

Learn Computer
Vision with CNN,
TensorFlow, and
PyTorch

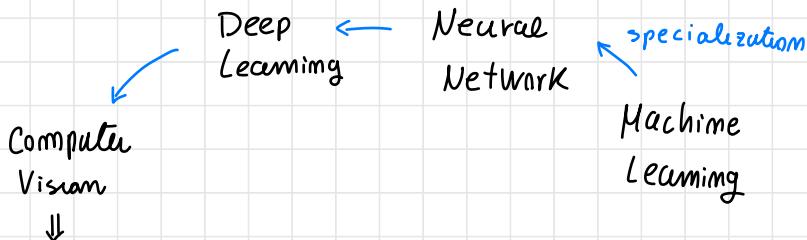
Master Object
Detection from
basics to advanced

Intro
+ Python
Set-up

COURSE INTRO

- Computer Vision is good way to learn programming + AI
↓
its Visual and Intuitive

- Computer Vision has constant demanding and improvements



↓
specialization of deep learning...

- from Python the course will focus on OpenCV, Deep Learning Theory + algorithms and applications.

- I {① Python (recap)
- ② Deep Learning Foundations (ANN from scratch, CNN, backpropagation...)
- ③ OpenCV (swiss army knife) → (image processing, transformation)
- ④ PyTorch (how models learn + optimize)
- II {⑤ visualize CNN's (architecture view)
- ⑥ Image Classification (build AI for object recognition + SOTA)
- ⑦ Data augmentation (flipping etc to enhance dataset)
- ⑧ Object Detection (find obj in images with YOLO + RCNN...)
- III {⑨ Image Segmentation (YOLO v11, other SOTA architecture)
- ⑩ Transformers (heart of NLP) and PROJECTS...

- Theory + Hands-on (PyTorch / Keras)

- SOTA Frameworks + Architectures

PYTHON PREREQUISITES (RECAP...)

Amacoda Installation

Open-source tool for sourcing, building and deploying DL tools

- best tool to work with multiple packages and environments

↑

easy management of packages + environment

fundamental to work with many Python versions... envs.

• Installation

On windows, remember to check "Add Amacoda3 to my PATH..."

(important to create environment variables to properly run `conda` commands)

↓

you have `Amacoda Navigator GUI`

providing many frameworks ready to use!

+ many envs created by default

+ Install VSCode to set-up programming IDE

Using VSCode

IDE that allow easily to create envs and work with python code

1) Create New environments... → using python 3.12 (Newest version)

↓

necessary in new projects to properly manage multiple libraries...

New features gets included always and maintenance to avoid conflict is required...

`conda create -p <env_name> python==<version>`

2) activate environment: `conda activate <env_name>`

3) install ipykernel to use jupyter notebooks in VS Code (with pip)

provide KERNEL to Notebook to run python code...

- 4) create a **requirements.txt** ... keep a note on required packages,
↳ responsible for future deployment
in future this will be
responsible to install all requirements

Python Basics, Syntax

basics of **Syntax**: set of rules that defines symbols to program
= correct code arrangement

Syntax errors are related
to variables bad usage or other

SEMANTIC: meaning / interpretation → how to manage to do a specific task... what it is supposed to do!

- Python delimit block of code with
Indentation → required for code structure and any operation of flow execution
(differently from other programming languages that rely on {} brackets)

...

coding practice with Jupyter Notebook on the main python concepts

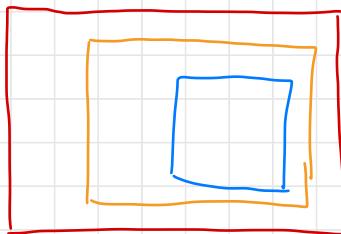
Data types, Data structures, functions, control flow, exception handling,
modules (libraries), object oriented Programming, generators, iterators,
magic methods/decorators...
+ Numpy, Pandas and Logging

(Notes on the Jupyter Notebooks `python_recap`)

Deep
Learning
Theory

INTRODUCTION TO DEEP LEARNING

brief idea of what is DL, and how it is related to AI and ML



Artificial Intelligence (AI): create app that performs their own tasks (such as Netflix recom., self-driving car, etc.)

Machine Learning (ML)

subset of AI that provides tools to analyze, visualize data + forecast and prediction

↳ supervised
↳ unsupervised

Deep Learning (DL)

subset of ML ..

main task is to mimic the human brain → make a machine learn as humans.



Neural Network (multilayered) (focus of this course)

↳ supervised

↳ unsupervised

Deep Learning

1) Artificial Neural Networks (ANNs) ↳ classification ↳ regression ↳ main problems in supervised ML

2) Convolutional Neural Networks (CNNs) → classification

Input data: images/video frame

↳ CNN are optimal with this inputs

+ RCNN, Masked RCNN,

help in Object Detection / segmentation

YOLOvX

↳ computer Vision: object detection

3) Recurrent Neural Network (RNN) → for NLP (Natural Language Processing) tasks

Input data: text/time-series data (forecast etc..)
(sequence of data with dependency)

+ Word Embedding, LSTM RNN, GRU RNN, Bidirectional LSTM RNN, Encoder/Decoder

+ Transformers, BERT... multiple use cases

FRAMEWORKS: • TensorFlow (made as open source by Google brain)
in end-to-end projects.

- big research is going on in NLP
- foundations are in ANN, CNN, RNN

Why DL is becoming so popular?

back in 2005... Facebook was launched ... than youtube, WhatsApp
... 2012... ↑ many online posts



social media platform
that created huge
amount of DATA

DATA generated EXPONENTIALLY

how to store those Data efficiently?

↓ { images, videos, texts... }

to be easily accessible

BIG DATA (MapReduce, Database...) → structured data
→ unstructured data



2011 Big data frameworks
created. store data and access efficiently

2011-2012 DL starts to become more popular

because of

1) HW requirements. GPU's cost decreasing
fundamental for model training



Train NVIDIA GPU price decrease + Cloud GPU
models possible

2) Huge amount of Data generated → DL models perform well
with huge data.

Using data in a smart way? → make product more similar to users

e.g. Netflix movie streaming platform

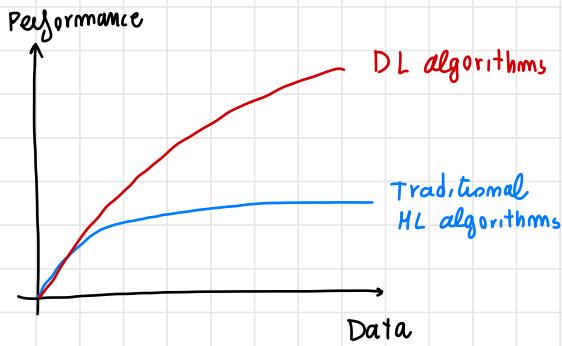
- how to keep users engaged? → based on Data available

↑

based on what I watch, ← Recommendation system
it suggest similar product.

my interaction with the platform is stored, and it is used to
create efficient models ⇒ increase review
↑

BIGGER DATA = more performative DL models



As Data increases, ML reach a limit, while DL continue to
increase accuracy performance with amount of Data

3) DL has been used in many domains

↓

- 1) Medical
- 2) E-commerce
- 3) Retail
- 4) Marketing
- ...

extensive application
thanks to the different
techniques / models
architecture

4) Open-source Frameworks for DL

↓

TensorFlow

PyTorch

TensorFlow open-sourced by Google brain team }
PyTorch open-sourced by Facebook }
} CNN, ANN, RNN ..
implementation

↓
community size increase! \Rightarrow more research in DL thanks to
an active community

↑
Using Big Data with this new algorithms to make more
review. Automating tasks impact fully

... Next \rightarrow how DL works in a Neural Network model

DEEP LEARNING CONCEPTS

{ ANN, LOSS FUNCTIONS, }
{ ACTIVATION FUNC., CNN }

PERCEPTRON

ARTIFICIAL NEURON / NN UNIT ← this is the basic unit of ANN

~ single Layered NN

What are its PROs/CONS ... Why we move to

MLP (Multi Layered Perceptron)
MLNN, ANN

- Input Layer
- Hidden Layer
- Weights
- Activation Function

Example **Perceptron** (single layer NN)

can be used for solving BINARY CLASSIFIER
to solve dataset problem of binary classification

DATASET

IQ	N° of study hours	% Pass/Fail
95	3	0
110	4	1
100	5	1

three rows of DATASET
how to use it in a
PERCEPTRON and train it ?

PERCEPTRON has the
following architecture

↳ ...

DATASET

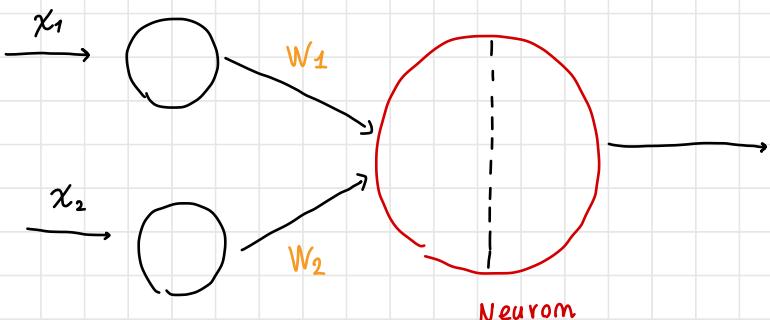
x_1 IQ	x_2 Nº of study hours	y % pass/Fail
95	3	0
110	4	1
100	5	1

INPUT LAYER
(same as input var.)

HIDDEN LAYER 1

simple NN Unit. Just one Neuron.
(single layer NN)

OUTPUT



you don't
count I layer
as separate
layer

[it performs two
tasks and return
an output]

↑
When input goes through Neuron, it is connected
using weights. (w_1, w_2)

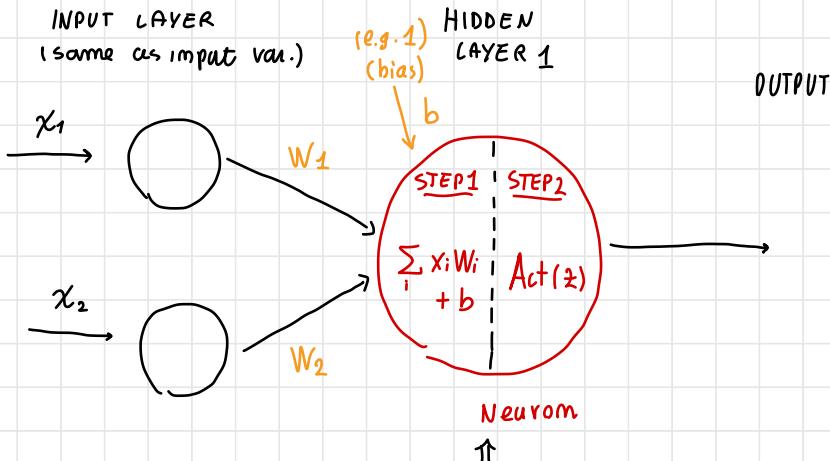
similalry to how our brain neurons process signals from environment and we respond... Neurons process information from data
• IF something important occurs there are HIGH WEIGHTS on some part of the processing pipeline...

↓

and our Neurons are well trained to handle data and responds
Weights are high if critical things occurs.

↑

ENTIRE NN mechanism consists in a good training of the weights w_i
to identify data information, possible by TRAINING!



this hidden layer consists of some hidden NEURON,
where two simple processing occurs to give a response (output)

✓
STEP 1

summation of
inputs weighted + Bias

$$z = w_1 x_1 + w_2 x_2 + b$$

$$= \sum w_i x_i + b$$

(I can have more
inputs x_i = features)

✓
STEP 2

On top of z computed we apply
an activation function

↑

it is the main aim of transforming
the value.

