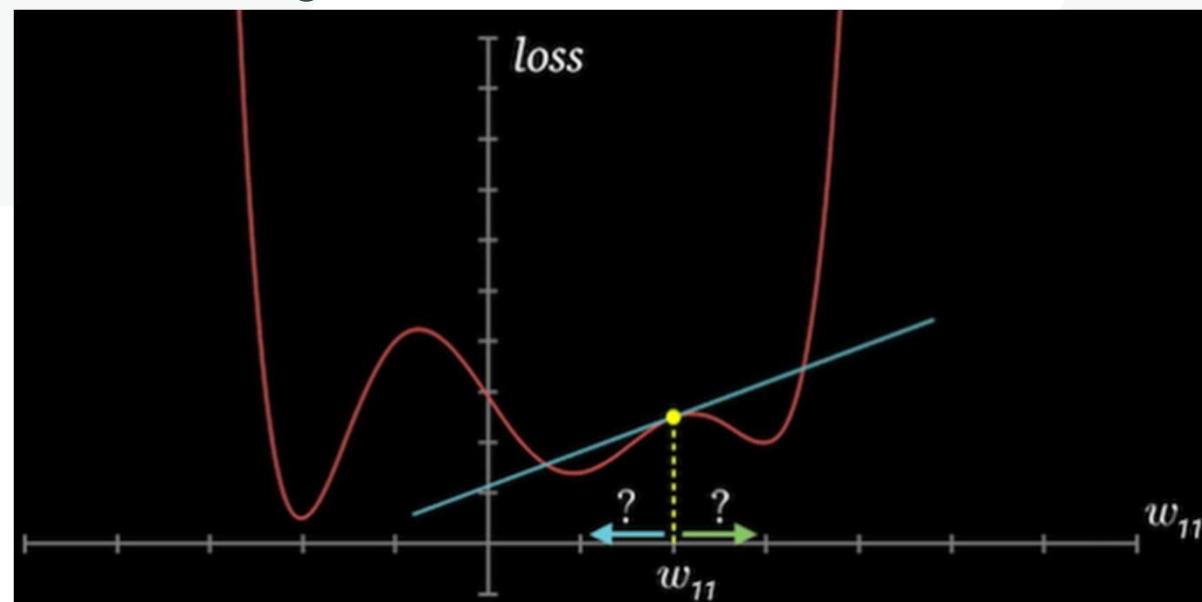
Gradient descent

Compute a gradient for a **single** weight w.r.t the loss

Average over the training data

Indicates the rate of change

$\text{weight} = \text{weights} - \text{weight.grad} * 1e-5$
 $\text{bias} = \text{bias} - \text{bias.grad} * 1e-5$



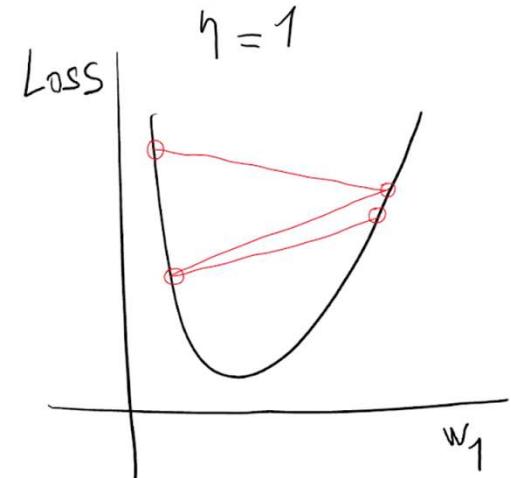
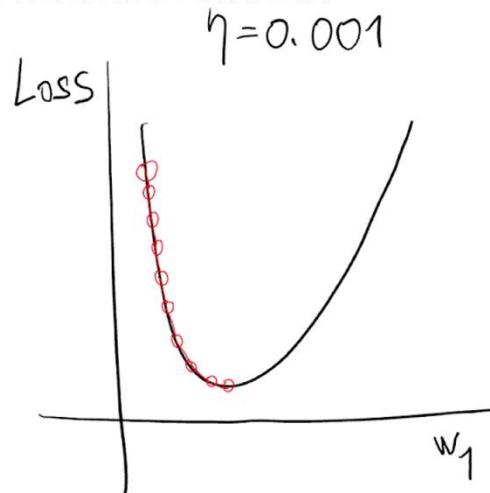
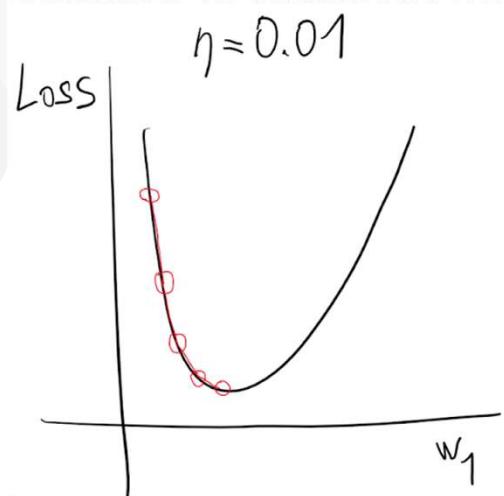
Epochs, learning rate, loss function

Learning rate - magnitude of the weight update

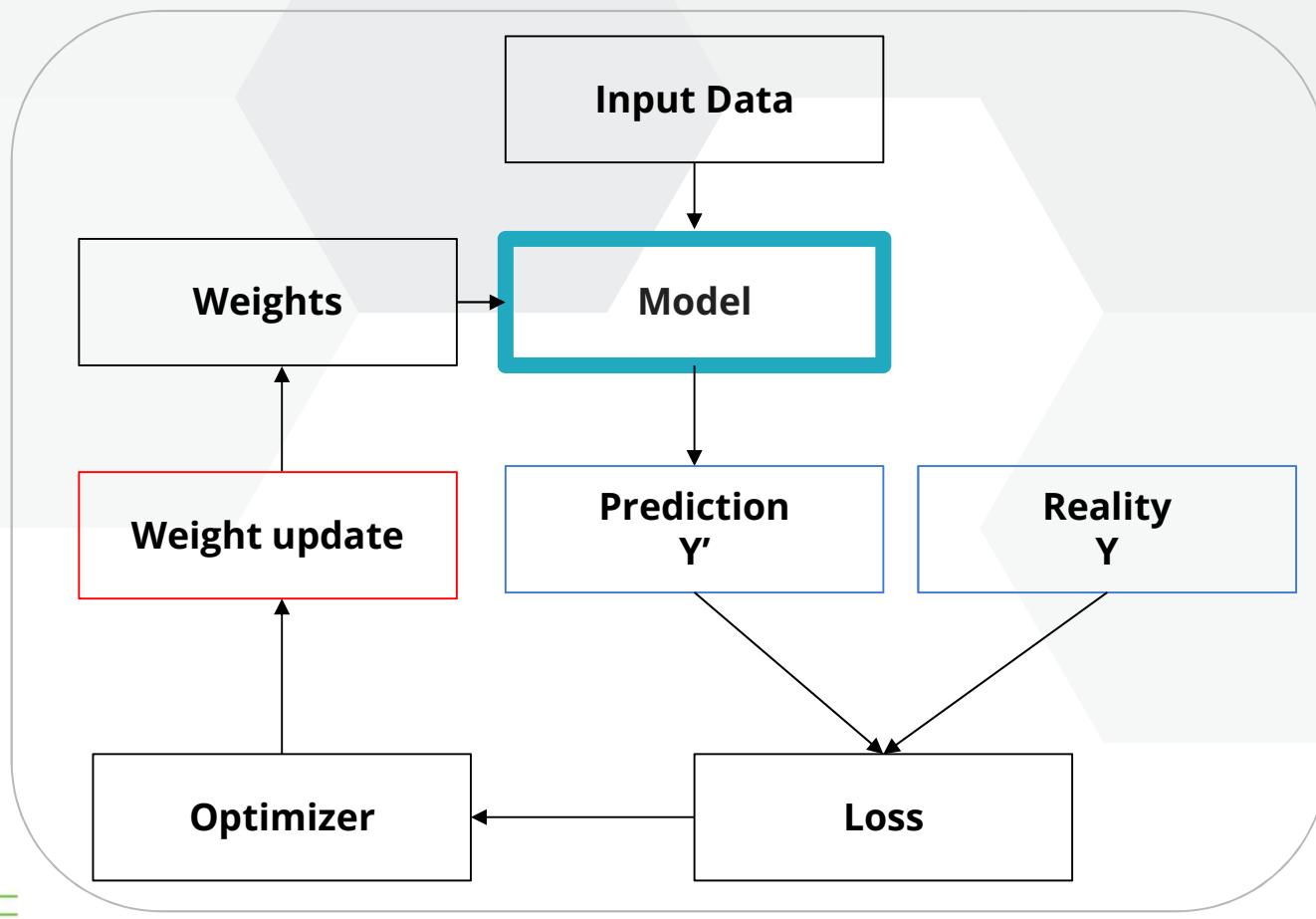
Loss function - objective to minimize

Tensor - pytorch array able to track gradient

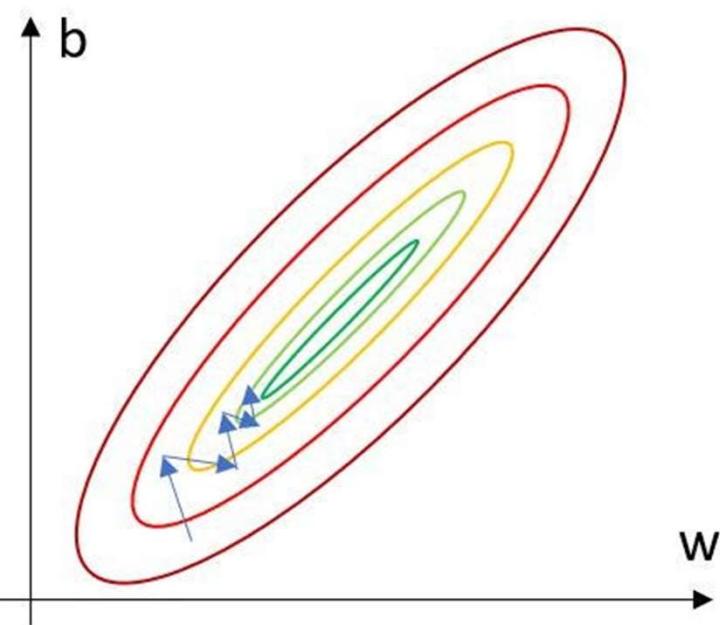
Epochs - number of cycles in the learning phase



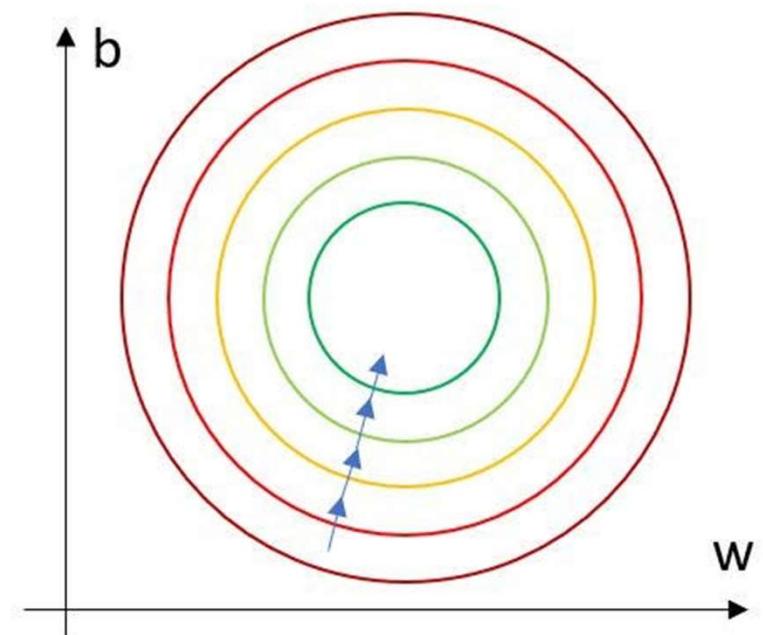
Task: Learn Optimal Weights to Minimize Loss



Unnormalized:



Normalized:

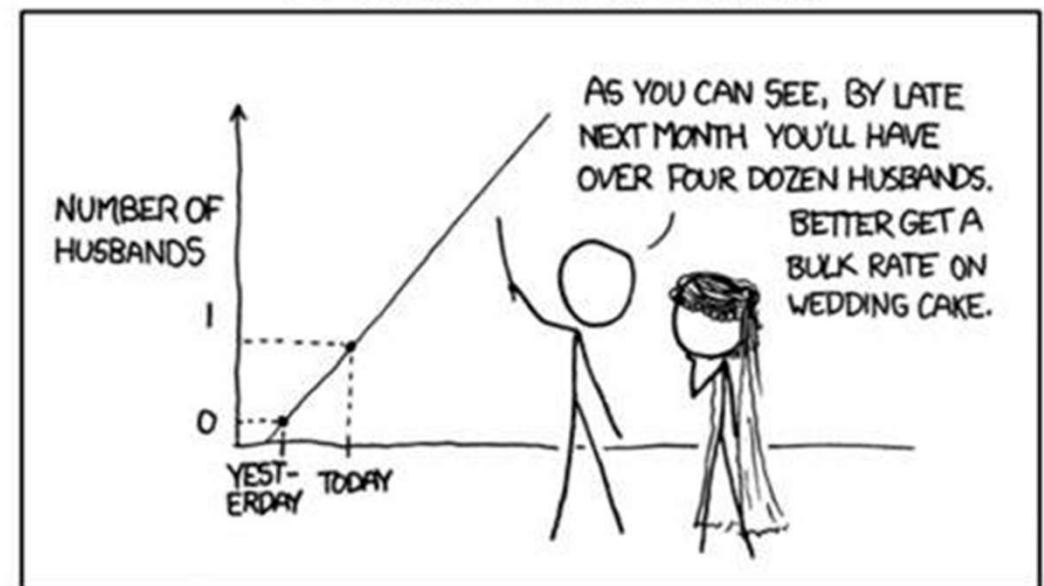


DS: How is model training going?
Intern: No so good. It's not learning anything...
DS: Show me your train loop
Intern: ...

```
1 import numpy as np
2 import pandas as pd
3 import torch
4
5 def train():
6     print('      o 0')
7     print(' _][__|o| |0 0 0 0| ')
8     print(' <_____| - |_____| ')
9     print(' /0 - 0 0 0 0 0 ')
10
11 train()
```