Transforms the output z between
some values, such as $0 \div 1$, $-1 \div 1$

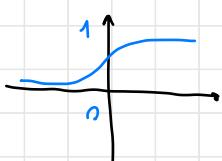
$0 \div 1$ I use SIGMOID Activation
or a STEP Activation

↓

STEP function



SIGMOID Function

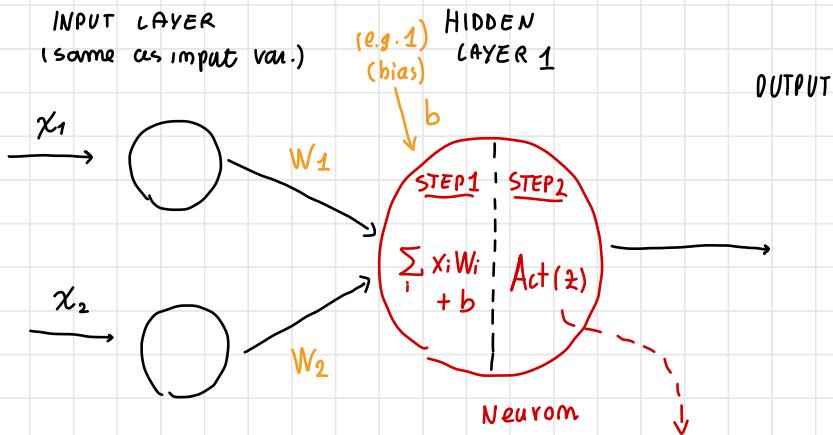


BIAS is important in this
STEP1. For example if w_i one
randomly initialized (as usual)
if $w_i = 0 \forall i \Rightarrow$ bias b need to
be $b \neq 0$ to avoid initial
deactivation of the NEURON

b : treated as a NOISE

= the bias handle generalization
of the model

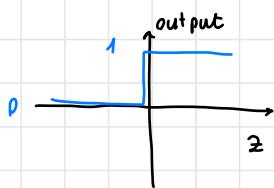
(In Linear Regression problems the
noise b is the intercept)



the importance of
ACTIVATION FUNCTION for

after performing STEP 1 and ← ex. in BINARY CLASSIFICATION
using data I want 0÷1 as output

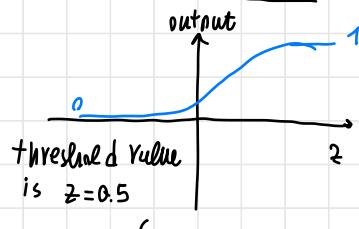
STEP Function



$$\begin{cases} 0 & z \leq 0 \\ 1 & z > 0 \end{cases}$$

threshold value
is $z=0$ in STEP

SIGMOID Function



threshold value
is $z=0.5$

$$\begin{cases} 0 & z \leq 0.5 \\ 1 & z > 0.5 \end{cases}$$

So the ACTIVATION

FUNCTION transforms output

of STEP 1 in a 0÷1 value

← solve BINARY CLASSIFICATION

require the mapping of

$$z = \sum_i x_i w_i + b \text{ in } [0, 1]$$

Ex) In the first data $x_1 = 95 \quad x_2 = 3 \quad (y = 0)$

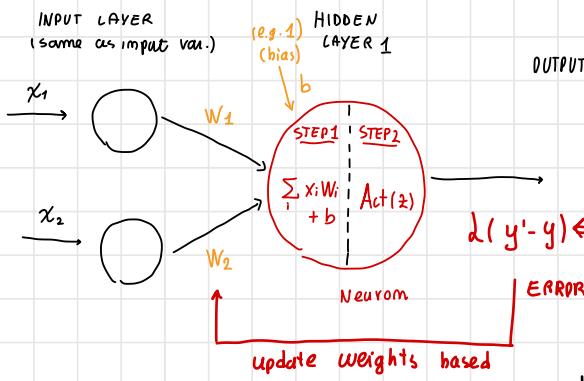
$$\text{output} = \text{Act}(x_1 w_1 + x_2 w_2 + b) = ? \rightarrow 0$$

If I get 0 as output of the single layer Perceptron
it is CORRECT, then

$$\text{ERROR} = 0$$

↓ good weights initialization!

IF ERROR ↑↑ I need to update the weights



the weights are updated by going FOREWARD / BACKWARD
(forward/backward propagation)

$\delta(y' - y) \leftarrow$ the ERROR used in the architecture
is the LOSS FUNCTION
(or COST FUNCTION)

update weights based on e... → weights update until the prediction is
satisfactory for all the records.(data)
good accuracy with respect to our data

Step-by-Step:

$$\text{STEP 1} \quad z = \sum_{i=1}^m w_i x_i + b = w_1 x_1 + w_2 x_2 + \dots + w_m x_m + b$$

↓ LINEAR PROBLEM STATEMENT

equation similar to $y = mx + c$ (LINEAR REGRESSION)

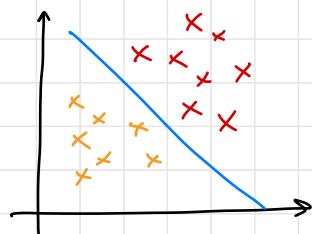
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m$$

+ STEP 2 manually selection

PERCEPTRON is
a linear
Classifier in
Step 1+step 2

overall with step 1 + step 2,
in a classification problem

• Perceptron create a
line that separate the
two classes!



At the end Single PERCEPTRON is a LINEAR CLASSIFIER
Combining STEP1 + STEP2

~ each ACTIVATION FUNCTION has its own functionality

↓
What are PROS and CONS of PERCEPTRON?

- + how to move to Multi layer Perception
- + how does a NN actually work

→ Single Neural Network
are NOT feasible for
Complex problems

PROS AND CONS OF PERCEPTRON

the one
discussed
before...

SINGLE LAYER PERCEPTRON

PERCEPTRON MODELS

⇒ What are its
PROS and CONS?

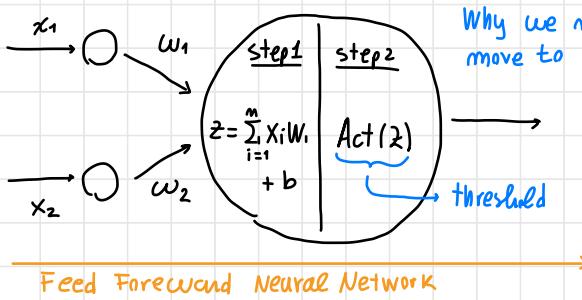
Why we need to
move to multi layer perception

MULTI LAYER

PERCEPTRON

= ANN

(multi layered
neural networks)



Feed Forward Neural Network

Performing all operations =
moving data from left to right

from a given input → we assign weights → we perform 2 steps... → after Activation the boolean classification guarantee 1 OR 0 output

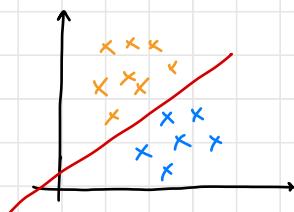
IF given x_1, x_2 I know $y=1$ in dataset

but I get $\hat{y}=0$ in the current model
← there is an ERROR!

the weights are updated
and we perform again
feed forward... until $\hat{y} \approx y$ (error negligible, good performance)

NOT efficient, but as PRO: Linearly separable Classification problems can be solved with Perceptron model.

e.g. linearly separable problem



using perceptron, the classification problem can be solved finding a best fit line.
BINARY CLASSIFIER

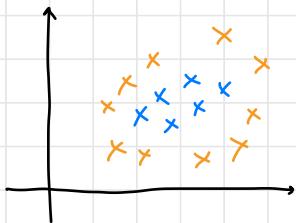
but changing w : randomly when $\hat{y} - y \approx 0$ is
NOT efficient...

In MLP models

we rely on **LOSS FUNCTIONS, BACK PROPAGATION** ← seamless weight update!

single layer NN perception are good just
for simple binary classifier **LINEARLY SEPARABLE**

IF instead



NON LINEAR ... a simple line doesn't exist
SEPARABLE

↓
MLP NN is required!

single layer is NOT enough

In MLP new techniques allow to achieve good performances

↓

- 1) Forward propagation
- 2) Backward propagation
- 3) Loss function
- 4) Various Activation functions → which one to use? When?
- 5) Optimizers!

After feed forward ⇒ a mechanism to update (weights)
of the architecture model is required

With MLP we can create DEEP NN models to solve complex problems

/

ARTIFICIAL NEURAL NETWORK (ANN)

MULTI LAYERED PERCEPTRON MODEL (ANN)

- 1) Forward Propagation
- 2) Backward Propagation → Geoffrey Hinton
(thanks to his research on back propagation techniques we can create DeepNN)
- 3) Loss functions
- 4) Optimizers
- 5) Activation function

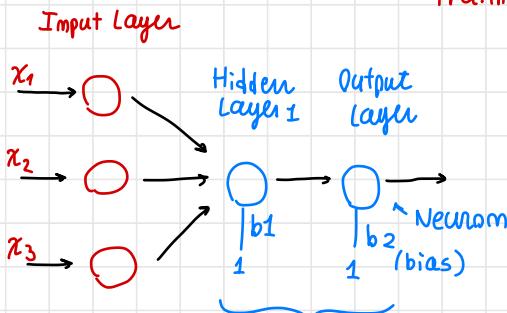
high level understanding of MLP concepts

MLP → Multi Layered Neural Network

{ While in a single layer Perceptron model we had feed forward NN + thresholding (activation function)
↳ that allow for binary classification }

In MLP we focus on a two step process Forward and Backward propagation (FP and BP)

in this 2 steps we can obtain efficient weight updates of the ANN, training the NN better



each layer can have multiple neurons
↳ (more than one hidden layer)

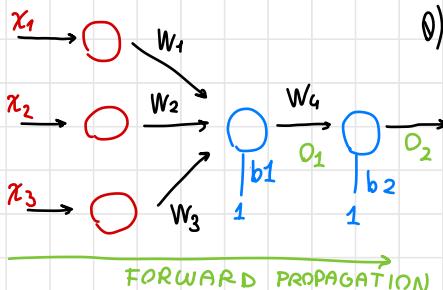
Dataset			output
IQ	Study Hours	Play Hours	Pass/Fail
95	4	4	1
100	5	2	1
95	2	7	0

$x_1 \quad x_2 \quad x_3$
3 INPUT features that pass to my input layer

↑ + bias $b_i \times (+1/-1)$ to avoid zero initialization
(random initialization)

2 Layered Neural Network (any number of hidden layers and neurons are feasible)

FORWARD PROPAGATION: sequence of steps through the NN

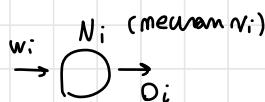


- ① • Feed Forward NN steps INITIALIZATION
 - 0.a) assign random bias $b_i = \text{rand}$ (some constant)
 - 0.b) assign random weights w_i (there are different initialization techniques)
+ each connection in the network require a weight

1) FORWARD PROPAGATION (FP)

In the neuron, two steps occurs

O_i is output of neuron N_i



$$\left\{ \begin{array}{l} \text{step1)} \sum_i x_i w_i + b_i = z_i \\ \text{step2)} O_i = \text{Act}(z_i) \end{array} \right.$$

$$\textcircled{1} \quad z_1 = x_1 w_1 + x_2 w_2 + x_3 w_3 + b_1 = \sum_{i=1}^3 w_i x_i + b \quad \textcircled{2} \quad O_1 = \text{Act}(z_1)$$

In FP, this operations is performed in ANY hidden neuron!

e.g. $w_1 = 0.01, w_2 = 0.02, w_3 = 0.03$ for $b_1 = 0.01$ HIDDEN LAYER 1

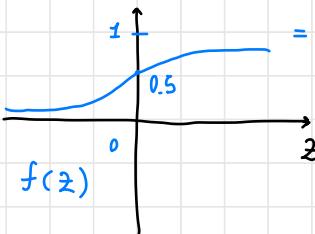
DATA x_1
 $x_{11} = 95 \quad x_{12} = 4 \quad x_{13} = 4$

$$\text{STEP1)} \quad z_1 = 95 \cdot 0.01 + 4 \cdot 0.02 + 4 \cdot 0.03 + 1 \cdot 0.01 = 1.151$$

$$\text{STEP2)} \quad O_1 = \text{Activation}(z_1) = \text{Activation}(1.151)$$

Consider for example

SIGMOID



many different activation functions exists (ReLU, ...)

↑ the activation $0 \div 1$ represent an

$$= \frac{1}{1 + e^{-z}} \in [0, 1] \forall z$$

important role in ACTIVATING and DEACTIVATING the Neuron ↑ when it has no more reason to respond (no role)

"ACTIVATION" state of the neuron based on the INPUT... (therefore response required next...)

This is possible only if correct weights are used.

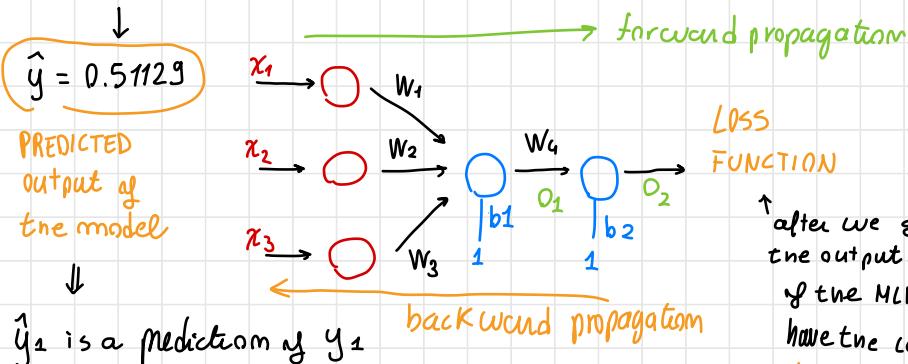
$$f(z) = \frac{1}{1+e^{-1.15}} = 0.759 \quad \text{hidden layer 1} \quad o_1 = 0.759$$

Hidden layer 2:

$$\text{initialized } w_4 = 0.02 \quad b_2 = 0.03$$

$$\text{STEP1}) z_2 = o_1 \cdot w_4 + b_2 = 0.759 \cdot 0.02 + 0.03 \cdot 1 = 0.04518$$

$$\text{STEP2}) o_2 = \text{Activation}(z_2) = \frac{1}{1+e^{-0.04518}} = 0.51129$$



$$\text{LOSS FUNCTION} = (y_1 - \hat{y}_1) = (1 - 0.51129) \approx 0.48$$

This is a big error!

$\hat{y} = 1$ would be $y - \hat{y} = 1$ OK!

ERROR WE GOT IN PREDICTION

my main objective is to reduce this error! \Rightarrow MINIMIZE THE ERROR OF PREDICTION

to reduce the error (and minimize it)

we update all the weights w_i : backword propagation process
from $w_4 \rightarrow w_3, w_2, w_1 \dots$

once backword propagation \leftarrow weights...
finish, w_i are updated for $x_1 = (x_{11}, x_{12}, x_{13})$ we update $w_i \forall i$

the next INPUT will be sent... $x_2 = (x_{21}, x_{22}, x_{23})$ and again forward/backward propagation occurs!

$\text{loss} = y_2 - \hat{y}_2$ and w_i updates... \uparrow the updates of w_i occurs if the loss needs to be reduced!

FP/BP process continue to reduce λ to satisfactory value.

LOSS FUNCTION

computing

$(y_i - \hat{y}_i)$ for every point
↑
calculate in each step updating the weights.

vs

COST FUNCTION

computing the error overall!

$$\sum_{i=1}^m (y_i - \hat{y}_i)^2$$

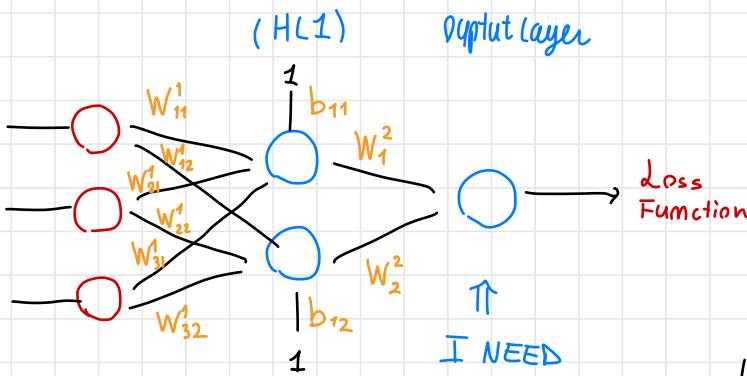
for all the dataset

calculate the error once + update weights just once.

In order to update w_i after the loss function is computed, during backward propagation **OPTIMIZERS** are used.

↑ helps in BP to reduce loss and updates weights

BACK PROPAGATION AND WEIGHT UPDATE FUNCTION



the number of neurons in the output layer depends on the problem solved

I NEED
to have
one neuron
in the output
layer IF I'm

solving a binary classification
problem $y = \{0, 1\}$

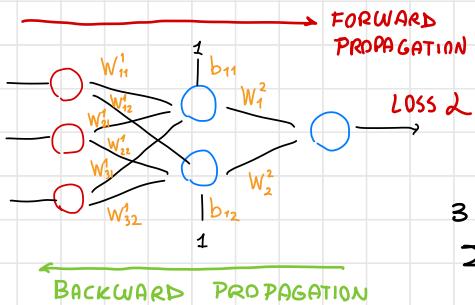
IF multi-class classification problem I can use more neurons

DATASET			Pass/Fail
IQ	Study Hours	Play Hours	
95	4	4	1
100	5	2	1
95	2	7	0
x_1	x_2	x_3	

INPUTS
↓
in dependent features

If this change the output change (dependent)

OUTPUT
↓
dependent feature



Having more neurons and input, now the operation during forward propagation changes.

$3 \text{ input layers} \times 2 \text{ hidden layers}$

$$\sum_i x_i W_{ij}^l = x_1 W_{11}^l + x_2 W_{21}^l + \dots$$

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \end{bmatrix} \quad W^l = \begin{bmatrix} W_{11}^l & W_{12}^l \\ W_{21}^l & W_{22}^l \\ W_{31}^l & W_{32}^l \end{bmatrix}$$

$\leftarrow x_i^T W^l : 1 \times 2 = 2, \dots$

than we forward
untie \hat{y} :

and being BINARY CLASSIFICATION PROBLEM the loss function depend

REGRESSION

- Mean Squared Error (MSE)
- Mean Absolute Error (MAE)
- Huber Loss

L

CLASSIFICATION

- Binary Cross Entropy
- Categorical Cross Entropy

basic loss: $(y - \hat{y})^2 \Rightarrow$ main aim is to reduce this loss

if this is BIG \rightarrow we update w_i during

BACKWARD PROPAGATION

Weight update formula

after loss computation, I start weight updates from the last layer W_1^2, W_2^2

$$W_1^2 \text{ NEW} = W_1^2 \text{ OLD} - \eta \frac{\partial L}{\partial W_1^2 \text{ OLD}}$$

positive or negative loss

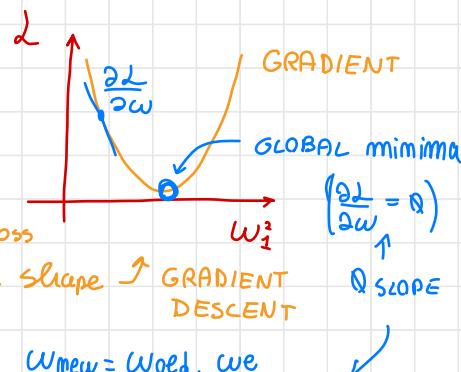
+ we update weights as this

SCOPE of the GRAPH $\Delta (W_1^2 \text{ OLD})$

usually the loss

has a parabolic shape \rightarrow GRADIENT DESCENT

where Δ slope $W_1^2 \text{ NEW} = W_1^2 \text{ OLD}$, we don't need to update weights



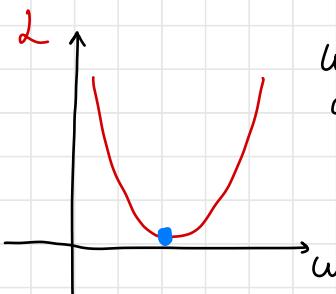
$$\text{in the same way } w_{2\text{new}}^2 = w_{2\text{old}}^2 - \eta \frac{\partial L}{\partial w_{2\text{old}}^2}$$

↓
this formula
work in updating the
Wi efficiently to reduce L

↑ LEARNING RATE (hyperparameter)

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}$$

Weight
update
formula



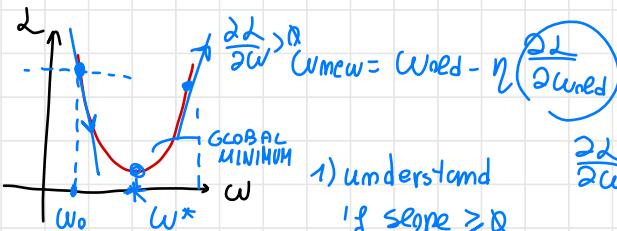
With the L used as MSE, MAE, ...
we get "gradient descent" shape L(w)

Main AIM is to update all Weights w;
With BACK PROPAGATION

In order to minimize the loss function $L = (y - \hat{y})^2$
we use OPTIMIZER S

GRADIENT DESCENT optimizer (for example ... other optimizers exists)

The ROLE of optimizers is to reduce the loss value L by
by updating the weights during BACKWARD PROPAGATION



1) understand

if slope ≥ 0

$\frac{\partial L}{\partial w_0} := \text{slope im } w_0 \text{ of } L$

right side of tangent point downward \Rightarrow slope < 0

$$\frac{\partial L}{\partial w_0} < 0 \text{ than } \Rightarrow w_{\text{new}} = w_{\text{old}} - \eta \left(-\frac{\partial L}{\partial w_{\text{old}}} \right)$$

$$= w_{\text{old}} + \eta \frac{\partial L}{\partial w_{\text{old}}}$$

We move
towards the
GLOBAL MINIMUM w^*

$w_{\text{new}} > w_{\text{old}}$
INCREASING!

In the same way if $w_{old} > w^*$ $\frac{\partial L}{\partial w_{old}} > 0$:

move on the left towards w^* $\Leftarrow w_{new} = w_{old} + \left(-\eta \left| \frac{\partial L}{\partial w_{old}} \right| \right) < w_{old}$



the main aim is to move toward w^* with update

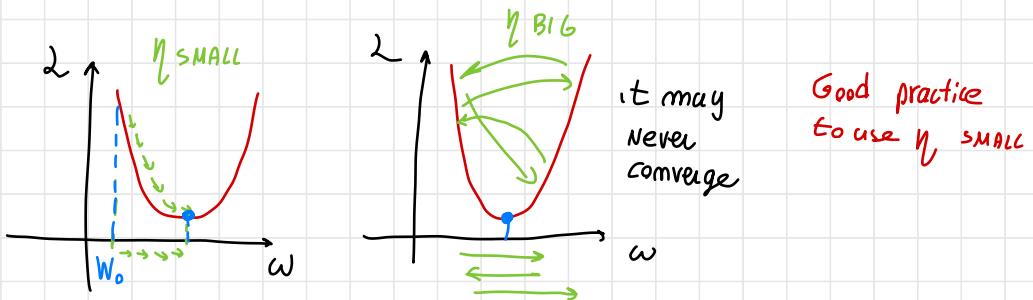
The **LEARNING RATE** η is an important HYPERPARAMETER

↓
Usually small initialization e.g. $\eta \approx 0.001$ SMALL STEP

define the "STEP SIZE" towards global minimum → HOW FAST CONVERGING

during ITERATIVE weight update

If η is BIG \Rightarrow We move with big steps! It can lose convergence



The optimizer knows that it reached the optimum when

$$\frac{\partial L}{\partial w} = 0 \Rightarrow \underline{w_{new} = w_{old}} \rightarrow w \text{ REACHES GLOBAL MINIM}$$

$L \approx 0$ no need to further change w !!

↑ IDEA of GRADIENT DESCENT OPTIMIZER

(other optimizers for weights update exists...)

BUT notice that w_i depends on w_j , $\forall j \neq i$, all weights are somehow connected since

$$d(w) = l(w_1, w_2, \dots, w_i, \dots, w_N)$$

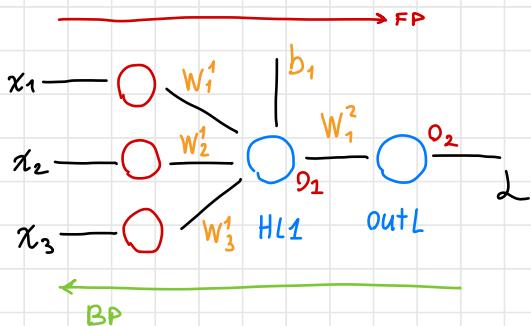
and each $\frac{\partial l}{\partial w_j}(w_i, w_{k\dots})$

→ a specific technique is considered

CHAIN RULE OF DERIVATIVE

✓ (used to compute those weights)

CHAIN RULE OF DERIVATIVES



related to w update formula
and how it can be computed

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial l}{\partial w_{\text{old}}}$$

the weights update start
from last ...

$$w_{1\text{NEW}} = w_{1\text{OLD}} - \eta \boxed{\frac{\partial l}{\partial w_{1\text{OLD}}}}$$

by expanding it l depends
on all w_i

$$\frac{\partial l}{\partial w_{1\text{OLD}}} = \frac{\partial l}{\partial O_2} \cdot \frac{\partial O_2}{\partial w_{1\text{OLD}}}$$

I can apply chain rule
of derivation and introduce O_2
CHAIN RULE of DERIVATIVES

this allow to compute each
component of the weights update

In the other terms

$$w_{1\text{NEW}} = w_{1\text{OLD}} - \eta \frac{\partial l}{\partial w_{1\text{OLD}}}$$

$O_1(w_1)$ as output of first hidden layer
+ $O_2(w_1^2)$ depends on O_1 also...