Short break. Back at 11:12

MY HOBBY: EXTRAPOLATING



Quiz

1. What shape should be the training data for our Linear layer?
2. What is a loss function?
3. What is mean squared error?
4. What happens when you invoke the .backward function on the result of the mean squared error loss function?
5. How do you update the weights and biases of a model?



model trained
for 1000 epochs



model trained
for 100 epochs



model trained
for 1 epoch

Quiz

6. Why do you reset gradients to zero after updating weights?
7. What is an epoch?
8. What is the benefit of training a model for multiple epochs?
9. What should you do if your model's loss doesn't decrease while training?
10. What are the inputs to a PyTorch optimizer?
11. How could we feed 2 features as input?



model trained
for 1000 epochs



model trained
for 100 epochs



model trained
for 1 epoch

Linear regression summary

$$Y = wX + b$$

Regression problem

Common with deep learning

- Needs input data
- Starts with random weights/parameters
- Computes loss
- Repeatedly optimizes loss to improve (a tiny bit) on training data (epochs)
- Resulting model has updated weights

Differences

- Linear only

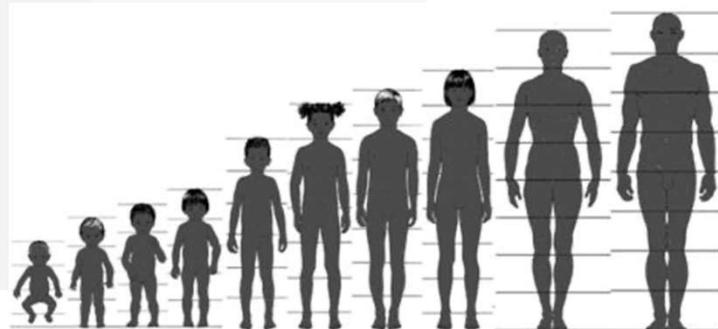
Classification Logistic regression

Is this DNA sequence a non-TATA promoter?

$f(X) = y$

$X = \text{height}$

$y = \text{weight?}$



$f(X) = y$

$X = \text{DNA sequence}$

$y = \text{contains a non-T}^{\wedge}\text{T A}$
 promoter?



Promoter

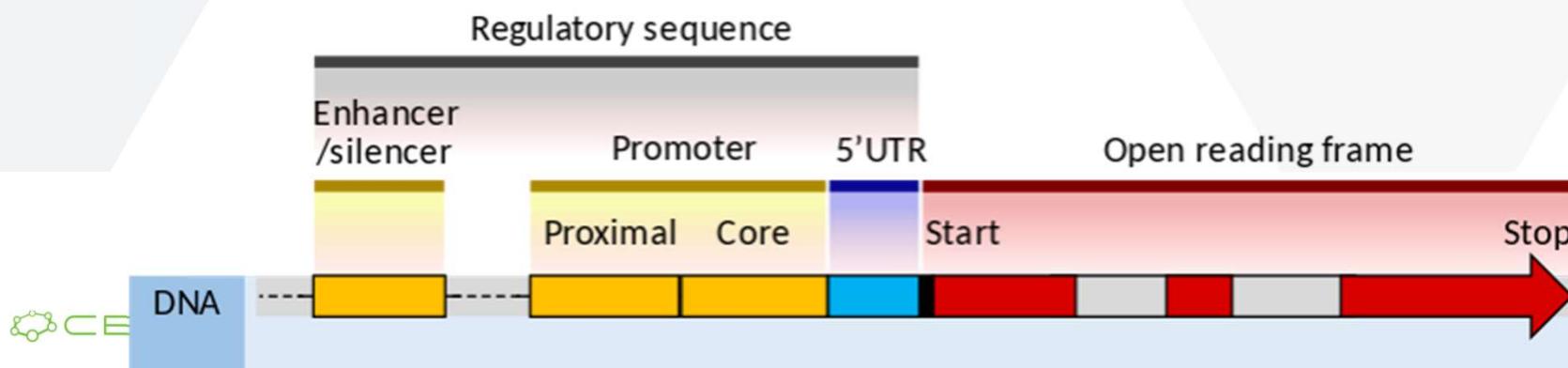
- Regulatory sequences
- Where proteins bind to initiate transcription of a gene

TATA

- drive **burst-like gene expression** (dosage)

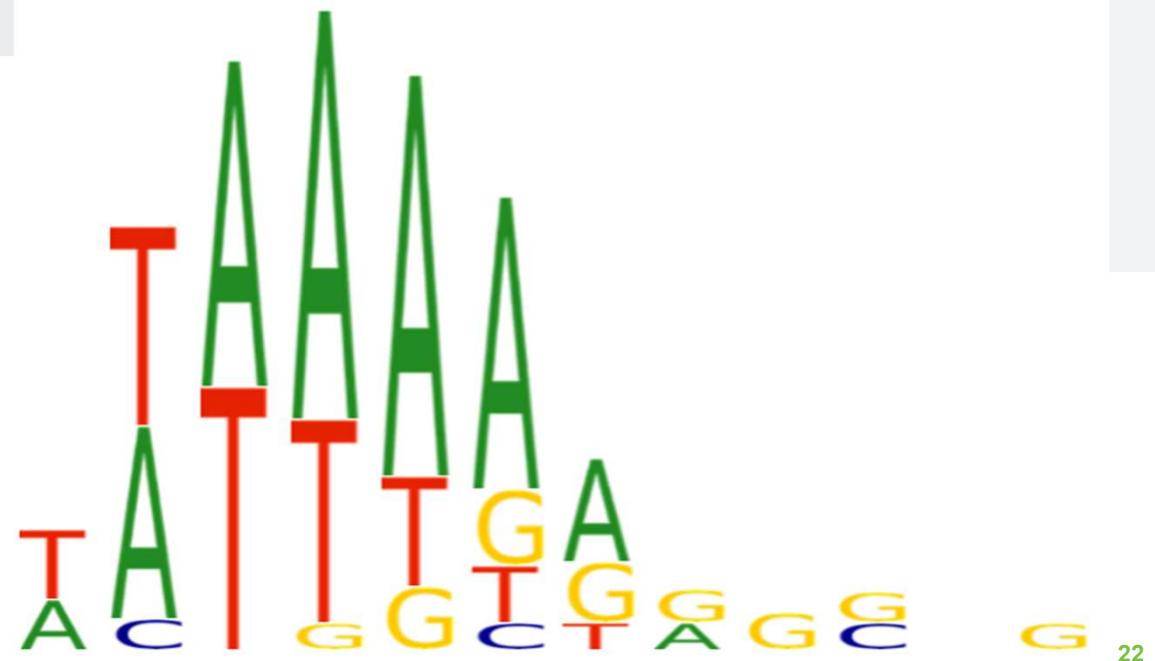
non-TATA

- govern housekeeping genes - **consistent expression** (continuous)



TATA-box region - motif

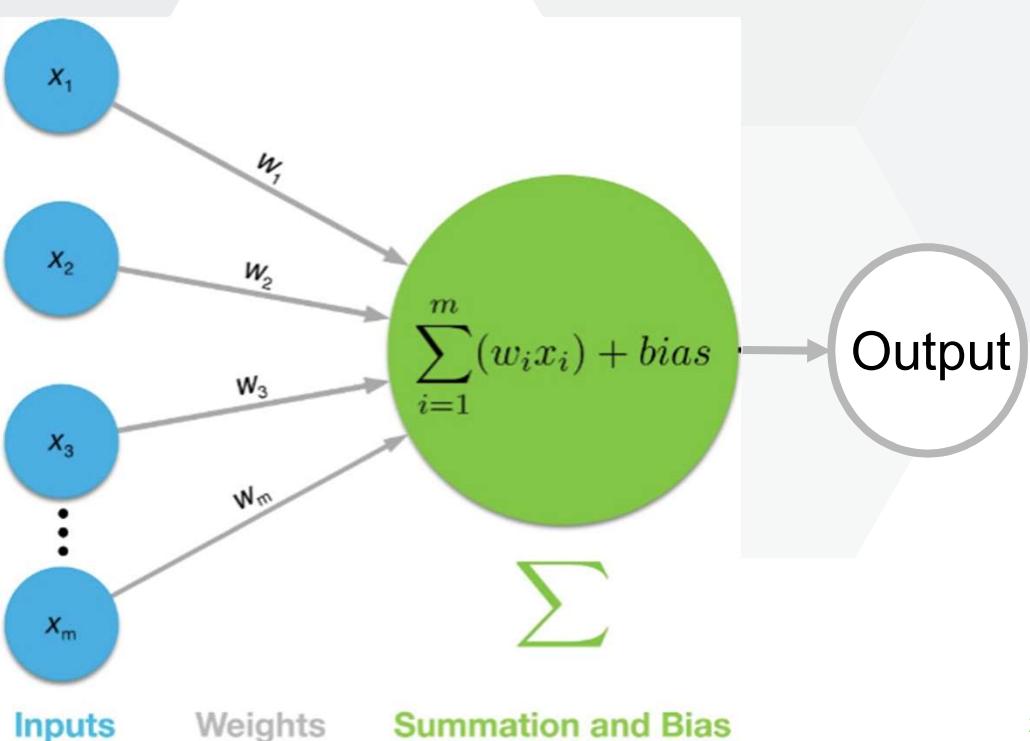
Sub sequence pattern (motif) for human TATA promoter sequences



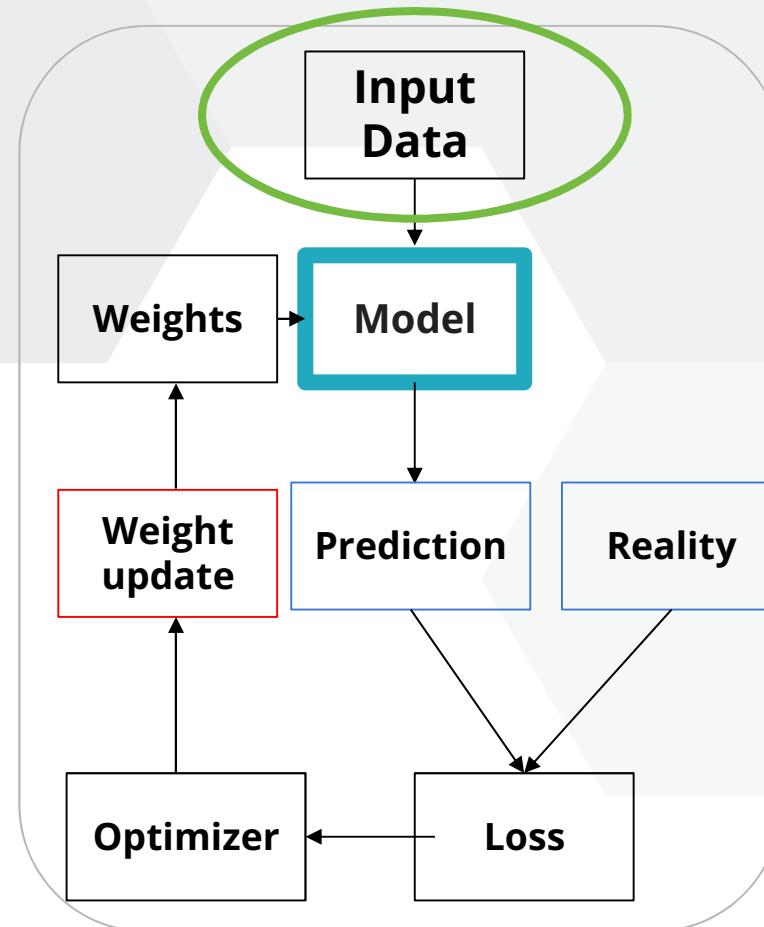
First problem: Numericalization -> Tokenization



ACC
G
T
C
A
C
C
T
A



Numericalization & Tokenization



Numericalization & Tokenization

ACCGT CACCTACGT CGATTACGGGGTCATGACC

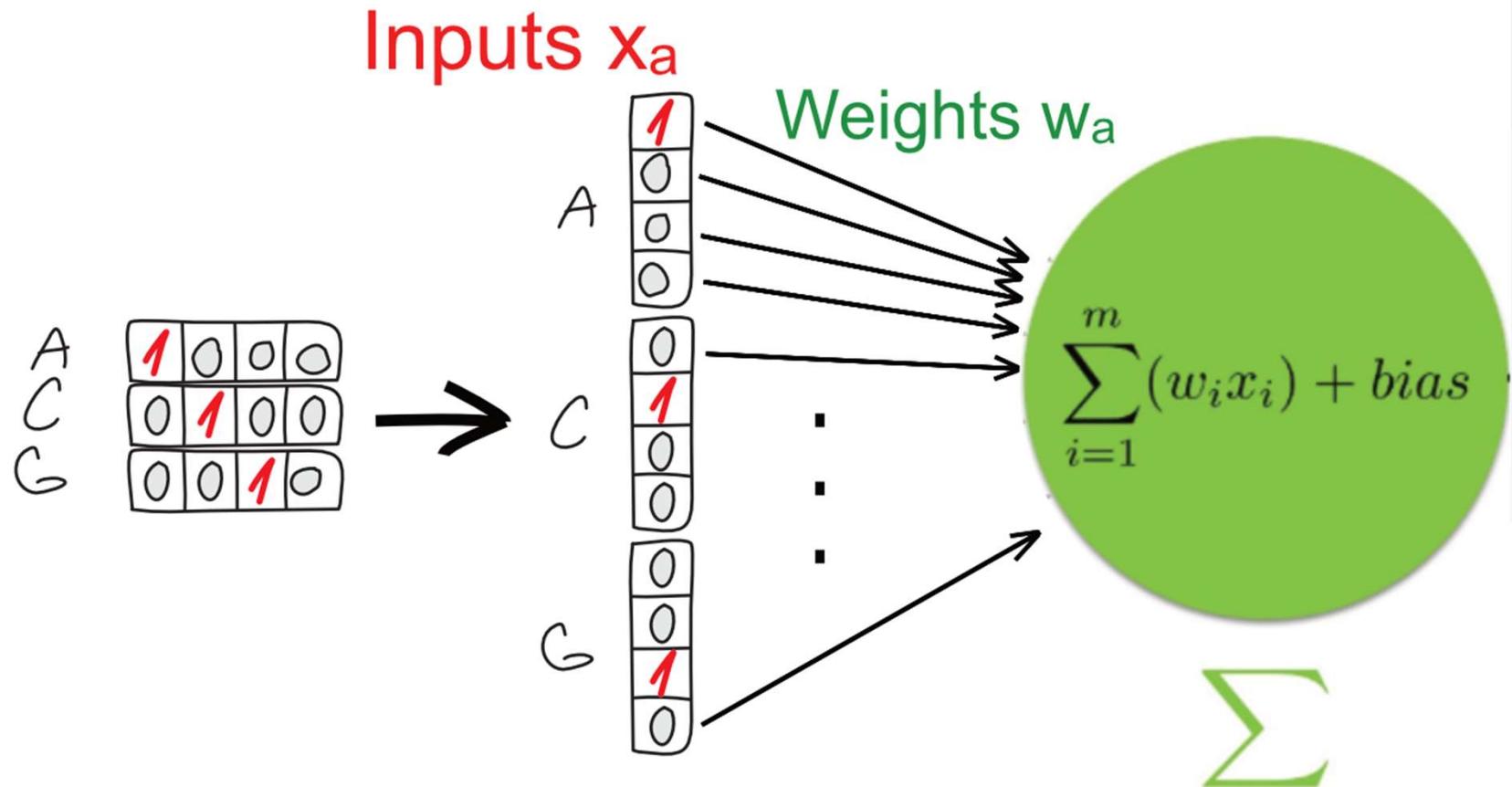
A	C	C	G	T	C	A	C	C	T	A
1	2	2	3	4	2	3	2	2	4	1