↙
chain rule apply to expand.

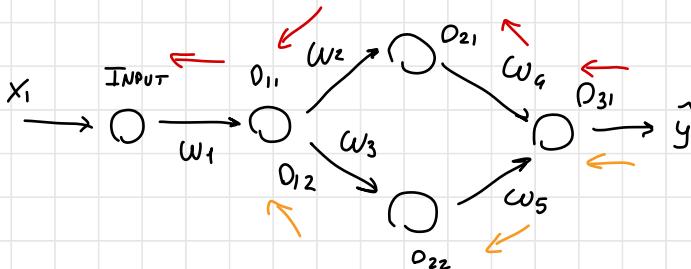
$$\frac{\partial L}{\partial W_{i,old}^j} = \frac{\partial L}{\partial O_2} \cdot \frac{\partial O_2}{\partial O_1} \cdot \frac{\partial O_1}{\partial W_{i,old}^j}$$

CHAIN RULE

$$L(O_2) = L(O_2(W_i)) = L(O_2(O_1(W_i))) = L(O_2(W_i^2(O_1))) \dots$$

↑ dependency chain!

Example



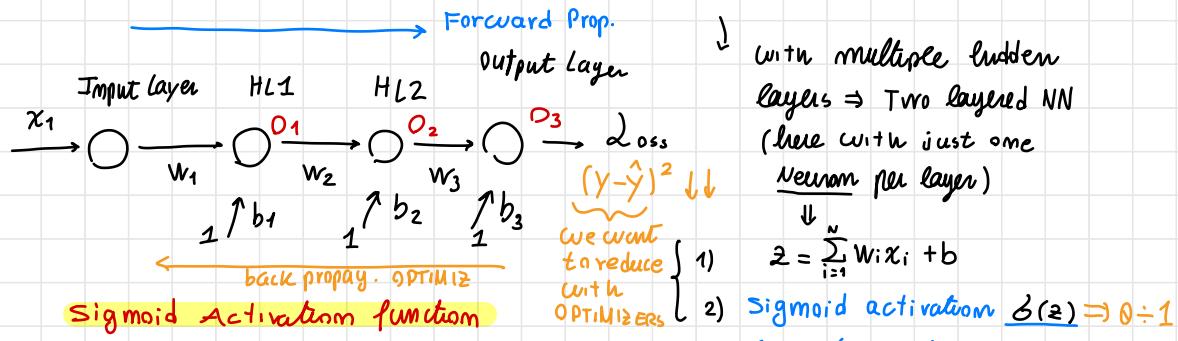
$$W_{i,new} = W_{i,old} - \gamma \left[\frac{\partial L}{\partial W_{i,old}} \right] \rightarrow \frac{\partial L}{\partial W_{i,old}} = \left[\frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial W_{i,old}} + \left(\frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial W_{i,old}} \right) \right]$$

to update
this chain rule
equation we will be used

This always happens with help of optimizer

$\frac{\partial L}{\partial O_i}$ can be computed efficiently as
loss wrt outputs of neurons

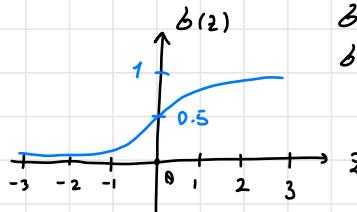
VANISHING GRADIENT PROBLEM AND ACTIVATION FUNCTIONS



Main aim is to transform its value as $\delta(z) \in [0, 1] \leftarrow$ useful for example in binary classification where the output should be between 0 and 1

↓
mathematically

$$\delta(z) = \frac{1}{1 + e^{-z}}$$



$$\begin{aligned} \delta(0) &= 0.5 \\ \delta(\infty) &\rightarrow 1 \quad \delta(-\infty) \rightarrow 0 \\ \text{always map } z &\text{ in } (0,1) \end{aligned}$$

here to update the weights

back propagation we start from behind to update pattern.

$$W_{\text{new}} = W_{\text{old}} - \eta \left| \frac{\partial L}{\partial W_{\text{old}}} \right|$$

\uparrow
 $\eta = 0.01$ random small initialization of the LEARNING RATE

$$\frac{\partial L}{\partial W_{\text{old}}} = \frac{\partial L}{\partial O_3} \cdot \left| \frac{\partial O_3}{\partial O_2} \right| \cdot \left| \frac{\partial O_2}{\partial O_1} \right| \cdot \left| \frac{\partial O_1}{\partial W_{\text{old}}} \right|$$

how can I solve this?

$$O_3 = \delta(W_3 \cdot O_2 + b_3) = \delta(z)$$

\downarrow
 $z = W_3 O_2 + b_3$

$$\leftarrow \frac{\partial O_3}{\partial O_2} = \frac{\partial \delta(z)}{\partial z} \cdot \frac{\partial z}{\partial O_2}$$

\uparrow
Sigmoid derivative

chain rule again.
to split $\partial O_3 / \partial O_2$

by applying this claim

since

$$\delta(z) \in (0, 1) \rightarrow 0 \leq \delta(z) \leq 1$$

always when apply activation

$$\frac{\partial \delta(z)}{\partial z} \in (0, 0.25) ! \text{ IMPORTANT} \Rightarrow 0 \leq \delta'(z) \leq 0.25 \quad (\text{math proof exists})$$

While $\frac{\partial z}{\partial \theta_2} = W_{3,0,ed}$

\downarrow

$$\frac{\partial \delta(z)}{\partial z} \cdot \frac{\partial z}{\partial \theta_2} = (0 \leq \delta'(z) \leq 0.25) \cdot W_{3,0,ed} = \underbrace{\frac{\partial \theta_3}{\partial \theta_2}}_{\in (0, 0.25)}$$

\leftarrow

$\in (0, 1)$ initialization in small value distribution

What happens is

that every derivative $0 \div 0.25$

$$\frac{\partial L}{\partial W_{1,0,ed}} = \underbrace{\frac{\partial L}{\partial \theta_3}}_{0 \div 0.25} \cdot \underbrace{\frac{\partial \theta_3}{\partial \theta_2}}_{0 \div 0.25} \cdot \underbrace{\frac{\partial \theta_2}{\partial \theta_1}}_{0 \div 0.25} \cdot \underbrace{\frac{\partial \theta_1}{\partial W_{1,0,ed}}}_{0 \div 0.25}$$

\rightarrow small value

small rate

\downarrow

$$W_{1,0,new} = W_{1,0,ed} - \frac{\eta}{2} \frac{\partial L}{\partial W_{1,0,ed}}$$

Negligible update of your weights!

\downarrow

almost null

you are not moving towards convergence $W_{1,new} \approx W_{1,0,ed} \dots$

$$\leftarrow W_{1,new} \approx W_{1,0,ed}$$

This "stack" occurs because of the sigmoid activation function

\downarrow PROBLEM of VANISHING GRADIENT



Very important problem is caused by sigmoid because $\delta(z)$ is s.t. $\delta'(z) \in (0, 0.25)$

In a Deep NN, with more layers ... the chains in the chain rule is longer \rightarrow even more vanishing

$$\frac{\partial L}{\partial W_{1,0,ed}} \approx 0 \quad W_{1,new} \approx W_{1,0,ed} \dots$$

In the same way to compute $w_{2\text{new}} = w_{2\text{old}} - \eta \frac{\partial L}{\partial w_{2\text{old}}}$

$$\frac{\partial L}{\partial w_{2\text{old}}} = \frac{\partial L}{\partial o_3} \left[\frac{\partial o_3}{\partial o_2} \right] \frac{\partial o_2}{\partial w_{2\text{old}}}$$

\rightarrow SMALL VALUE Keep appending

$$\frac{\partial o_3}{\partial o_2} = \frac{\partial(\delta(z))}{\partial z} \cdot \frac{\partial z}{\partial o_2}$$

$$\simeq [0 \div 0.25] \cdot w_3 \text{ SMALL VALUE (since } w_3 \simeq 0 \div 1)$$

as initialization techniques

$w_{2\text{new}} \simeq w_{2\text{old}}$ \rightarrow NO weight update!

We are back propagating but NOT converge!
it get stucked..

SIGMOID

PRO

- Suitable for Binary Classification: (0,1) output
- Clear prediction since it has output close to 0 or 1

CONS

- prone to VANISHING GRADIENT problem
- Function output is NOT zero-centered*
- The Exponential e^{-z} in $\delta(z)$ is demanding for computers (relatively time consuming)

this allow efficient
Weight update!!
↑
the zero-centering is done when
training models with
gradient descent

this helps in
weight updates

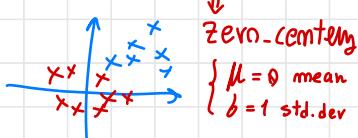
* ZERO-CENTERED

$\theta_{\text{mean}} \equiv \text{zero-centered}$

\Rightarrow efficient
Weights-update

In gradient
descent optimizers

This zero centering
in ML linear Regression
and classification we
perform STANDARDIZATION



To Fix this problem \Rightarrow Researchers Worked im Exploring other Activation functions

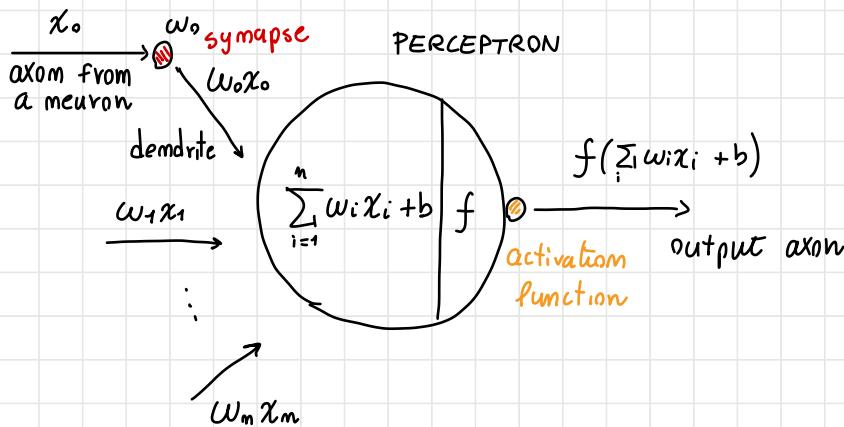
- ReLU
 - Tanh
 - Pre ReLU
 - Swiss
- } Which Activation should be used im HL? and which one im output layer.

Sigmoid can be good im small NN, BUT im a Deep NN you get a small $\Delta z/w$ that keep $w_{new} \approx w_{old}$ due to gradient vanishing

ACTIVATION FUNCTION

Attached to each neuron in the network, determines whether it should activate or not, based on whether each neuron's input is relevant for the model's prediction

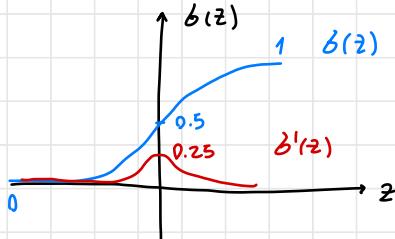
+ helps to normalize output im $0 \div 1$ or $-1 \div 1$ for each neuron



Sigmoid function

(RECAP)

transform $z \rightarrow \delta(z) \in (0,1)$ as $\delta(z) = \frac{1}{1+e^{-z}}$



because of this reason
we rely on other activations...

$\delta'(z) \in (0,0.25)$

cause of Vanishing Gradient problem

↓
during back propagation
the value of $\partial L / \partial w$
is small and the
update is stopped

the claim
rule of thumb
small number
cause this issue

- It was the most used originally, for SMALL NN

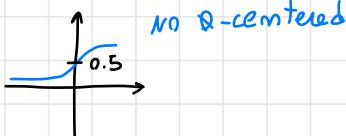
↓

but with DNN is a PROBLEM! but in principle is smooth and
easy to derive...