Numericalization & Tokenization

ACCGT CACCTACGT CGATTACGGGGTCATGACC

A	C	C	G	T	C	A	C	C	T	A
1	2	2	3	4	2	3	2	2	4	1

<table><tbody><tr><td>1</td></tr><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>0</td></tr></tbody></table>	1	0	0	0	<table><tbody><tr><td>0</td></tr><tr><td>1</td></tr><tr><td>0</td></tr><tr><td>0</td></tr></tbody></table>	0	1	0	0	<table><tbody><tr><td>0</td></tr><tr><td>1</td></tr><tr><td>0</td></tr><tr><td>0</td></tr></tbody></table>	0	1	0	0	<table><tbody><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>1</td></tr><tr><td>0</td></tr></tbody></table>	0	0	1	0	<table><tbody><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>1</td></tr></tbody></table>	0	0	0	1	<table><tbody><tr><td>0</td></tr><tr><td>1</td></tr><tr><td>0</td></tr><tr><td>0</td></tr></tbody></table>	0	1	0	0	<table><tbody><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>1</td></tr><tr><td>0</td></tr></tbody></table>	0	0	1	0	<table><tbody><tr><td>0</td></tr><tr><td>1</td></tr><tr><td>0</td></tr><tr><td>0</td></tr></tbody></table>	0	1	0	0	<table><tbody><tr><td>0</td></tr><tr><td>1</td></tr><tr><td>0</td></tr><tr><td>0</td></tr></tbody></table>	0	1	0	0	<table><tbody><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>1</td></tr></tbody></table>	0	0	0	1	<table><tbody><tr><td>1</td></tr><tr><td>0</td></tr><tr><td>0</td></tr><tr><td>0</td></tr></tbody></table>	1	0	0	0
1																																																						
0																																																						
0																																																						
0																																																						
0																																																						
1																																																						
0																																																						
0																																																						
0																																																						
1																																																						
0																																																						
0																																																						
0																																																						
0																																																						
1																																																						
0																																																						
0																																																						
0																																																						
0																																																						
1																																																						
0																																																						
1																																																						
0																																																						
0																																																						
0																																																						
0																																																						
1																																																						
0																																																						
0																																																						
1																																																						
0																																																						
0																																																						
0																																																						
1																																																						
0																																																						
0																																																						
0																																																						
0																																																						
0																																																						
1																																																						
1																																																						
0																																																						
0																																																						
0																																																						

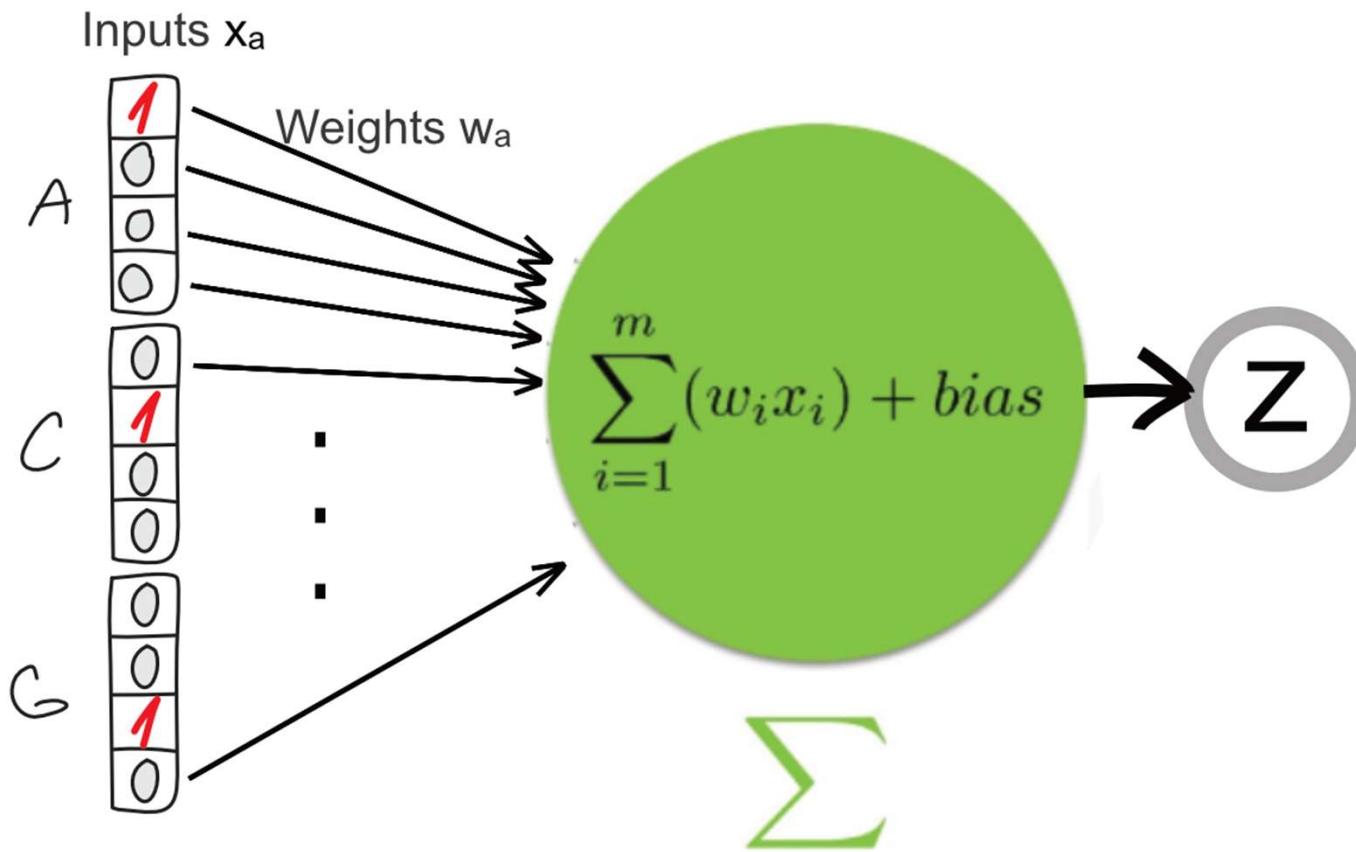


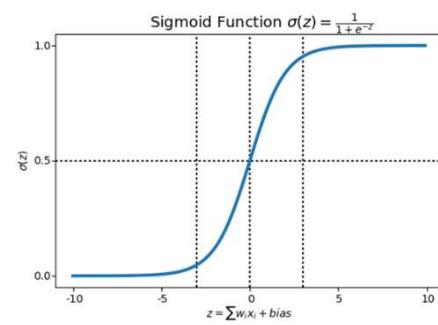
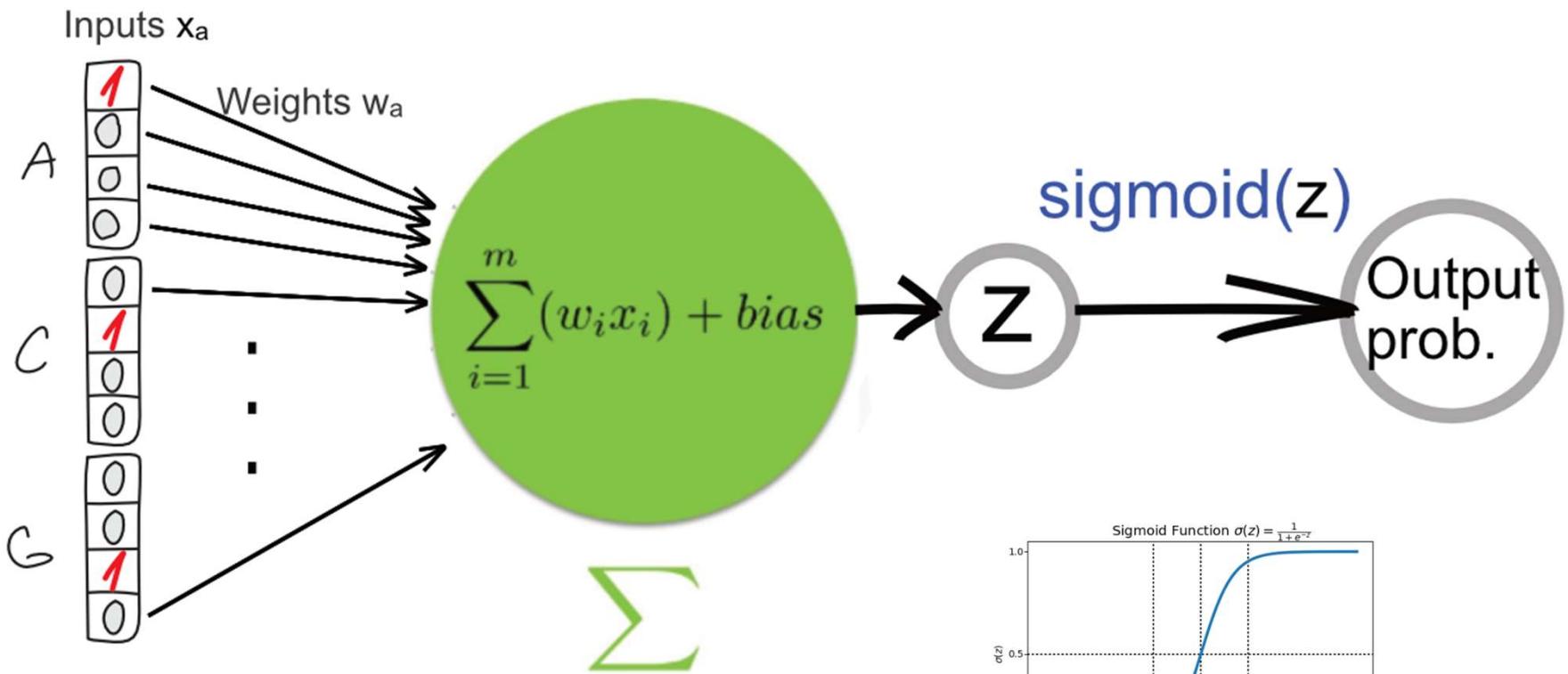
Logistic regression - DNA promoters classification

[https://github.com/BioGeMT/
MALTAomics-Summer-
School/](https://github.com/BioGeMT/MALTAomics-Summer-School/)