↖

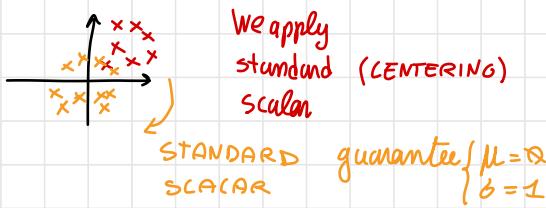
- Sigmoid ~ firing rate of a neuron (how much active/inactive)

- function output is NOT centered in \mathbb{R} ! Reducing efficiency of weight update



In ML algorithms

from a dataset



→ having a \mathbb{R} -centered function is good since, as proved, if points are zero centered (function zero centered) there will be an efficient update in weights

↓
helping us to compute
the derivative quickly

- the sigmoid function performs EXPONENTIAL operations which is slow for computers.
↳ takes time!

but $\delta(z)$ has advantages



- it is smooth, prevent "jumpy" in output
- good output normalization in $(0, 1)$
- clear predictions (very close to 1 or 0)

The main disadvantages

- prone to gradient vanishing
- Not-zero centered output
- time consuming exponential operation



because of this other possible activation functions exists, again with their pros and cons

TANH ACTIVATION FUNCTION

While sigmoid $\delta(z) \in (0, 1)$, $\delta' \in (0, 0.25) \leftarrow$ backpropagation face the gradient vanishing + disadvantages discussed...



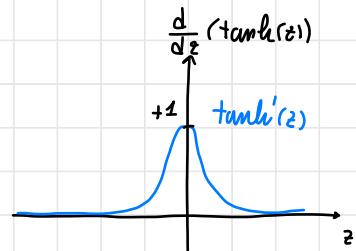
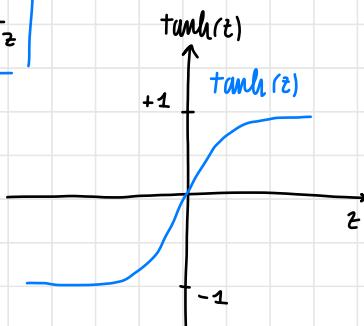
tanh solves some of the gradient descent problems

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$(z = \sum_{i=1}^m w_i x_i + b)$$

$$\tanh(z) \in (-1, 1)$$

∴ the activation by tanh map between -1 and 1



While $\frac{d}{dz} \tanh(z) \in (0, 1)$ always

$\tanh'(z) \in (0,1)$ \rightarrow advantage with respect to vanishing gradient problem

\tanh is an hyperbolic tangent function

and it has a curve similar to the sigmoid, but with $-1,1$ value

+ \tanh is θ -centric (better than sigmoid)

mean θ

being $\frac{d}{dz} \tanh \in (0,1)$ in a medium size DNN you don't face
gradient vanishing... **but in a very Deep NN you face some vanishing...**

↓

multiply a chain of small values...

For generic binary classifier \tanh is used for hidden layers, while
sigmoid for output layer

- + Notice as disadvantage that
- e^z is computationally expensive
 - vanishing problem still present in DNN
- advantage
- θ -centric = efficient weight update

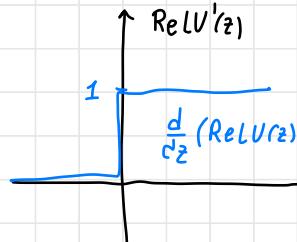
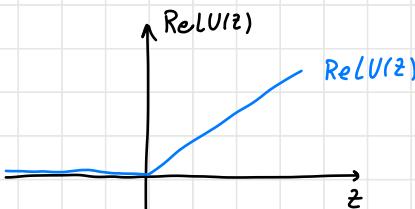


to fix some of this problems, other activation function exists... **to solve the MAJOR problem related to VANISHING GRADIENT PROBLEM** the ReLU activation is the solution...

ReLU ACTIVATION FUNCTION

~ this will solve the Vanishing gradient
BUT still one problem remain...

$$\text{ReLU}(z) = \max(0, z)$$



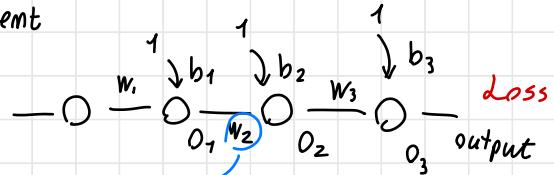
$$\begin{cases} \text{IF } z \leq 0 \Rightarrow \text{ReLU}(z) = 0 \\ \text{IF } z > 0 \Rightarrow \text{ReLU}(z) = z \end{cases}$$

during backward propagation

$$\frac{d}{dz} \text{ReLU}(z) = \begin{cases} 0 & \text{IF } z \leq 0 \\ 1 & \text{IF } z > 0 \end{cases}$$

this solves
vanishing gradient

in a deep NN



to update w_2 ... using previous equations

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial o_3} \underbrace{\frac{\partial o_3}{\partial o_2}}_{\{0,1\}} \frac{\partial o_2}{\partial w_2} \quad \{0,1\}$$

$$\hookrightarrow \frac{\partial o_3}{\partial o_2} = \frac{d(\text{ReLU}(z))}{dz} \frac{d z}{\partial o_2} = \{0,1\} \cdot w_3$$

IF $\frac{d}{dz} \text{ReLU}(z) = 1$: weight update will happen!
 $w_{\text{new}} \approx w_{\text{old}}$ NOT stacked!

IF $\text{ReLU}(z) = 0 \Rightarrow$ this creates a DEAD NEURON
if derivative of $\text{ReLU}(z)$ is 0: $w_{\text{new}} \approx w_{\text{old}}$

neuron without functionality ...

the entire neuron will pass the information,
without updating its value...

$$\frac{\partial L}{\partial w} = 0$$

$$\Downarrow$$

$$w_{\text{new}} = w_{\text{old}}$$

the $\frac{\partial \text{ReLU}(z)}{\partial z} = 0$ if $z \leq 0$

While if $z > 0$: $\text{ReLU}(z) = z \rightarrow \frac{\partial \text{ReLU}(z)}{\partial z} = 1$



$$\begin{cases} z > 0 : \text{OK } \checkmark \\ z \leq 0 : \text{dead Neuron} \end{cases}$$

even if this solve the vanishing gradient, when $z \leq 0$, it create a dead neuron..

ReLU: advantages

- It is solving the vanishing gradient problem
- Activation $\max(0, x)$ has a FAST calculation
- The ReLU function has a Linear Relationship
- It is much faster than sigmoid or tanh

disadvantages

- It leads to dead neuron problem
- ReLU function output is $(0, z)$



↑
It is NOT zero-centric

there are scenarios in which it is possible to handle that $z \leq 0$, but also other Activation functions can handle it

The next version of ReLU is Leaky ReLU...

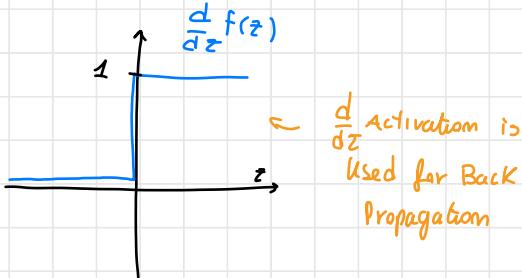
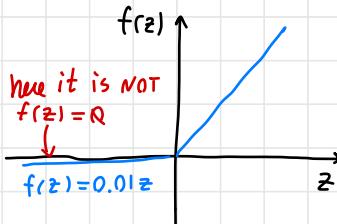
ReLU := "Rectified Linear Unit"

↑
Specific type of Linear Activation..

Another is ELU (Exponential Linear Unit)

Leaky ReLU and Parametric ReLU

$$f(z) = \max(0.01z, z)$$



- Since ReLU problem was related to the Dead Neuron Problem
 \downarrow
Dead ReLU Problem (entire Neuron deactivated)

\downarrow
 Solved by Leaky ReLU and parametric ReLU

by using a small constant instead of $\max(0, z)$ we use
 $\max(Kz, z)$

I'm **PARAMETRIC ReLU** this small $K \approx 0.01$

$$\max(\alpha z, z) = f(z)$$

\downarrow
 hyperparameter to tune \leftarrow by hyperparameter tuning

$\alpha = 0.01, 0.02, \dots$ depending on the specific problem

When $\alpha = 0.01$ is Leaky ReLU

Thanks to this parameter α (or 0.01 fixed) $\frac{d}{dz} f(z) \neq 0 \quad \forall z$

\downarrow
 always $\frac{d}{dz} f(z) > 0$ even for $z \leq 0 \dots$

\downarrow
 therefore in backpropagation we never get the dead Neuron!

\uparrow
 fixed this problem!

Advantages

- Leaky ReLU has all the advantages of ReLU
- There is NO dead ReLU problem
 - (it removes this problem thanks to small constant)
 - + it solves also vanishing gradient problem

↓

other Activations solves other problem...

Another Linear Unit is the Exponential Linear Unit activation...

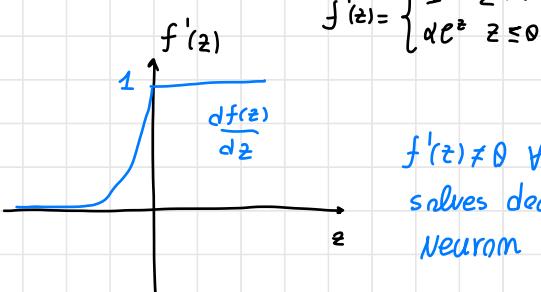
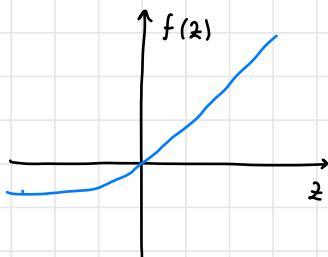
↑ Some Disadvantage of Leaky and Parametric ReLU

↓ still NOT zero-centric! => weight update will NOT be done efficiently

↓ with ELU we try to solve the problem of ReLU about dead Neuron + additional advantages...

ELU Activation Function (Exponential Linear Unit)

$$f(z) = \begin{cases} z & z > 0 \\ \alpha(e^z - 1) & z \leq 0 \end{cases}$$



$$f'(z) = \begin{cases} 1 & z > 0 \\ \alpha e^z & z \leq 0 \end{cases}$$

$f'(z) \neq 0 \forall z$ it
solves dead
Neuron problem!

Advantage

- No dead ReLU issues
- zero-centric Activation (the mean is 0)

Disadvantage

- Mathematically ... is slightly more intense computationally (Respect ReLU)

↑

Used to solve ReLU problems...

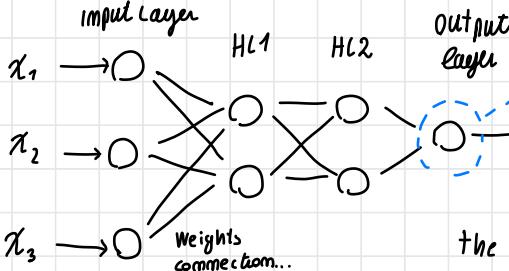
The usage of the Activation function is based on the problem and dataset

... the best combination of ReLU, Sigmoid, etc... is a "magic"

SOFTMAX ACTIVATION FUNCTION

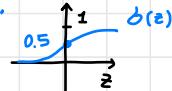
specifically used in the output layer for a multi-class classification problem

back propagation



with one output layer medium

BINARY CLASSIFICATION problem...

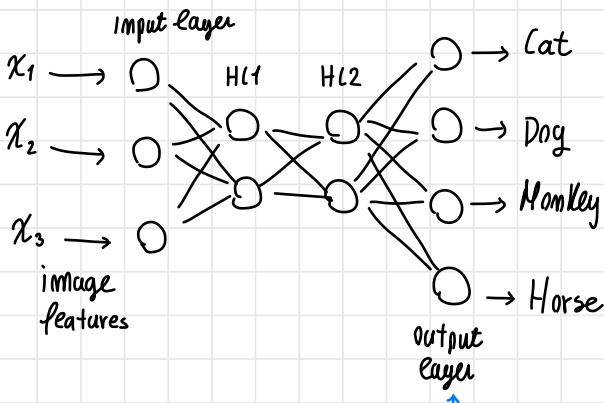


the output layer in binary classification has a SIGMOID ACTIVATION

it transform $\delta(z) \in (0,1)$ \Leftarrow that map the values in 0 to 1 map the output of last layer...