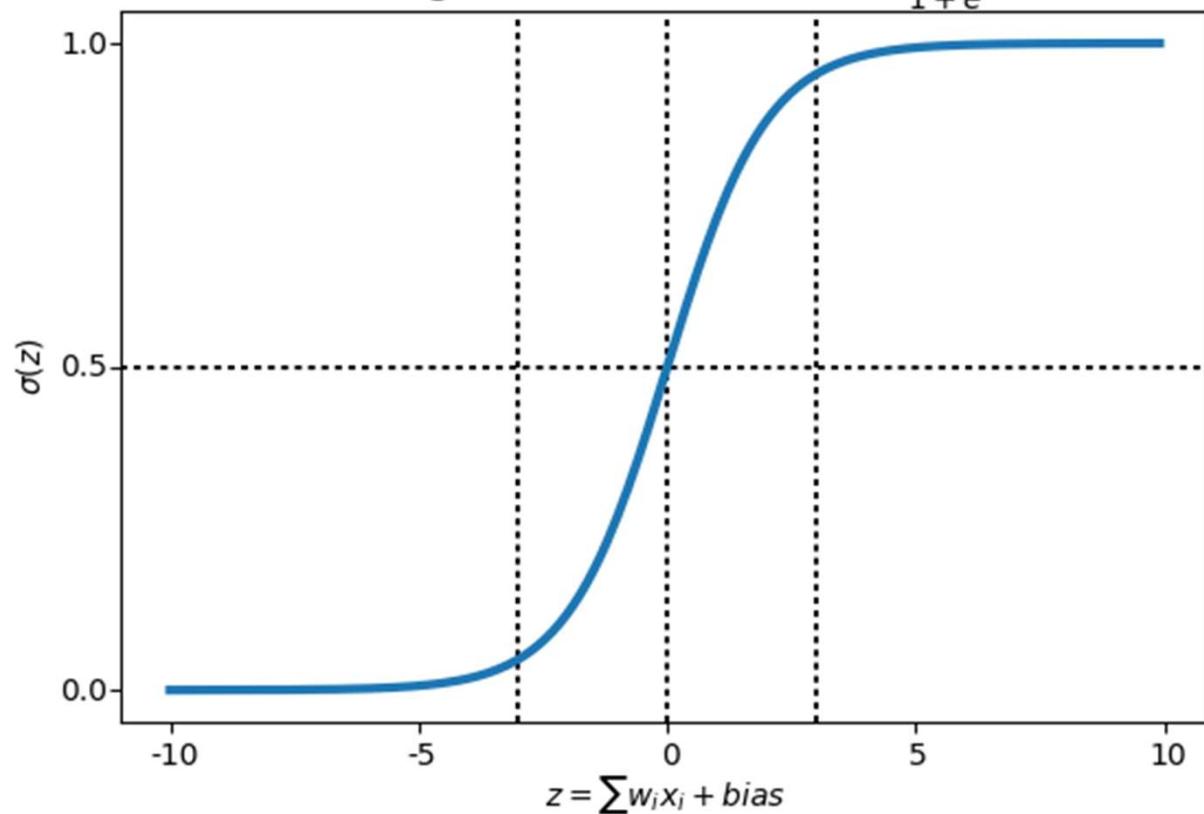
Day 2 Workshop

MALTA_02_DNA_enhancers.ipynb

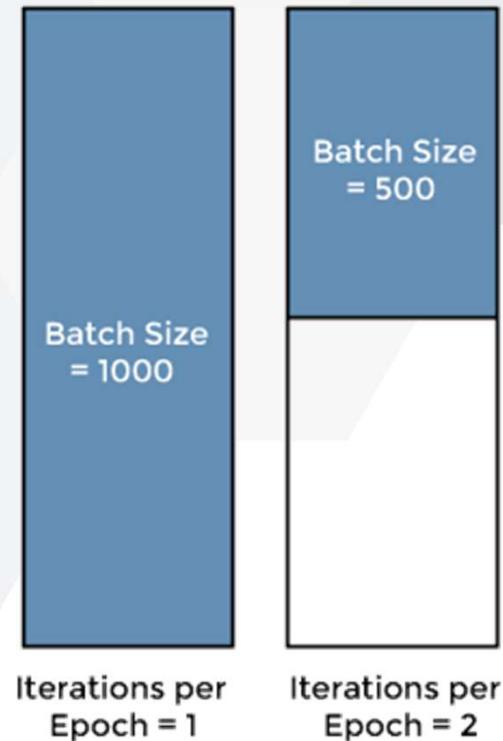




Sigmoid Function $\sigma(z) = \frac{1}{1+e^{-z}}$



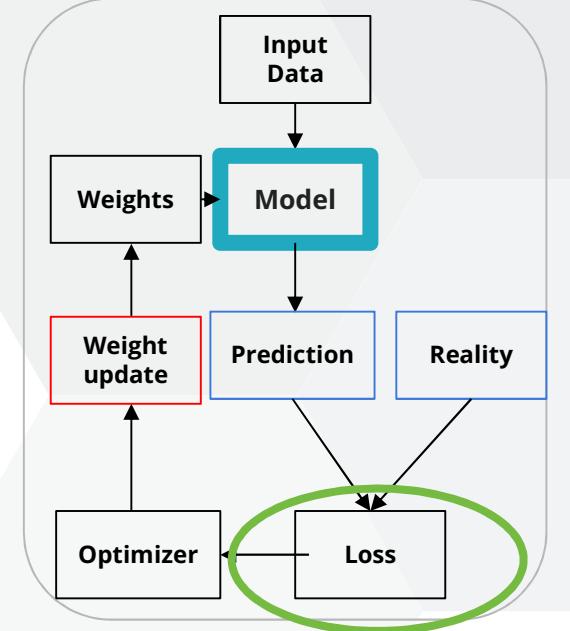
Batch size and shuffling



Batch size - speed VS stability of the learning process

Loss functions and activations

	Regression	Binary Classification
Loss function	MSE	Binary CE
Activation	None	Sigmoid

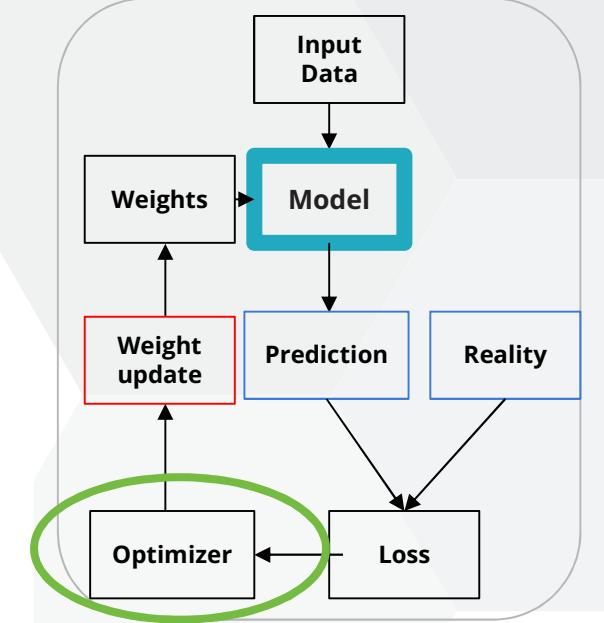
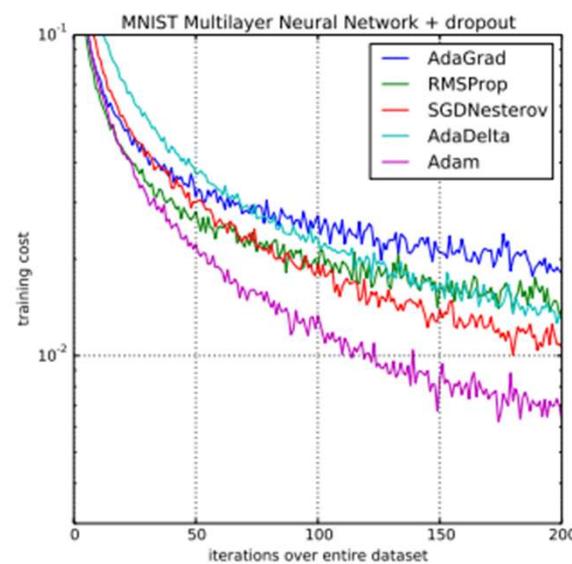


Adam optimizer

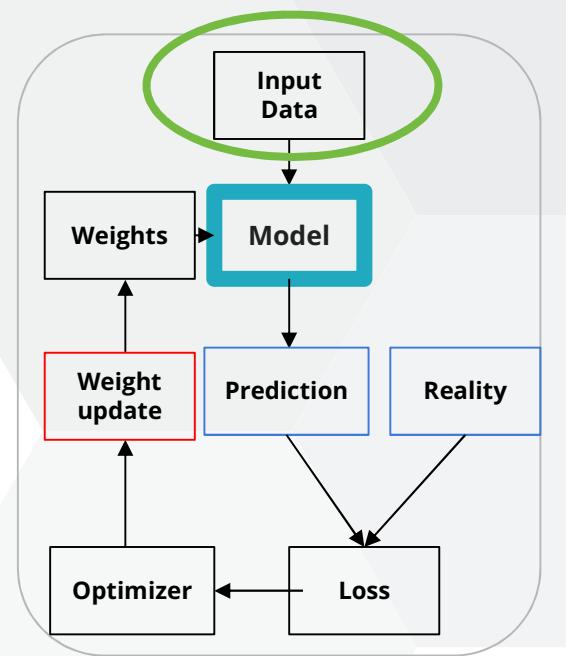
Converges faster than SGD

Uses momentum and moving average of the gradient

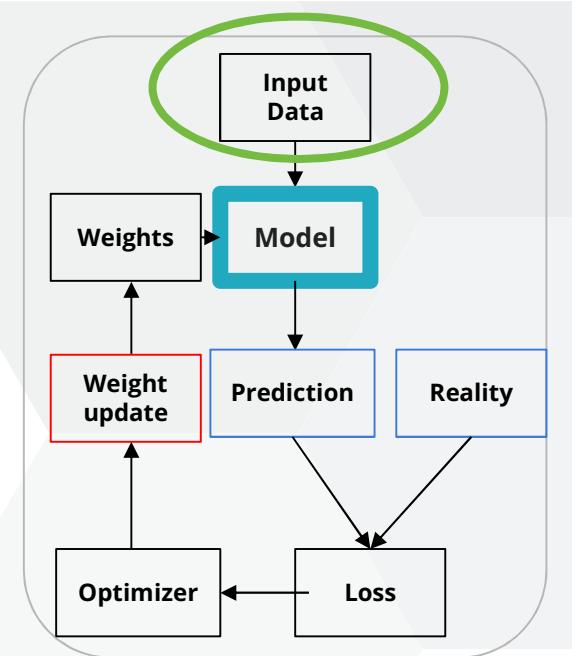
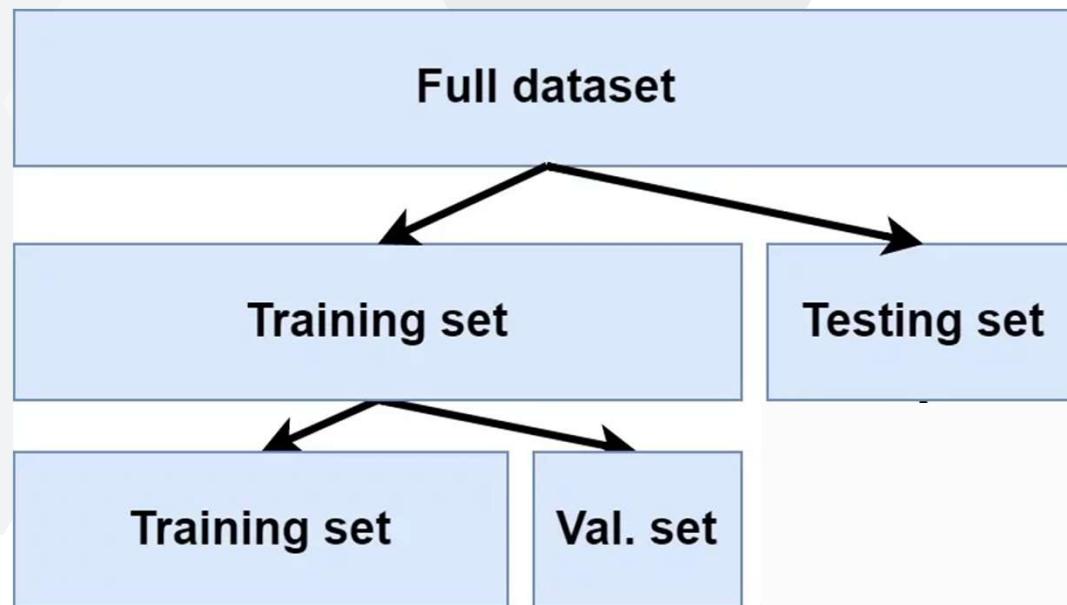
Computer per-weight learning rates