I can set a threshold. $\begin{cases} \delta(z) \geq 0.5 \Rightarrow 1 \\ \delta(z) < 0.5 \Rightarrow 0 \end{cases}$

BUT if instead I want to have multiple outputs:



I need to predict if the input features (x_1, x_2, x_3, x_4) are related to one of the output classes

In this scenario SIGMOID Activation is NOT the correct choice... (sigmoid give BINARY output)
Having 4 nodes in output layer each medium give Prob(Cat) ...

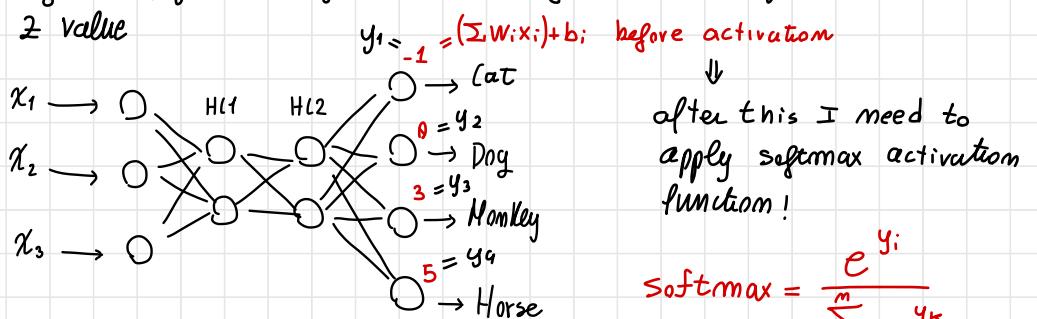
We want the PROBABILITY of each class...

↓ we need to use

Softmax activation function:

before apply this softmax, each of the last layer node has

z value



$$\text{softmax} = \frac{e^{y_i}}{\sum_{k=0}^m e^{y_k}}$$

In this case

$$i = \{1, 2, 3, 4\} \quad m = 4$$

↙ $i :=$ mode of the output layer
 $m :=$ # output mode

$$\text{softmax: } \text{Cat } (y_1) := \frac{e^{y_1}}{\sum_{k=0}^4 e^{y_k}} = \frac{e^{-1}}{e^{-1} + e^0 + e^3 + e^5} \approx 0.00033$$

($y_i = \text{output} * W + b$)

$$\text{Dog: } \text{Softmax}(y_2) = \frac{e^0}{e^{-1} + e^0 + e^3 + e^5} = 0.0024$$

$$\text{Monkey} = \frac{e^3}{e^{-1} + e^0 + e^3 + e^5} = 0.0183$$

$$\text{Horse} = \frac{e^5}{e^{-1} + e^0 + e^3 + e^5} = 0.1353$$

↖ this is the maximum...

↓
then the probability is found by normalizing

$$P(\text{Horse}) = 0.1353 / (0.00033 + 0.0024 + 0.0183 + 0.1353) \approx 86\%$$

$$P(\text{cat}) = 0.00033 / (...) \dots$$

↖ highest! → the MCP predict Horse

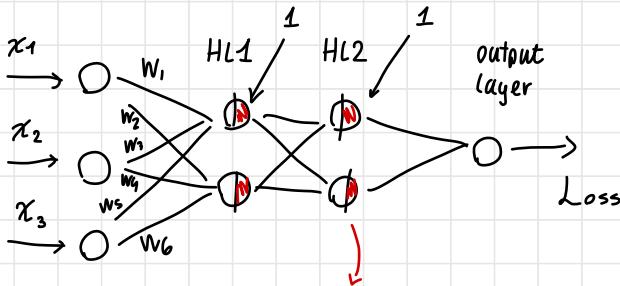
The Softmax is applied at output layer to calculate multi-class classification probability with respect to all categories present

Binary Classification \Rightarrow Sigmoid

Multi-Class Classification \Rightarrow Softmax

↓ But how to choose in the Hidden layers?
and in other problems...

WHICH ACTIVATION FUNCTION TO USE?



usually after weights initialization... and bias added, we apply activation function...

Input layer

When apply activation function...

When the NN is deep (many hidden layers...)

HL $i \Rightarrow$ Sigmoid \Rightarrow Vanishing Gradient Problem!

OR tanh
No

(In HL we avoid Sigmoid / Tanh)



ReLU or its variant

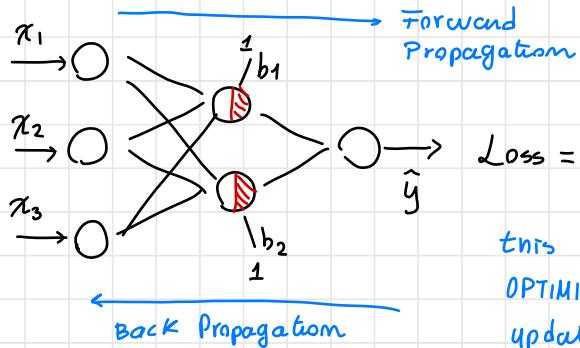
(Parametric ReLU, Leaky ReLU, ELU)

→ which solve Vanishing Gradient

While in OUTPUT Layer: $\begin{cases} \text{Binary classification: sigmoid} \\ \text{Multi class classification: softmax} \end{cases}$

(This holds in any ANN of any depth/length with any HL)

LOSS FUNCTION AND COST FUNCTION



$\text{Loss} = (y - \hat{y})^2$ ↓ goal to reduce the loss!
 ↓ (error on prediction)
 this is reduced with
 OPTIMIZERS... gradient descent for example
 update the weights during backpropagation

What are the different types ↴
 of LOSS FUNCTION? (and what about the COST FUNCTION)

Given data points:

(dataset)	x_1	x_2	x_3	Output
"	"	"	"	0
"	"	"	"	1
"	"	"	"	0
"	"	"	"	1

Loss function

we take ALL data points and do forward + backward propagation

e.g. MSE (mean square Error)

$$\text{MSE} = (y - \hat{y})^2 \quad \leftarrow \text{this happens } \forall \text{ datapoint } i$$

$$\text{MSE}_i = (y_i - \hat{y}_i)^2$$

Cost function

in this case we refer to a batch

$$\text{MSE} = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

instead of one point we pass all / batch of points of dataset...

We propagate through the net multiple datapoints × simultaneously
 ↓
 all \hat{y}_i are used to compute it...

so \forall output \hat{y} of datapoint the error is averaged... and again the objective is cost function reduction

So.. I'm loss function ↓
 for datapoint $x_i = (x_{i1}, x_{i2}, \dots, x_{if})$ in dataset
 we compute \hat{y}_i and calculate
 $L = (y_i - \hat{y}_i)^2$ and backpropagate + update weights...

While in the cost a batch $i = k_1 \dots k_b$ of points $x_{k_1}, x_{k_2}, \dots, x_{k_b}$
 and estimate $\hat{y}_{k_1}, \hat{y}_{k_2}, \dots, \hat{y}_{k_b}$ then the cost is
 update for set of points
 $\frac{1}{b} \sum_{i=k_1}^{k_b} (y_i - \hat{y}_i)^2$ and backpropagating we update
 the weight once

So, what are the main Loss and Cost functions?

(Focus on Regression / Classification)



With the help of ANN, it is possible to solve ↗ Classification

↘ Regression

in this two main problems we use different loss functions..

REGRESSION

- o Mean Square Error (MSE)
- o Mean Absolute Error (MAE)
- o Huber Loss
- o Root Mean Square Error (RMSE)

MSE

\hat{y} is the prediction output
 of the ANN after forward propagate

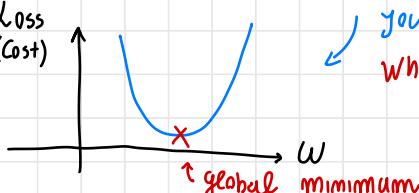
$$\text{Loss function} := \underbrace{(y - \hat{y})^2}_{\text{MSE}}$$

$$\text{Cost function} := \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

QUADRATIC EQUATION!

this is very important, meaning full choice

Loss
(Cost)



you have a paraboloid curve...
 which has a GLOBAL MINIMUM

And the objective is to update the weight s.t. the loss decrease towards global minimum

Parabola curve ~ Gradient Descent \leftarrow has ONE GLOBAL MINIMUM
+ knows how to update formulae

In COST FUNCTION also I get multi-dimensional paraboloid to update w_i ... the same in higher dimension

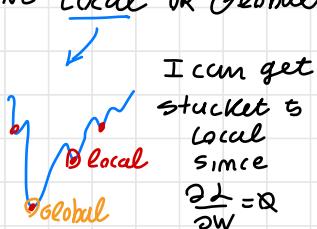
Advantages

- MSE is Differentiable
(@any time $d\text{loss}/dw$ can be computed)
as done during weight update

$$W_{\text{new}} = W_{\text{old}} - \eta \left(\frac{\partial L}{\partial W_{\text{old}}} \right) \xrightarrow{\text{slope of loss}}$$

If differentiable we can compute it

- MSE has ONE LOCAL OR Global Minimum



- MSE converges faster thanks to the loss shape

Disadvantages

- NOT Robust to outliers



ANN finds best fit line...

$(y_i - \hat{y}_i)^2$ error is squared.
because of outliers error is increasing by squaring out this $x_{\text{noise}} \dots$ PENALIZING the COST FUNCTION
the line move a lot!
(penalizing cost with errors)

MAE

$$\text{Loss function} = |y_i - \hat{y}_i|$$

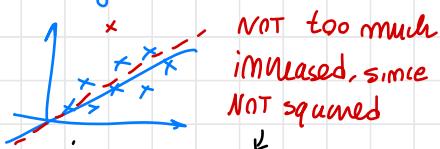
$$\text{Cost function} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

(Not squared, just absolute value..)

Advantages

- it is ROBUST to outliers

(NO more penalizing errors... even if I have outliers)



therefore MAE is better than MSE if outliers

(IF you don't have outliers MSE is better..) ←

for big outliers you can remove it? but lose data

HUBER LOSS (Combination of MSE + MAE)

$$\text{Cost function} := \begin{cases} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 & \text{if } |y_i - \hat{y}_i| \leq \delta \\ \delta |y_i - \hat{y}_i| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases}$$

← threshold hyperparameter ← if NOT outlier (low error)
← if outlier you scale up

↑

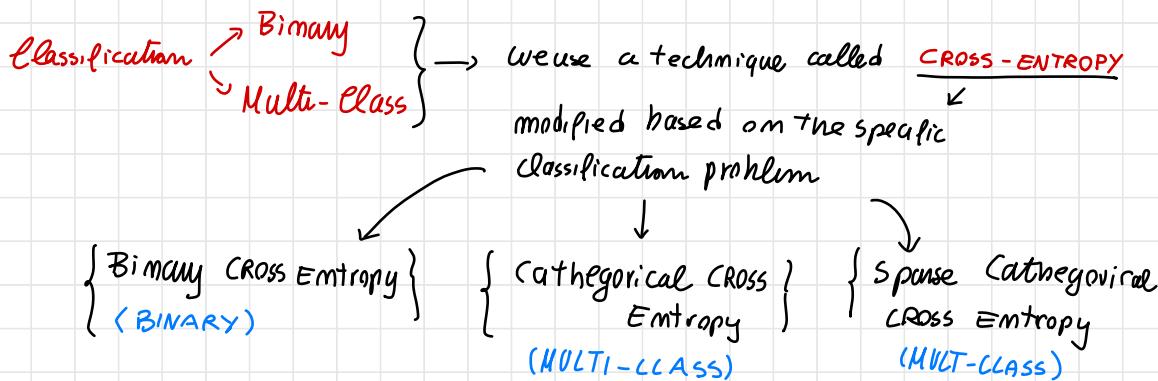
IF NO outlier you use MSE, otherwise you use a modification of MAE

RHSE is an update of MSE with square root

$$\text{Cost function} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Loss / Cost for CLASSIFICATION PROBLEMS

that we will minimize with optimizers



BINARY CROSS ENTROPY (used in Logistic Regression...)

$$\text{Loss} = -y \cdot \log(\hat{y}) - (1-y) \log(1-\hat{y}) \Leftarrow \text{log-loss formula}$$

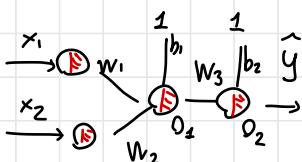
↓

$$\text{Loss} = \begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log(\hat{y}) & \text{if } y=1 \end{cases}$$

{ y := Actual value
 \hat{y} := Predicted value

two scenarios in BINARY CASE

↓ example with ANN



$$O_2 = \sigma(z_2)$$

SIGMOID since
BINARY CLASSIFICATION
(output can be 0/1)

$$O_2 = \sigma(\underbrace{O_1 * w_3 + b_2}_z) = \hat{y}$$

$$\hat{y} = \frac{1}{1+e^{-z}}$$

that loss is computed on this \hat{y} of simple sigmoid..