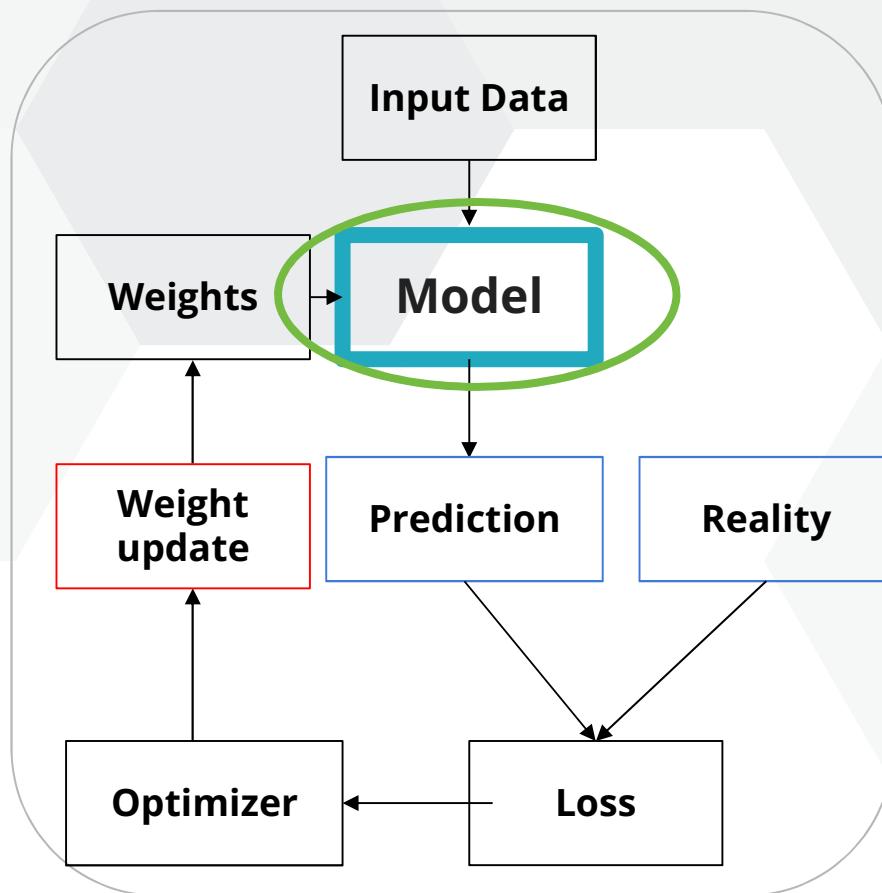
Dataset management



Dataset management



Multilayer Perceptron



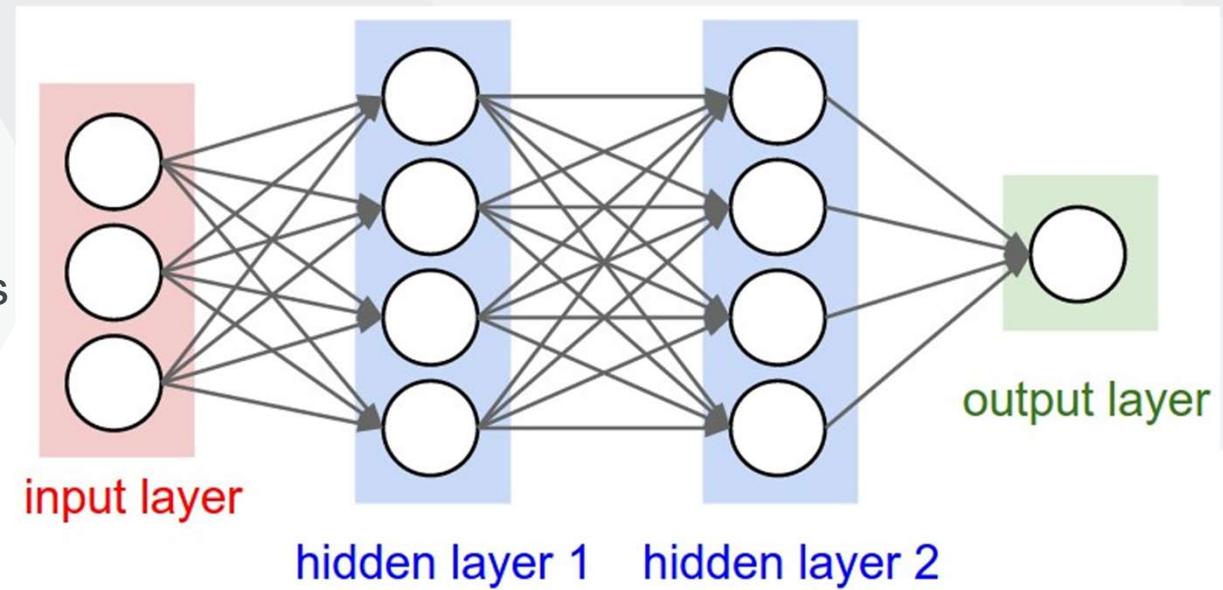
Nonlinearity and hidden layers

We could stack Linear Regression

No more powerful

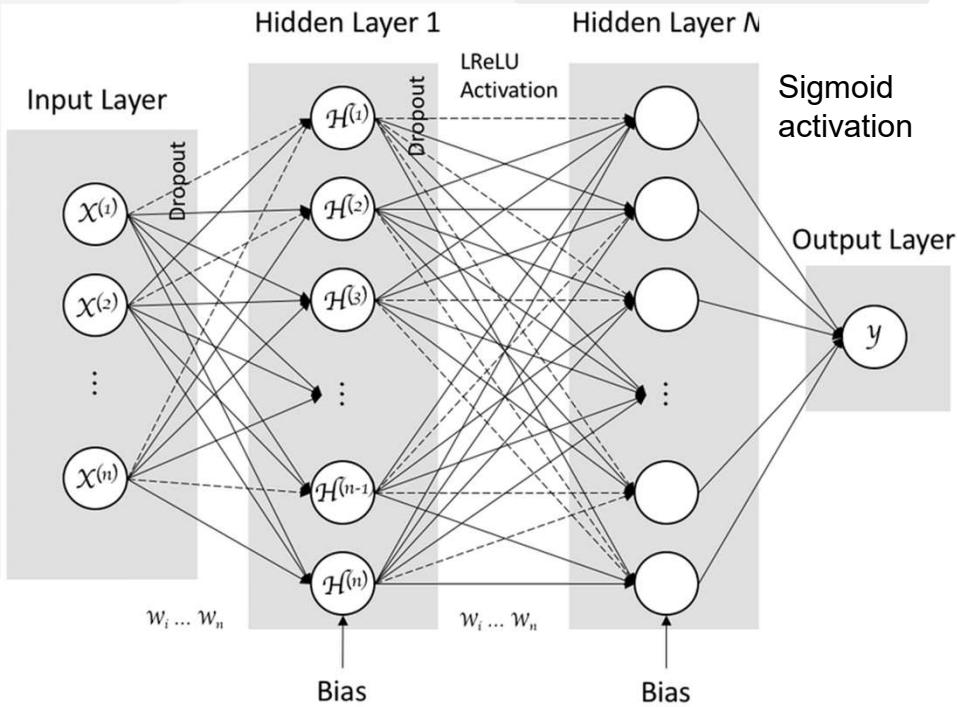
Only linear combinations

We need nonlinearity



Dense layer = every neuron is connected to all neuron from the previous layer

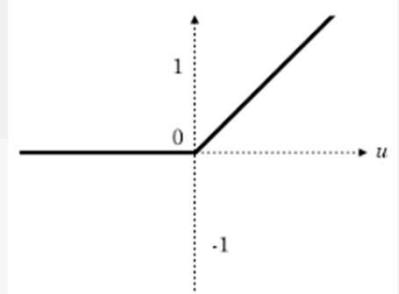
Multilayer Perceptron (MLP)



Hidden Layer

Activation:
ReLU

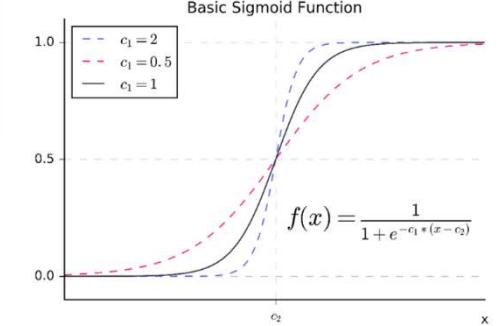
$$f(u) = \max(0, u)$$

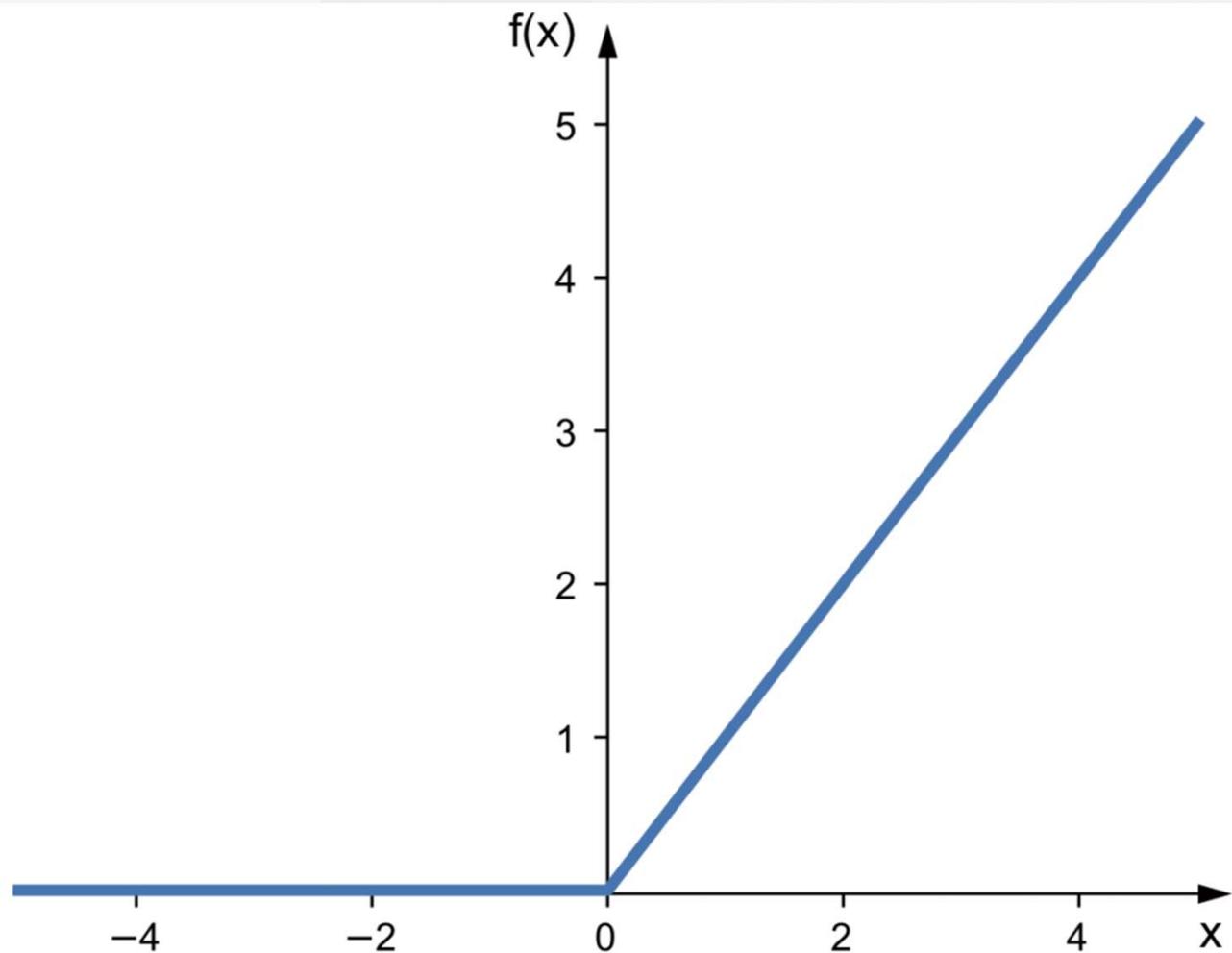


Output Layer

Activation:
Sigmoid

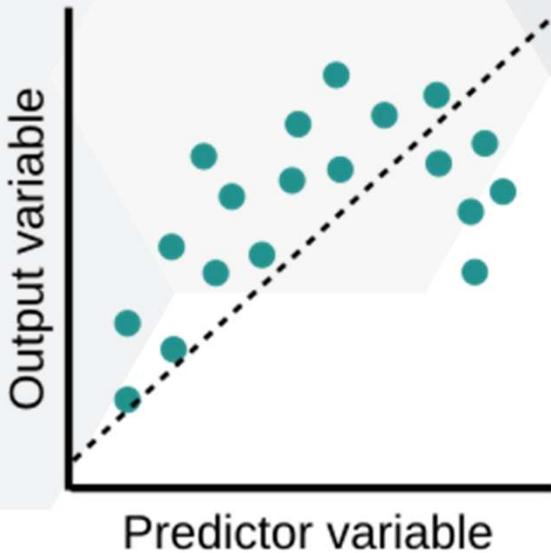
Basic Sigmoid Function



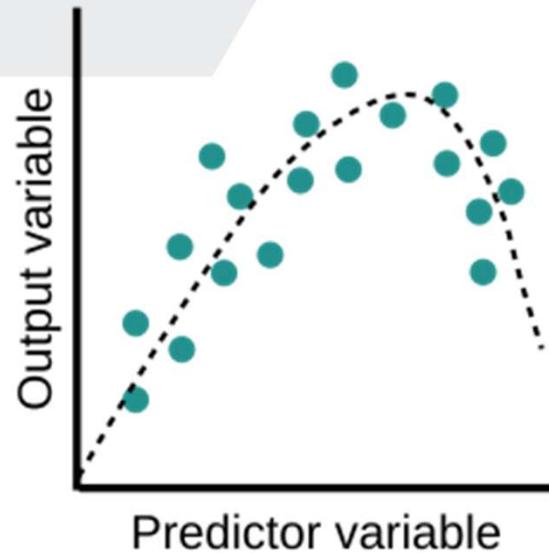


Overfitting

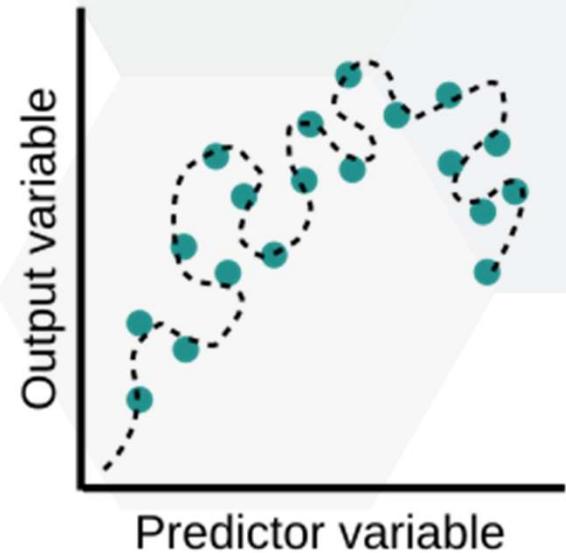
Underfit



Optimal



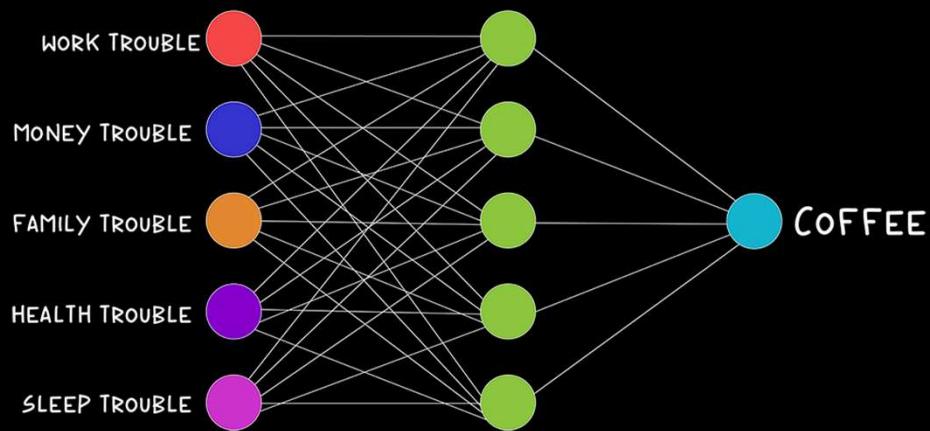
Overfit



Quiz

1. What are the benefits of shuffling the training data before creating batches?
2. What is the reason for adding nonlinearity?
3. What should be the output of an NN if we want to do classification?

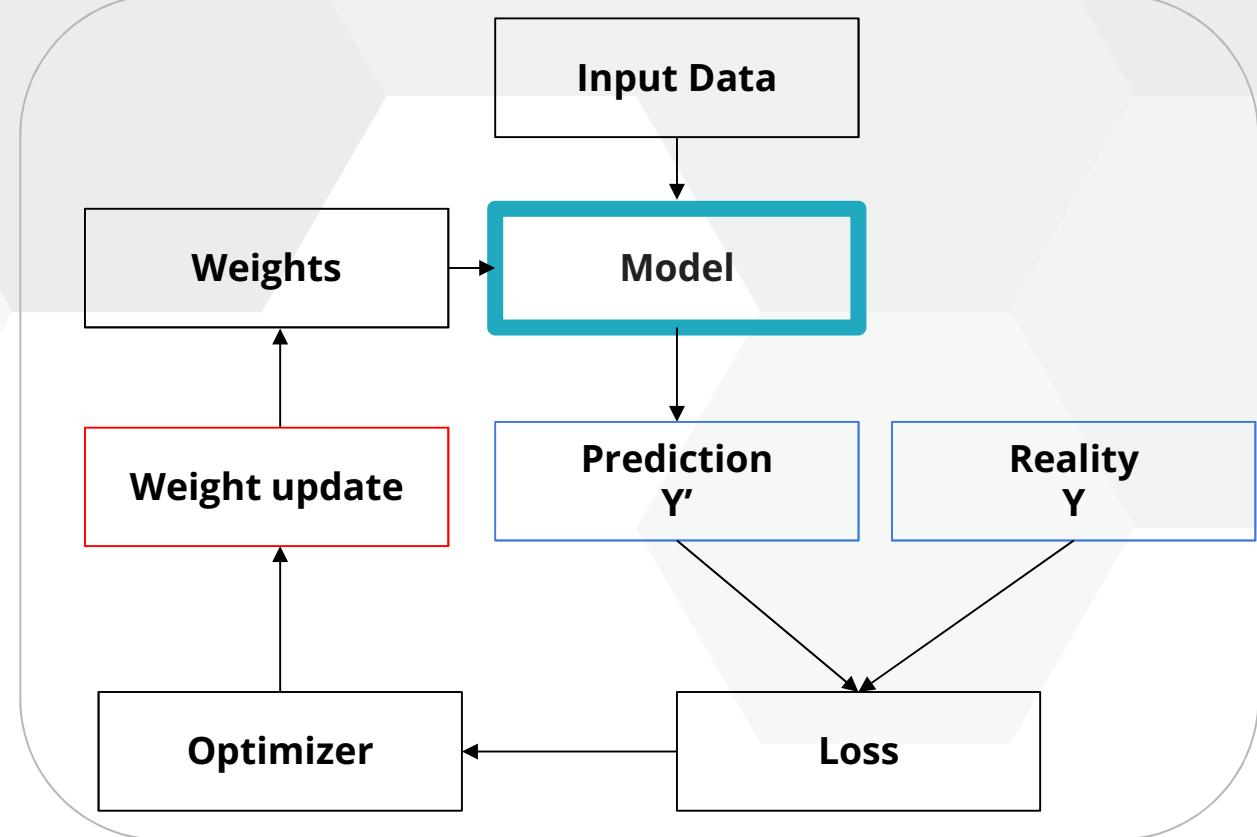
MY BRAIN'S NEURAL NETWOK



**Short break.
Back at**

Review

Linear regression
Gradient
Parameter
Hyperparameter
Epoch
Batch size
Learning rate
Regression
Classification
Loss function
Activation function
MLP



Break

Data scientists while their model is training



*A convolutional neural network taking a coffee break
Panos x DALLE2*

Convolutional Neural Networks

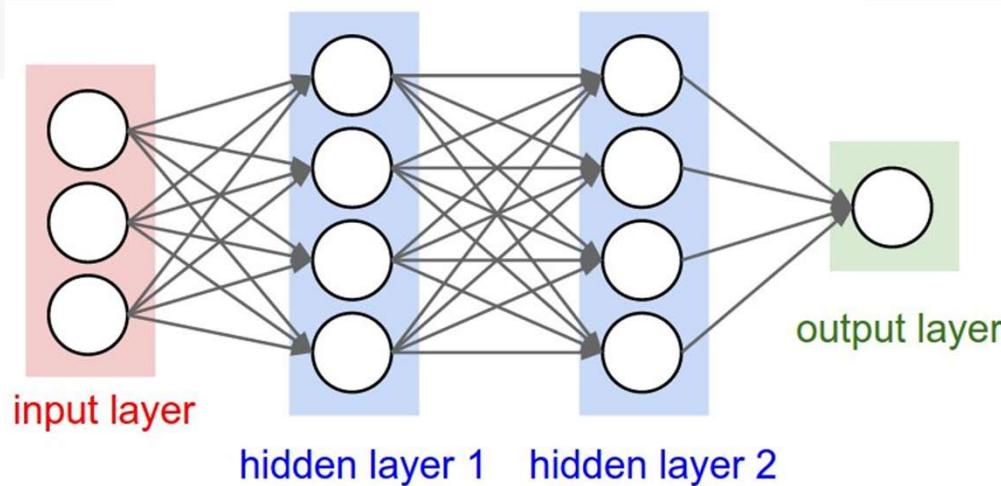
MLP drawbacks

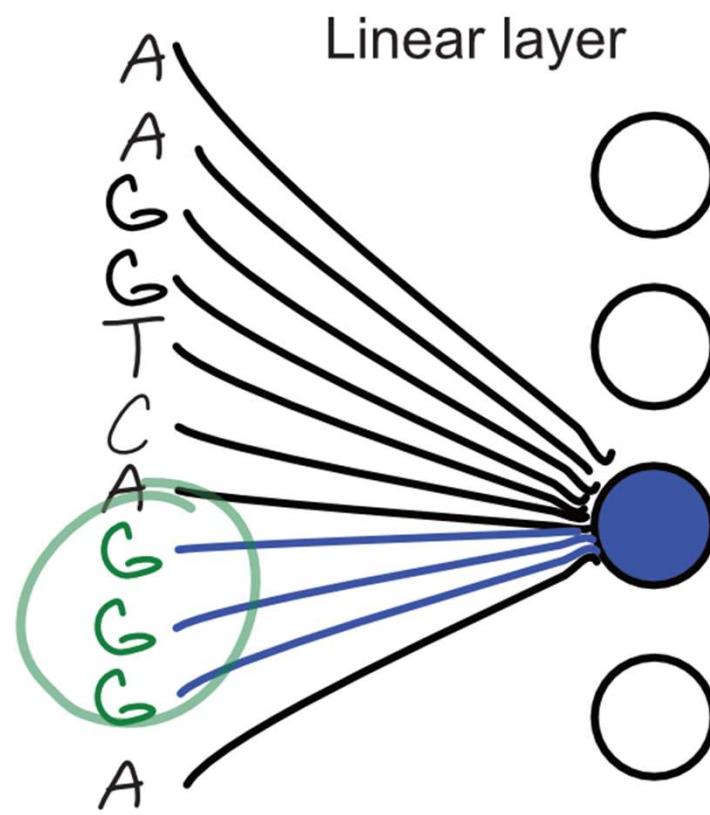
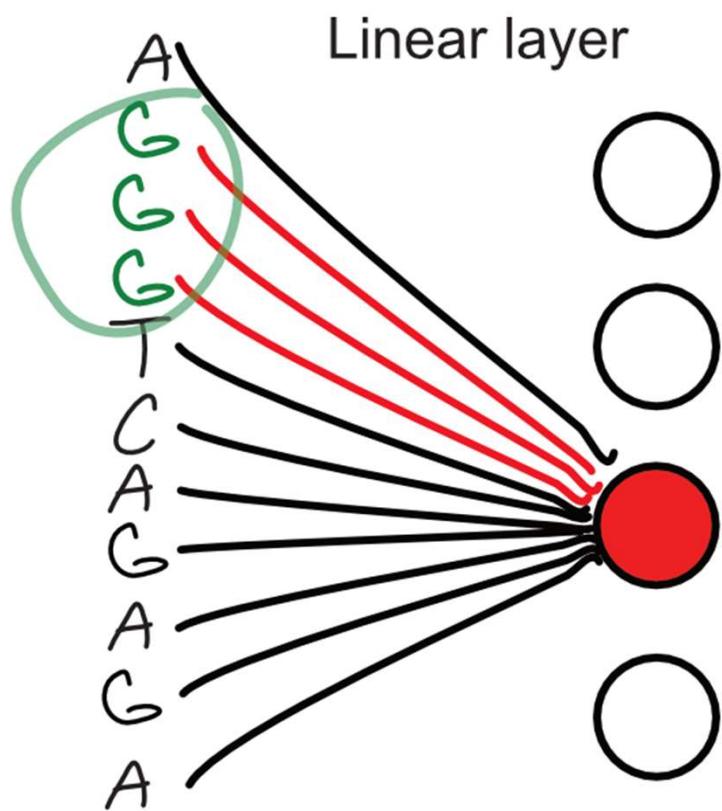
Weight for each value - too many parameters

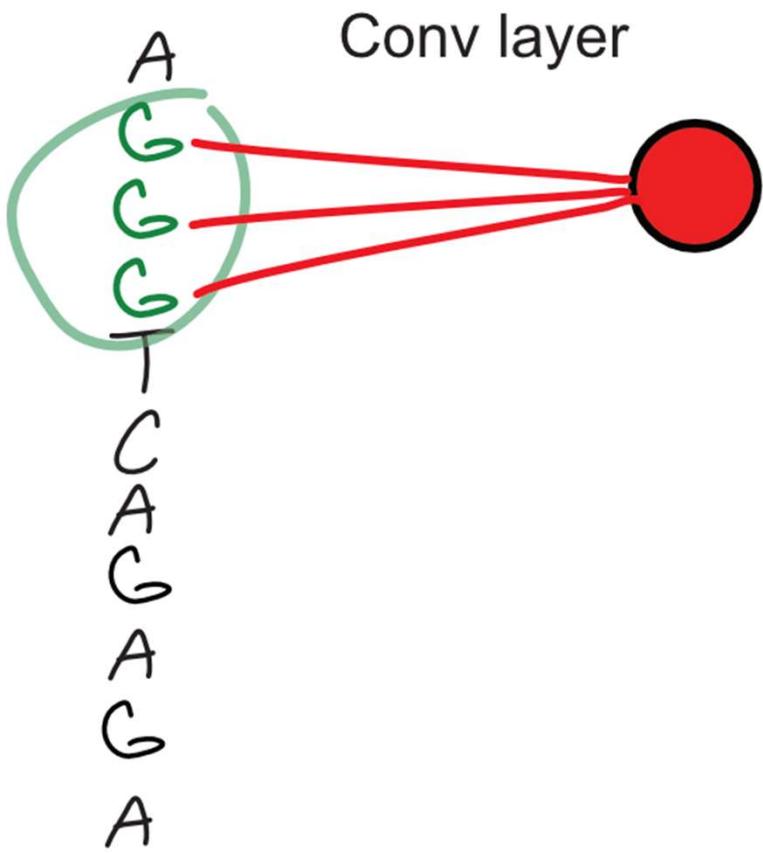
The **same pattern in different place** behaves differently

Treats the input as a row of values

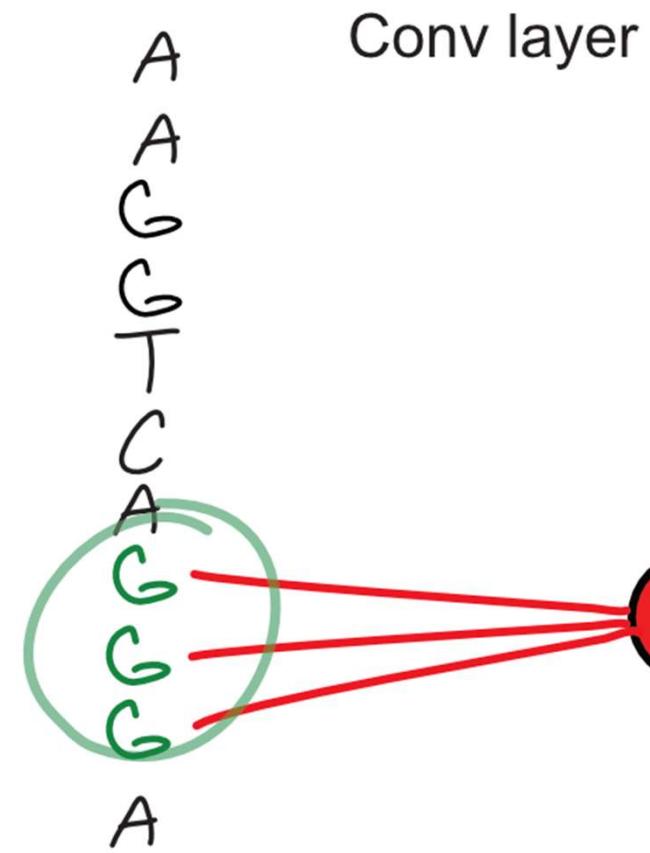
- No proximity, no sense of local neighbourhood