CATEGORICAL CROSS ENTROPY

(used in multi-class classification)

from a dataset with features f_1, f_2, f_3 :

f_1	f_2	f_3	Output	multi-class!
2	3	4	Good	(more than two)
5	6	7	Bad	
8	9	10	Neutral	

the categorical cross entropy

start from output processing \Rightarrow apply OHE (One Hot Encoding)

Output	$j=1$	$j=2$	$j=3$	on the output data
Good	1	0	0	
Bad	0	1	0	
Neutral	0	0	1	

then the loss function

$$\mathcal{L}(x_i, y_i) = - \sum_{j=1}^c y_{ij} \times \ell_m(\hat{y}_{ij}) \quad | \quad i=1, \dots, N \text{ (all data points)} \\ \text{---}$$

$c := \# \text{ of categories in the OHE}$

Where $i := \text{datapoint's index}$

$$y_{ij} = \begin{cases} 1 & \text{if the element is in the class } j \\ 0 & \text{otherwise} \end{cases} \quad | \quad \leftarrow$$

↑ actual value from the dataset ...

How I compute \hat{y}_{ij} ? (the PREDICTION) \rightarrow in a multi-class classification

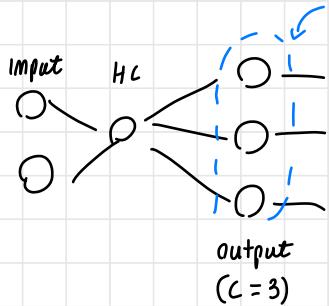
↓
the Activation Function applied is

therefore \swarrow

$$\text{Softmax activation} \simeq \text{Soft}(z) = \frac{e^{z_i}}{\sum_{j=1}^c e^{z_j}}$$

...

Softmax activation function in output
output:



$\hat{y}_{ij} \in [0,1] :=$ probability of
data i to be of class j
 $\Rightarrow [0.2, 0.3, 0.5]$

categorical
cross entropy output...

If $C=5$ I get $[0.1, 0.2, 0.3, 0.2, 0.2]$ for example... $(\sum_{j=1}^5 \hat{y}_{ij} = 1)$
↑ probability of belonging
to category j

↓

this also gives the probability of other categories
(this differentiate from sparse...)

SPARSE CATEGORICAL CROSS ENTROPY

Still probabilities output... $[0.2, 0.3, 0.5]$

Disadvantage:

losing info about
the probability of
other category

↓
highest index...
2nd Index output

We focus just on

+ the class with the highest probability (NOT in the overall)

from $C=5$ $[0.2, 0.3, 0.1, 0.2, 0.2]$

↑ final output tell me P_{\max} is given by 1st prob.

this differentiate from categorical standard...

To know the overall probability on all classes we
rely on categorical cross-entropy...

To get only the highest probability output you rely on sparse one

WHICH LOSS FUNCTION TO USE...

↓

What about Right Combination? of the problem statement
and ANN architecture...

Hidden Layer	activation	Output Layer	Problem statement	Loss function
ReLU or its variants	Sigmoid		BINARY CLASSIFICATION	Binary Cross Entropy
ReLU or its variants	Softmax		MULTI-CLASS CLASSIFIC.	Categorical/Sparse Cross Entropy (CE)
ReLU or its variants	Linear (specific activ.)		REGRESSION	MSE, MAE, Huber loss, RMSE

↑ activation function
[applied in output layer]

Used to solve Regression

... optimizers, weight initialization,... other components to describe in DL..

OPTIMIZERS

once defined and computed the **LOSS FUNCTION**
both in Regression/Classification problems

How the weights ↴

are updated to minimize the Loss? **optimizers are responsible of that**

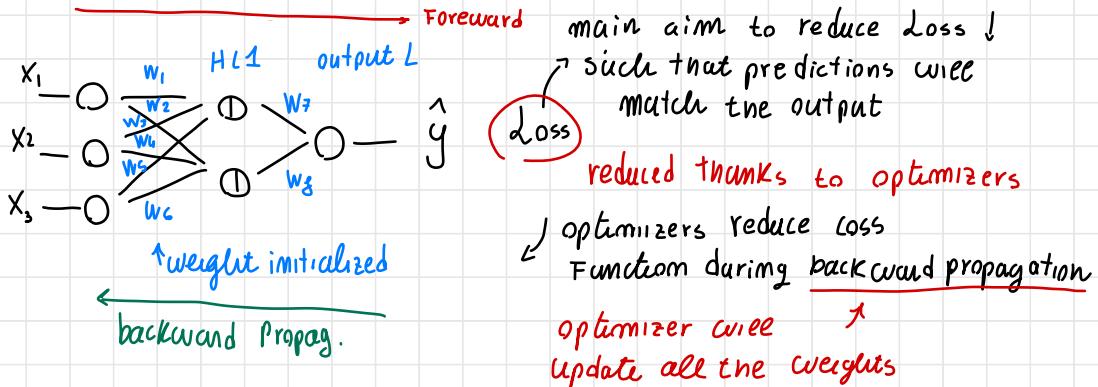
- Gradient Descent
- Stochastic Gradient Descent (SGD)
- Mini Batch SGD
- SGD with Momentum
- Adagrad and RMSProp
- Adam optimizer

↳ start optimizers description...

GRADIENT DESCENT

OPTIMIZER

Optimizers are responsible of reducing the loss function! fundamental part



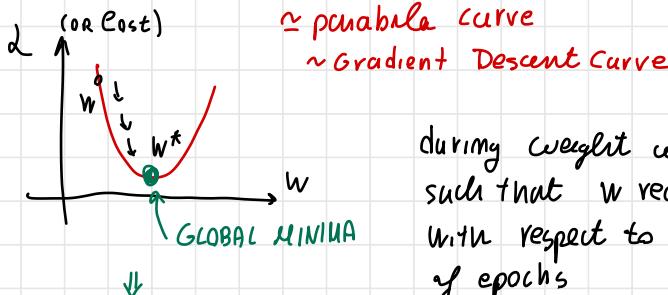
Weight Update

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial L}{\partial W_{\text{old}}}$$

each optimizer update the weights in a different way... different weight update

Learning rate \approx speed of convergence to GLOBAL MINIMUM

plotting $L(w)$ after Forward



during weight update w changes such that w reaches w^* minimum with respect to different number of epochs

When calculate

$$\frac{\partial L}{\partial w} @ \text{global minima} \quad \frac{\partial L}{\partial w_{\text{old}}} = 0 \rightarrow \quad W_{\text{new}} = W_{\text{old}} \quad \text{CONVERGENCE to optimal point}$$

• What does Gradient Descent do?

Considering my MSE loss := $(y - \hat{y})^2$ cost func = $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

In GRADIENT DESCENT:

two fundamental concepts are { Epochs
Iterations

forward $\rightarrow \hat{y}$ = cost function.. take more data points \vec{x}_i , ($\vec{x}_1, \dots, \vec{x}_n$)
compute \hat{y} ($\hat{y}_1, \dots, \hat{y}_n$)

Having 1000 datapoints in dataset

1 Epoch { 1000 Datapoints $\rightarrow \hat{y}_i$ ($i=1 \dots 1000$) \Rightarrow cost function ℓ
backward propagation
Weights update \otimes

than to Reduce this we
rely on GRADIENT DESCENT
optimizer

{ one forward propagation with
all datapoints + backward + weights update } = 1 Epochs

continue this process until $\ell \downarrow$ reduce enough

2nd Epoch { \rightarrow
 $\leftarrow \otimes$ } perform multiple
Epochs...
in each epoch
GRADIENT
DESCENT
change the
Weights with same weight update formula

In only Epoch $\leftarrow \otimes$

$\ell \downarrow$ decrease...

Gradient Descent optimizers
works in this way!

Iteration

✓ Epoch I perform back-propagation ... In gradient descent as defined above $1 \text{ epoch} = \underbrace{1 \text{ iteration}}_{\uparrow}$

If instead I divide datapoints in 10% batch

$$\text{Data}_1 = 100 \text{ records} = 1000 \cdot 0,1$$

$$\text{Data}_2 = 100$$

:

$$\text{Data}_{10} = 100$$

$$\left. \begin{array}{l} \\ \\ \end{array} \right\} 1 \text{ epoch} = 10 \text{ iterations}$$

NOT in Gradient Descent, where you take all dataset in each epoch

What are PROs and CONS?

Advantage

- 1) convergence will happen
(due to how defined W_{new})

Disadvantage

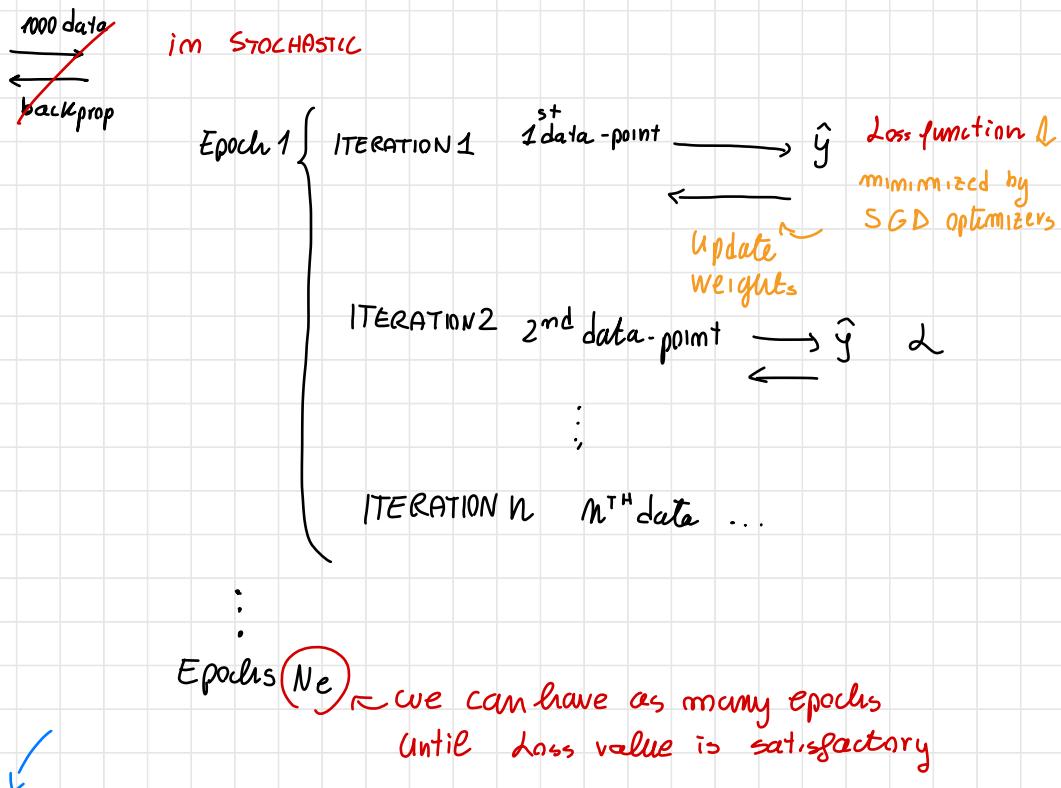
- 1) you require huge resource RAM, GPU since you need to compute forward for all datapoints!
↓
huge RAM to store and update all!....

NOT good in limited computational resources system...
Resource intensive!

STOCHASTIC GRADIENT DESCENT (SGD)

Optimizer

This helps solving the "Resource Intensive" problem by addressing EPOCHS / ITERATIONS



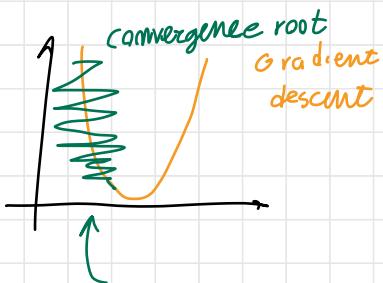
:) Advantage

- Now even if limited RAM, GPU,
you don't allocate all
simultaneously...
SOLVE Resource Issue

Disadvantages

- Time complexity is Huge!
(takes lots time to finish...)
 1000000 datapoints $\times 100$ epochs...
take lots of time

- convergence will take more time
- NOISE GETS INTRODUCED



as Disadvantage of SGD

While classical GD smoothly converge towards minimum, since you use entire RECORDS of dataset, guarantee Smoothing

While in SGD taking single data points at time, you update less smoothly

Time complexity increases + introduce noise

due to the
SINCE data usage

Model training and convergence takes time...