Conv layer



Conv layer

A
C
C
G
T
C
A
C
C
C
T
A

Search for
the pattern

C
C
C

A
C
C
G
T
C
A
C
C
T
A

* CCC

A
C
C
G

T
C
A
C
C
T

A

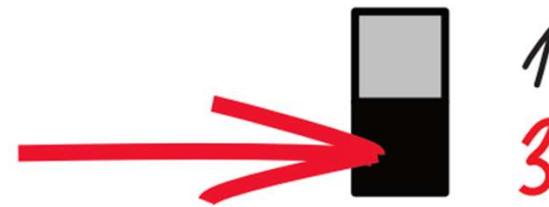
* C
C
G



1

A
C
C
G
T
C
A
C
C
T
A

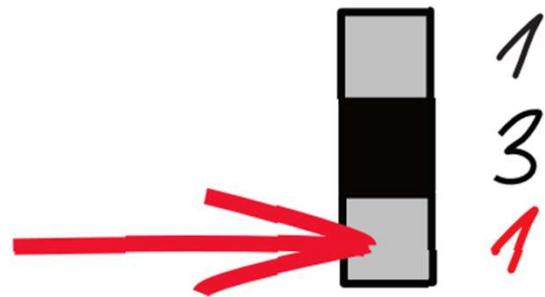
* C
C
G



1
3

A
C
C
G
T
C
A
C
C
T
A

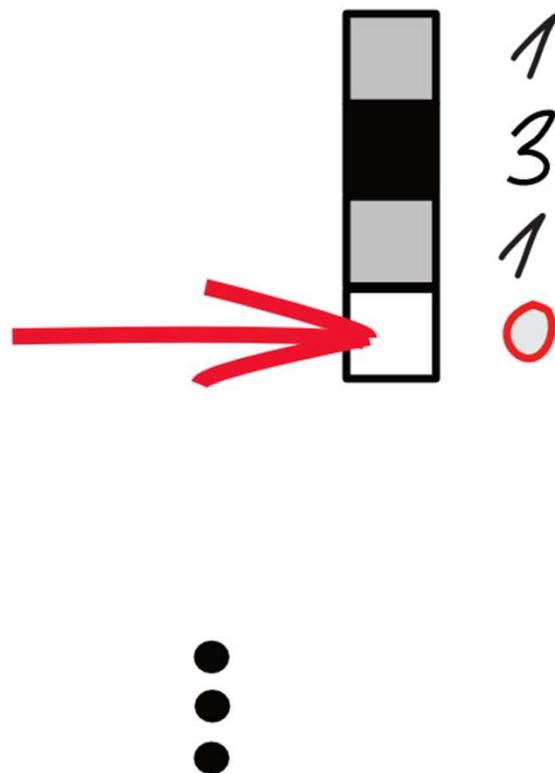
* C
C
G



1
3
1

A
C
C

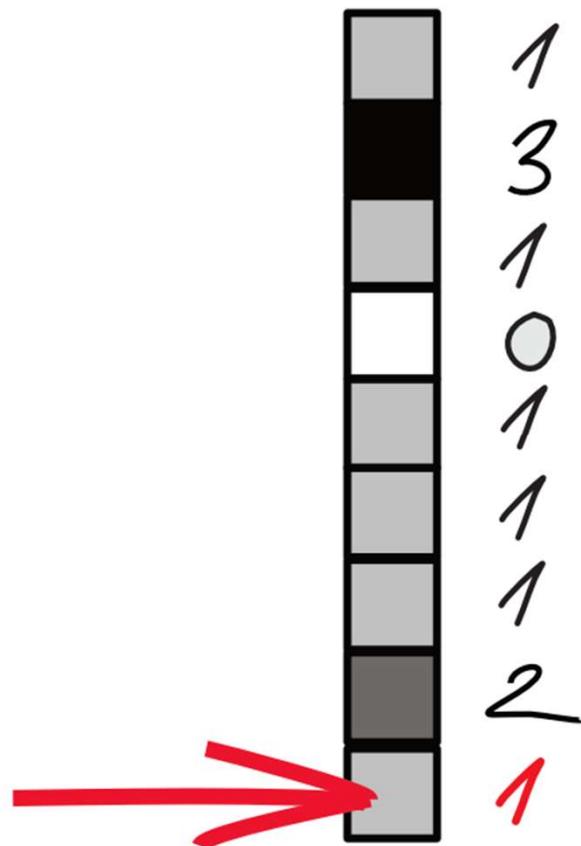
G T C A C C T A
T C G G
* CCC G



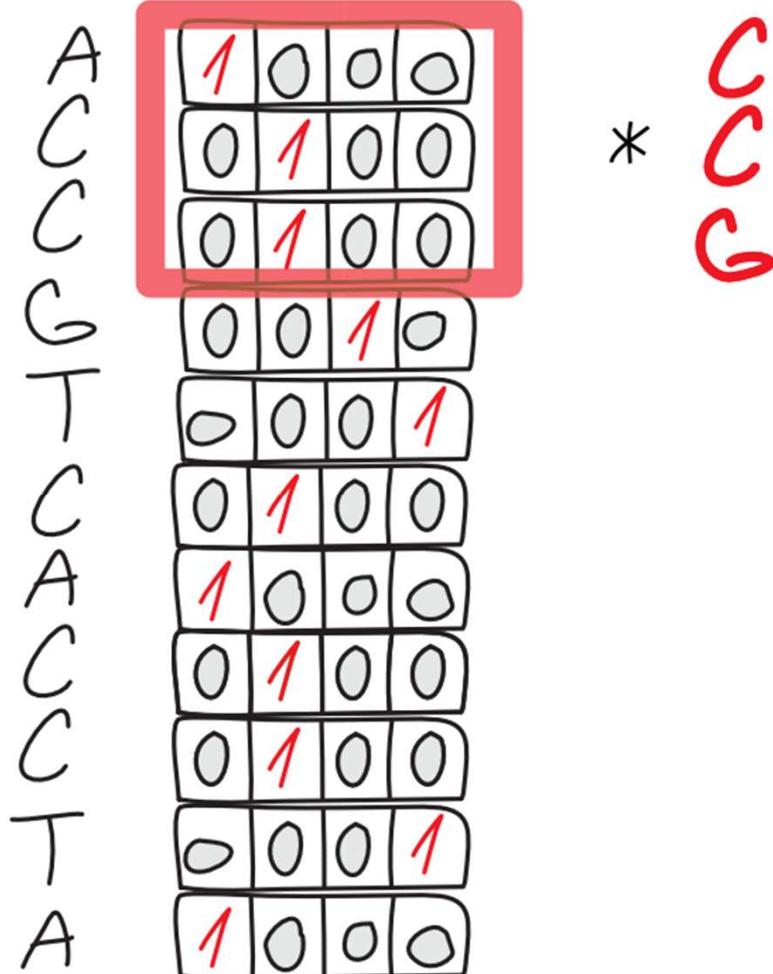
A
C
C
G
T
C
A
C

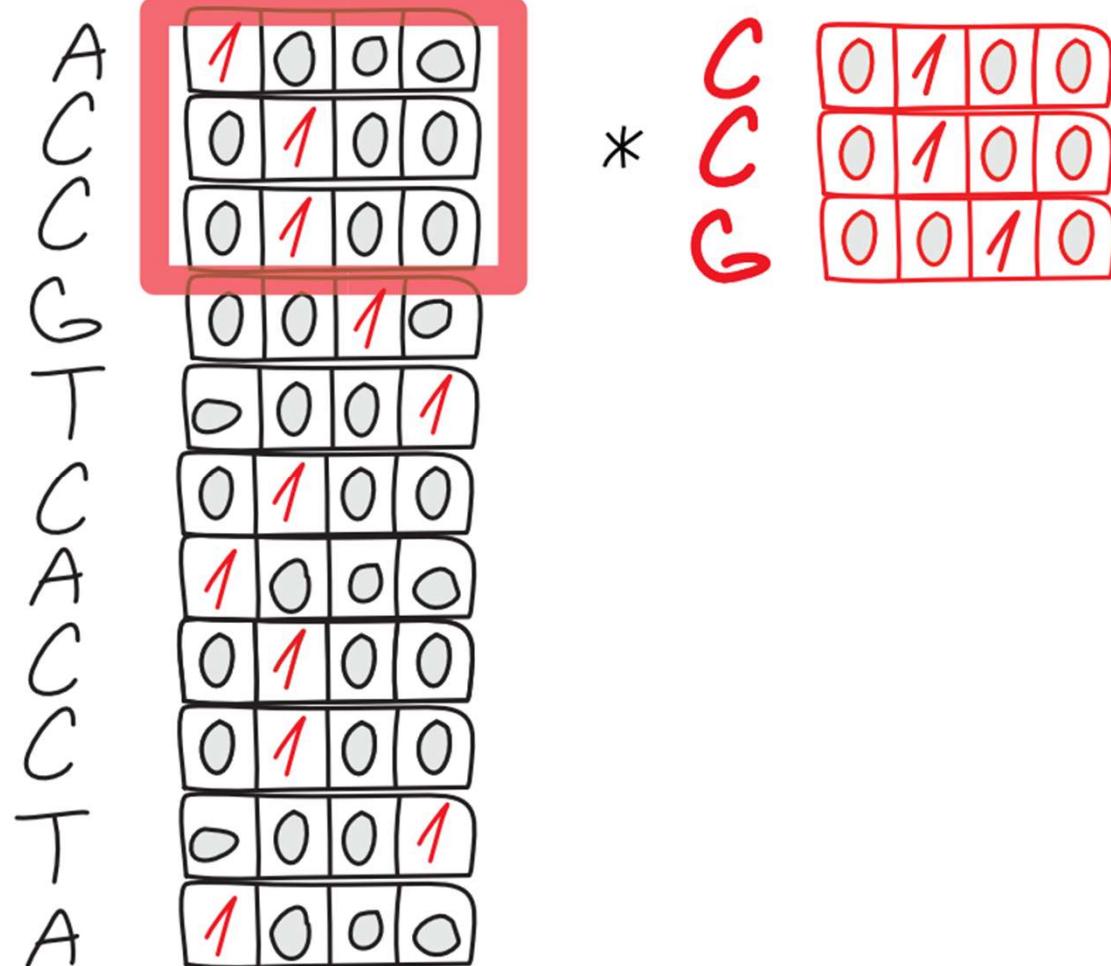
C
T
A

* C
C
G



A	1	0	0	0
C	0	1	0	0
C	0	1	0	0
G	0	0	1	0
T	0	0	0	1
C	0	1	0	0
A	1	0	0	0
C	0	1	0	0
C	0	1	0	0
T	0	0	0	1
A	1	0	0	0



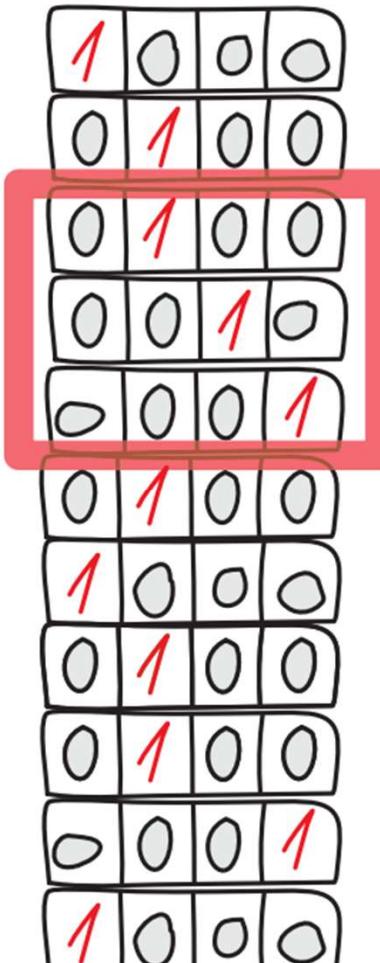


A
C
C
G
T
C
A
C
C
C
T
A

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

$$\begin{array}{l} * \quad \left. \begin{array}{l} 0 \\ 1 \\ 0 \\ 0 \end{array} \right\} = 0 \\ \left. \begin{array}{l} 0 \\ 1 \\ 0 \\ 0 \end{array} \right\} = 1 \\ \left. \begin{array}{l} 0 \\ 0 \\ 1 \\ 0 \end{array} \right\} = 0 \end{array} + 1 = 1$$

A	1	0	0	0
C	0	1	0	0
C	0	1	0	0
G	0	0	1	0
T	0	0	0	1
C	0	1	0	0
A	1	0	0	0
C	0	1	0	0
C	0	1	0	0
T	0	0	0	1
A	1	0	0	0

A 

 C

 C

 G

 T

 C

 A

 C

 C

 T

 A

$\begin{matrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{matrix}$

$*$

0	1	0	0	$= 1$
0	1	0	0	$= 0$
0	0	1	0	$= 0$

$= 1 + 0 + 0 = 1$

A	1	0	0	0
C	0	1	0	0
C	0	1	0	0
G	0	0	1	0
T	0	0	0	1
C	0	1	0	0
A	1	0	0	0
C	0	1	0	0
C	0	1	0	0
T	0	0	0	1
A	1	0	0	0

$$\begin{array}{l}
 \begin{array}{c}
 \begin{array}{ccccc}
 & 1 & 0 & 0 & 0 \\
 & 0 & 1 & 0 & 0 \\
 & 0 & 1 & 0 & 0 \\
 \text{G} & 0 & 0 & 1 & 0 \\
 \text{T} & 0 & 0 & 0 & 1 \\
 \text{C} & 0 & 1 & 0 & 0 \\
 \text{A} & 1 & 0 & 0 & 0 \\
 \text{C} & 0 & 1 & 0 & 0 \\
 \text{C} & 0 & 1 & 0 & 0 \\
 \text{T} & 0 & 0 & 0 & 1 \\
 \text{A} & 1 & 0 & 0 & 0
 \end{array}
 \\ \times \quad \begin{array}{ccccc}
 & 0 & 1 & 0 & 0 \\
 & 0 & 1 & 0 & 0 \\
 & 0 & 0 & 1 & 0
 \end{array}
 \end{array}
 \end{array}
 = 0 + 0 + 0 = 0$$

⋮

A
C
C
G
T
C
A
C
C
T
A

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

*

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

$$\begin{aligned} 0 &| 1 & 0 & 0 = 1 \\ 0 &| 1 & 0 & 0 = 0 \\ 0 &| 0 & 0 & 1 = 0 \end{aligned} \quad \left. \begin{aligned} 1 \\ 0 \\ 1 \end{aligned} \right\} = 1$$

1
3
1
0
1
1
1
1

A
C
C
G
T
C
A
C
C
T
A



Convolutional layer
with 1 filter - [C C G]

$$\begin{array}{c} 0 \ 1 \ 0 \ 0 \\ 0 \ 1 \ 0 \ 0 \\ 0 \ 0 \ 1 \ 0 \end{array} = - \quad } \\ \quad = - + \quad } = - \\ \quad = - \end{array}$$

A
C
C
G
T
C
A
C
C
C
T
A



Convolutional layer
with 1 filter - [C C G]

$$\begin{array}{c} 0 \boxed{1} 0 0 \\ 0 \boxed{1} 0 0 \\ 0 0 \boxed{1} 0 \end{array} = - \quad \} = - + -$$



1
3
1
0
1
1
1
1
2
1

A
C
C
G
T
C
A
C
C
T
A



Convolutional layer
with 1 filter - [C C G]

1 feature
map



1
3
1
0
1
1
1
2
1

A
C
C
G
T
C
A
C
C
C
T
A

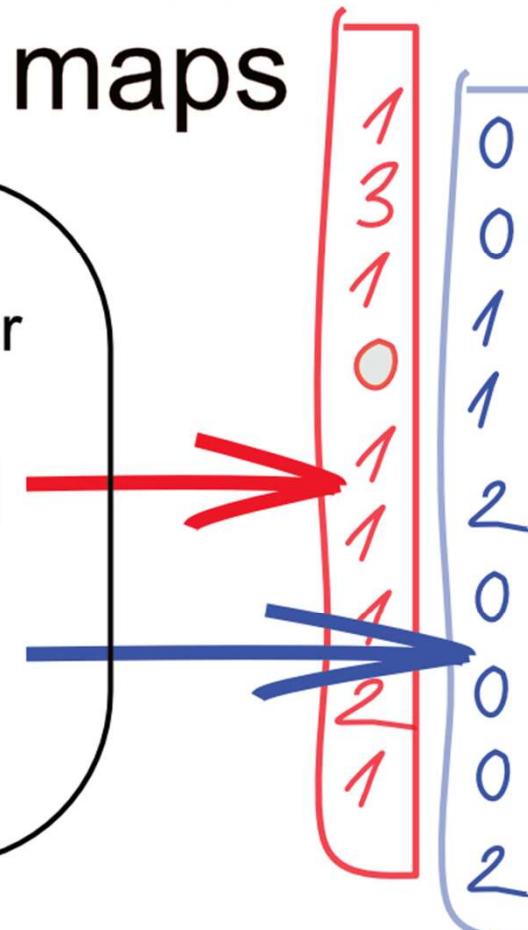


Convolutional layer
with 2 filters

1st filter [C C G]

2nd filter [T T A]

2 feature
maps



convolutional filters **learn their weights**

1st filter

w_1	w_2	w_3	w_4
w_5	w_6	w_7	w_8
w_9	w_{10}	w_{11}	w_{12}

2nd filter

w_1	w_2	w_3	w_4
w_5	w_6	w_7	w_8
w_9	w_{10}	w_{11}	w_{12}

convolutional filters learn their weights

w_1	w_2	w_3	w_4
w_5	w_6	w_7	w_8
w_9	w_{10}	w_{11}	w_{12}



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

C
C
G

w_1	w_2	w_3	w_4
w_5	w_6	w_7	w_8
w_9	w_{10}	w_{11}	w_{12}



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

T
T
A

convolutional filters **learn their weights**

DNA codons:

ATT
ATC
ATA

Amino Acid:

Isoleucine

w_1	w_2	w_3	w_4
w_5	w_6	w_7	w_8
w_9	w_{10}	w_{11}	w_{12}



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

A
T
A/C/T

Convolutional layer

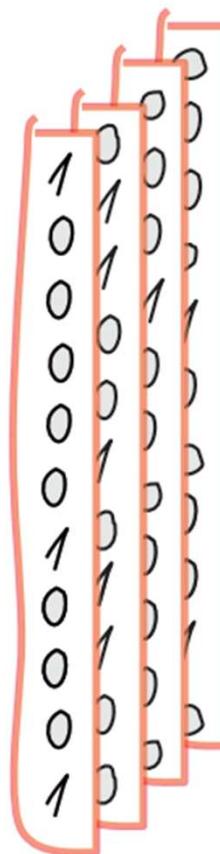
- Any number of filters
- Weights can be any number
 - 0.123
 - -0.003
 - Usually very close to 0 for “math reasons”

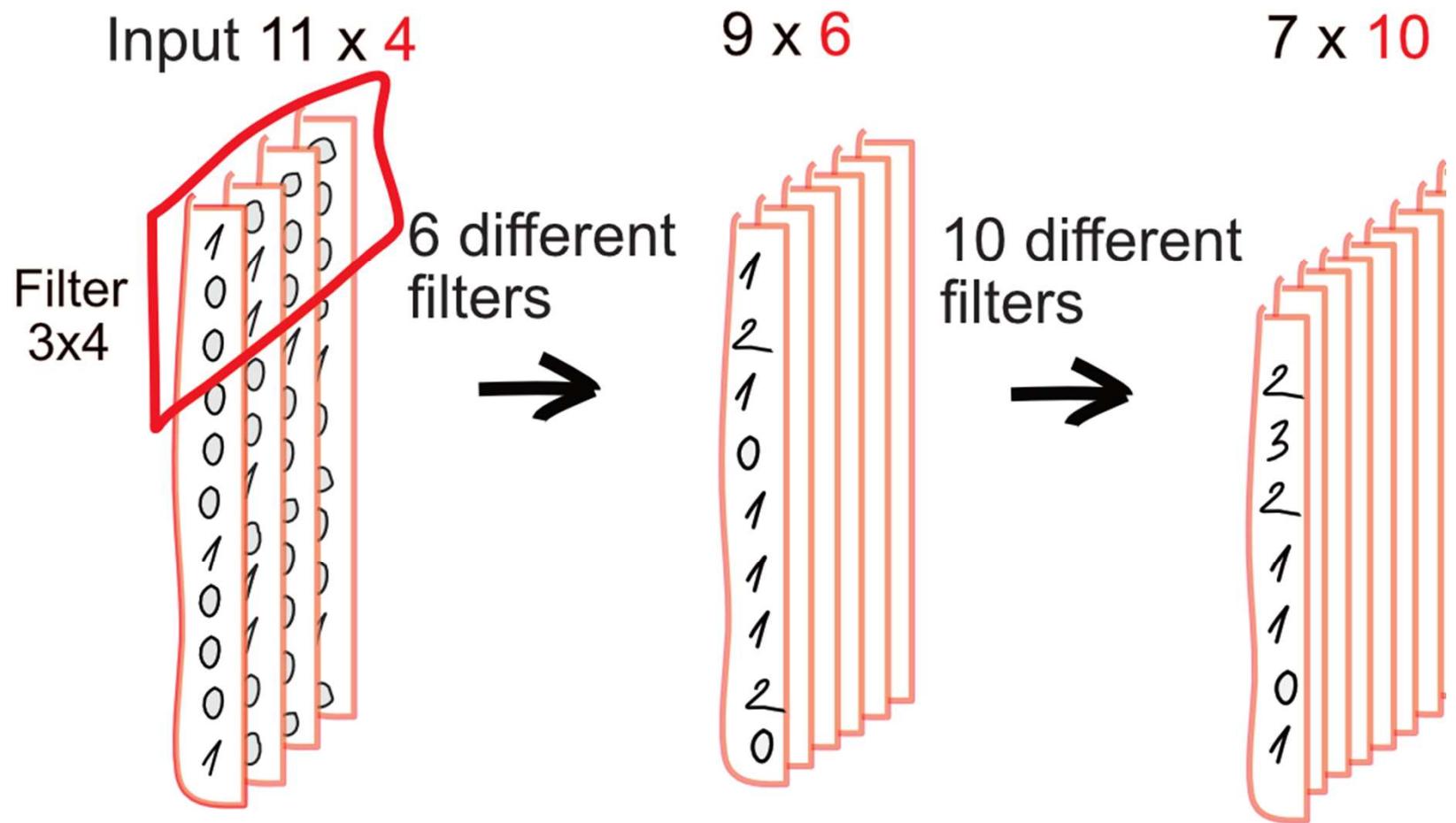
Input 11 x 4

A	1	0	0	0
C	0	1	0	0
C	0	1	0	0
G	0	0	1	0
T	0	0	0	1
C	0	1	0	0
A	1	0	0	0
C	0	1	0	0
C	0	1	0	0
T	0	0	0	1
A	1	0	0	0

4 channels
(feature maps)

Imagine as
→





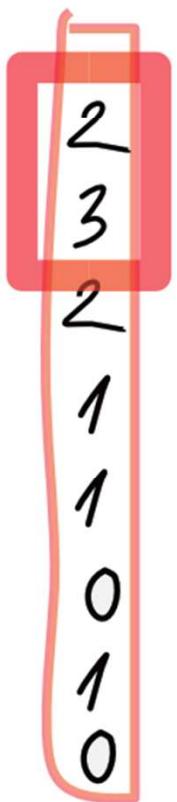
(Max) pooling layer

Decrease the spatial dimensions of the input (feature maps)

=> less computational steps

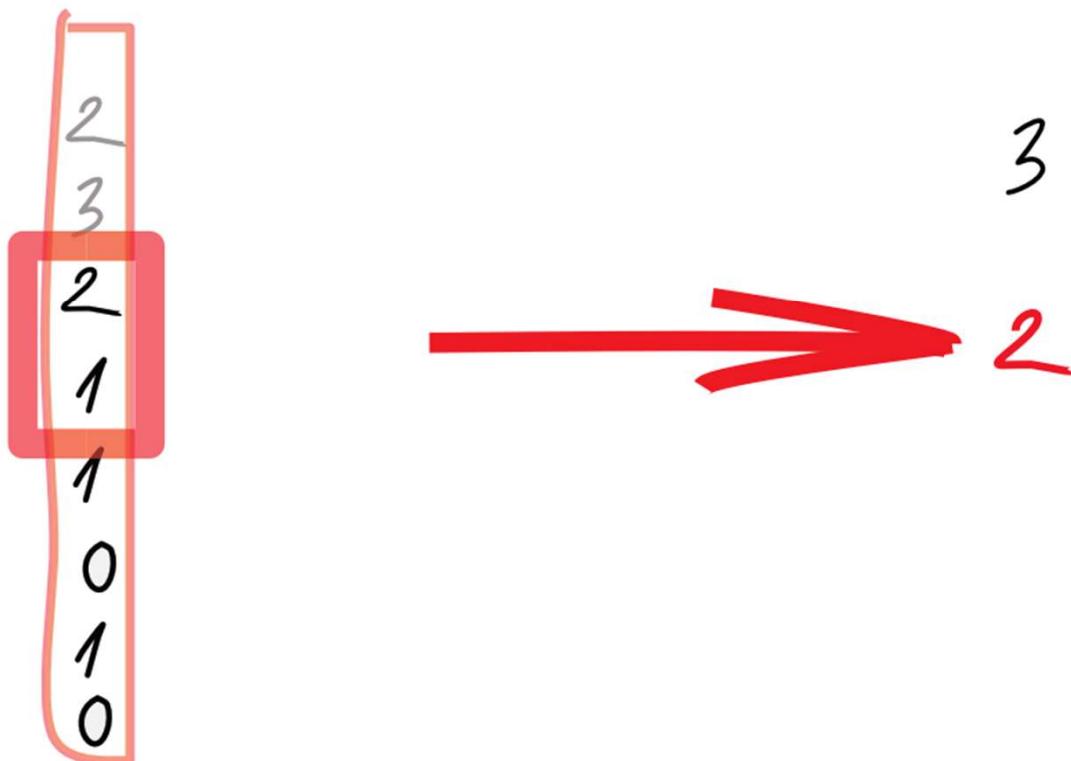
Translation invariance

Max pooling layer of size 2



$\Rightarrow 3$

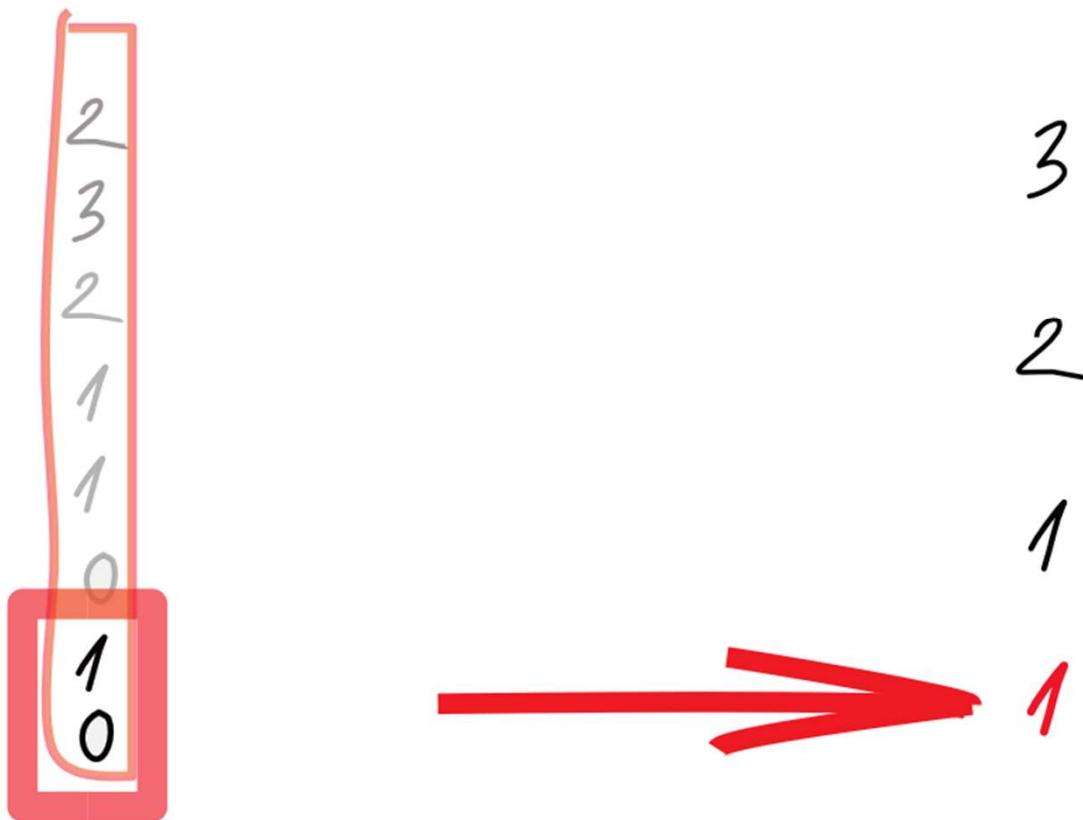
Max pooling layer of size 2



Max pooling layer of size 2

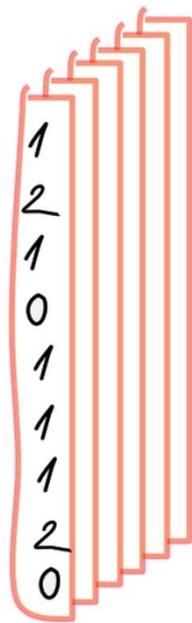


Max pooling layer of size 2



Flatten layer - spread the values

Length 9 x 6 feature maps



Flatten layer



length 54
- ready for
a linear layer



Convolutional NN - DNA enhancers classification

[https://github.com/BioGeMT/
MALTAomics-Summer-
School/](https://github.com/BioGeMT/MALTAomics-Summer-School/)

Day 2 Workshop

MALTA_02_DNA_enhancers.ipynb

Quiz CNNs

- How many filters can a convolutional layer have?
- Did we just work with a 2-dimensional or 1-dimensional convolution layer?
- What is a kernel?

The Team



Panos

Panagiotis
Alexiou



Petr

Petr
Simecek



Vlasta

Vlastimil
Martinek



David

David
Cechak



Katarina

Katarina
Gresova

martinekvlastimil95@gmail.com

Whatsapp +420 732 425 407

davidcechak@gmail.com





Thank you for your Attention