How to REDUCE this Noise? → concept of BATCH

Instead of taking single data points, you take batches of data points and perform forward/backward propagation.

The noise can be easily reduced through the concept of another optimizer

MINI BATCH SGD Optimizer

"Noisy" because noise in the

Weight convergence...

Path not smooth

(lot of ZIG ZAG towards minimum)

while SGD is very slow, with lots of

iterations and Noisy update since each iteration take one data

↳ In Mini Batch SGD with Epoch and Iterations we introduce BATCH and Batch-size

having # Data points = 100000 and batch-size = 1000

In Epoch 1: Iteration 1 $\xrightarrow{1000}$ Cost func $\sum_{i=1}^{1000} (y_i - \hat{y}_i)^2$ A epoch I take batch-size 1000 and use it in each iteration
optimizer i.e., Mini batch (MSE) \downarrow # iteration = $\frac{\text{Data points}}{\text{batch size}} = \frac{10^5}{10^3} = 100$
optimize weight \leftarrow SGD changing \hat{y}_i , t since with 1 data-iteration

Iteration 2 $\xrightarrow{1000}$ $C = \sum_{i=1}^{1000} (y_i - \hat{y}_i)^2$ I have lot of NOISE... this
Mini batch SGD concept of batch size help to reduce the noise

:
Iteration 100

\uparrow I use cost function wrt the batch size data

$$C = \sum_{i=1}^{\text{batch}} (y_i - \hat{y}_i)^2$$

Epoch 2 Iteration 1

:
Iteration 100

With this approach... Having $\mathcal{O}(GB)$ RAM,

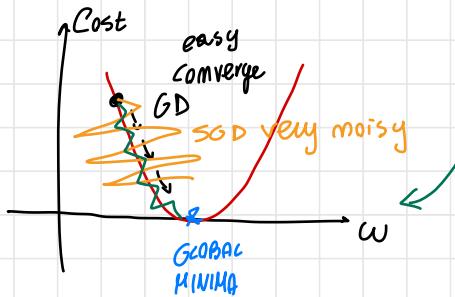
With 1000 data-points it will NOT be a problem
in resources

While 5000 batch-size RAM $\uparrow\uparrow$

With batch-size $\uparrow\uparrow$ the

iteration $\downarrow\downarrow$

and Reduce the NOISE!



thanks to batch of batch SGD
the noise get reduced

\downarrow
No more single record at
each iteration! \Rightarrow reduce noise
towards min.

Advantages

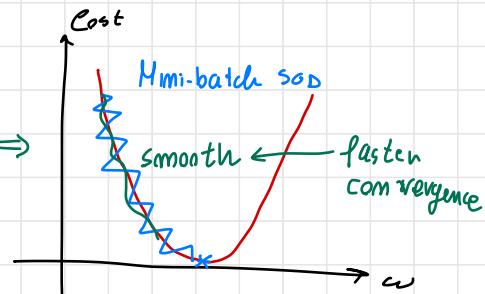
- Convergence speed increase w.r.t SGD (+thanks to batchsize)
- Noise will be less compared to SGD
- It solves efficiency resource problem efficient resource usage (RAM)

Disadvantage

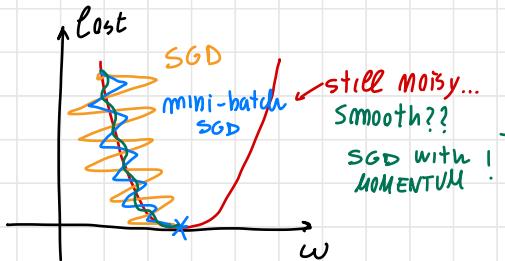
- The noise still exists!
still slow down a bit convergence

\downarrow
the next objective
is to smooth the noise
help to reach global
minima faster

this is achieved by
SGD with MOMENTUM
(smoothening) \Rightarrow



SGD with MOMENTUM Optimizer



How the smoothening can happen?

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial L}{\partial W_{\text{old}}}$$

Change this formula to incorporate momentum (smoothening)

$$b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial L}{\partial b_{\text{old}}}$$

Also, during back propagation the bias gets updated

Rewrite in time series

Weight update formula

$$W_t = W_{t-1} - \eta \frac{\partial L}{\partial W_{t-1}}$$

by using time dependency we can understand how to smooth

↓ to guarantee smoothening, we rely on...

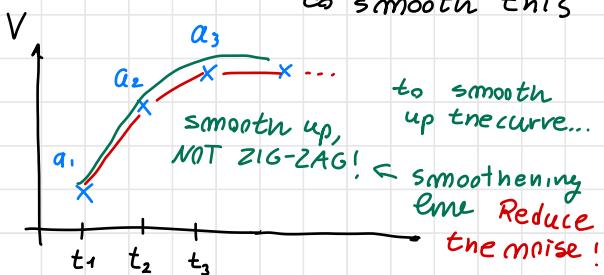
EXPONENTIAL WEIGHT AVERAGE (SMOOHENING)

Time:	t_1	t_2	t_3	t_4	...
Value:	a_1	a_2	a_3	a_4	...

t_n a_n ↗ this is a "time-series" problem... we want to smooth this

↓
Exp Weighted Average } $V_{t_1} = a_1$
Smoothening

$$V_{t_2} = \beta V_{t_1} + (1-\beta)a_2$$



β control the smoothening function ex. $\beta = 0.95$ $\beta \in [0, 1]$

$$V_{t_2} = 0.95 \times a_1 + 0.05 a_2$$

give control to the a_2 ← a_2 has less effect
Value

With this $\beta \approx 1$ the value a_2 control MORE, if $\beta \gg 0.5$ a_1 control the smoothening

Exp weighted avg allow smoothening \rightarrow REDUCE THE NOISE!

$$V_{t_1} = a_1$$

$V_{t_2} \dots \text{not } a_2 \leftarrow$ you smooth giving weight to a_2/a_2

:

the β value control how the smoothness
is controlled by current - previous point...

$$\begin{matrix} x \\ a_1 \\ a_2 \end{matrix}$$

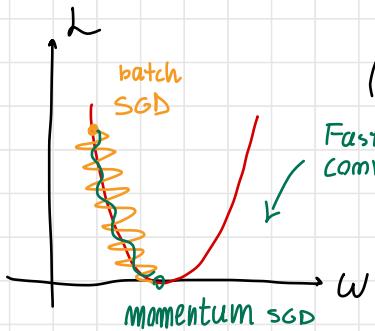
If β is big, a previous control more..

at same way $V_{t_3} = \beta V_{t_2} + (1-\beta) a_3 = 0.95 [0.95 a_1 + 0.05 a_2] + 0.05 a_3$

\Downarrow
 $\beta V_{t_1} + (1-\beta) a_2$

^t the second point
control smoothening
of third point

↑
This Exponential weighted average
can be used on $\partial L / \partial W_{t-1}$ \downarrow smoothening the entire curve



→ exp. weighted average helps in
finding out the weights...
it's introduced in
weight update

$$W_t = W_{t-1} - \eta \frac{\partial L}{\partial W_{t-1}}$$

↪ chain Rule we apply
exp. weighted avg
smoothening guarantee \Downarrow

Advantage

- Reduces noise
- Quick convergence

this exp weighted avg happens
With techniques as ARIMA, SARIMAX...
for time series

- importance of Learning rate η ? \Rightarrow we want a dynamic Learning rate...

ADAGRAD: ADaptive GRADient Descent Optimizer

until now

$$W_t = W_{t-1} - \left(\eta \frac{\partial L}{\partial W_{t-1}} \right) \text{ update weight}$$

Learning rate $\in [0,1]$ slow speed of convergence,
usually initialized small at beginning

$\eta = \text{FIXED}$ during convergence in GD, SGD...

DYNAMIC? What if η is dynamic?



As the convergence happen
the learning Rate should
change!

initially y big, fast, then we
reduce y when close to min

ADAGRAD make η_t
dynamic controlling convergence speed

$$W_t = W_{t-1} - \left(\eta' \frac{\partial L}{\partial W_{t-1}} \right)$$

η' : dynamic LEARNING RATE (Not Fixed!)

$$\eta' = \frac{\eta}{\sqrt{d_t + \epsilon}}$$

$\epsilon := \text{small value}$
(NEVER $\epsilon = 0$, even
when $d_t = 0$ is fine!)
 $d_t + \epsilon \neq 0$ always

$$d_t := \sum_{i=1}^t \left(\frac{\partial L}{\partial W_i} \right)^2$$

as back propagation wrt t , $d_t \uparrow$ therefore $\eta' \downarrow$

at beginning d_t is small and η' is big, encn opposite...

η' get decreasingly with a finer small step convergence

$$\begin{array}{lll} t=1 & t=2 & t=K \\ \eta = 0.01 & \eta = 0.005 & \dots \quad \eta = 0.0001 \\ \text{initialize} & \text{decrease} & \end{array}$$

↳ as Advantage you get dynamic η ...

but Disadvantage in a very deep NN, $d_t \uparrow$ and $\eta' \downarrow$ very small!

- η' can become a very small value approximately equal to 0
 $W_t \approx W_{t-1}$ in this bad case... Weight update doesn't occur
 ↗ this disadv. is addressed by a new optimizer...

ADADELTA and RMSPROP optimizers

While ADAGRAD uses a dynamic Learning rate $\eta' = \eta / \sqrt{\alpha_t + \epsilon}$

but in Deep Layered Neural Network $d_t \uparrow$ since chain rule is costly...
 $\frac{\partial L}{\partial w_t}$ over entire depth... $\Rightarrow \eta' \downarrow \rightarrow 0$ $\Rightarrow W_t \approx W_{t-1}$, No weight update!

To fix the problem on $\eta' \rightarrow 0$ converging to 0 and weight update stucked... New formula

$$\eta' = \frac{\eta}{\sqrt{Sd_w + \epsilon}}$$

exponential weighted average
SMOOTHENING the update of d_t

$$Sd_w = \beta \cdot Sd_{w-1} + (1-\beta) \left(\frac{\partial L}{\partial w_{t-1}} \right)^2$$

IF $\beta = 0.95$, even if $\frac{\partial L}{\partial w_t} \uparrow$
 restricting by $(1-\beta) = 0.05$ I
 restrict Sd_w

↳ Restricting $\frac{\partial L}{\partial w_{t-1}}$ with β ...

With $Sd_w = 0$ at beginning .. then update

AdaDelta and RMSprop has dynamic learning rate η

↓

Exponential Weighted Average
with EWA

Restricting α_t to be small $\alpha_t = Sdw_t + \text{smoothing}$

AdaDelta / RMSprop = Dynamic LR + Smoothing EWA

↑ ↑

two optimizers with slight differences but same concept
against same update

$$W_t = W_{t-1} - \eta' \frac{\partial L}{\partial W_{t-1}}$$

$$\eta' = \frac{\eta}{\sqrt{Sdw + \epsilon}}$$

→ see this technique
doesn't have
momentum in weights
etc.. combination

η more adv.

→ Adam Optimizer bring mainly more advantages
considered very good

↔

ADAM OPTIMIZER

It collect all the advantages of previous
discussed optimizers

Architecture..

SGD with MOMENTUM
{less noise}

+ RMSprop

{dynamic LR}
+ smoothing
(so α_t never
too big)

Adam optimizer
mix this two
features...

Weight update

weight update

$$\begin{cases} W_t = W_{t-1} - \eta' V_{dw} \\ b_t = b_{t-1} - \eta' V_{db} \end{cases}$$

b as update

As SGD with Momentum
you do EWA

(major updates) dynamic learning rate

$$\eta' = \frac{\eta}{\sqrt{Sdw_t + \epsilon}}$$

(Dynamic LR)

⊗

$Sdw_0 = 0$ initialize

$$Sdw_t = \beta \cdot Sdw_{t-1} + (1-\beta) \left[\frac{\partial L}{\partial W_t} \right]^2$$

EWA

SMOOTHING (by Sdw_t)

$$Vdw_t = \beta Vdw_{t-1} + (1-\beta) \frac{\partial L}{\partial W_{t-1}} \quad \text{smoothening of weights}$$

(EWA)

$$Vdb_t = \beta Vdb_{t-1} + (1-\beta) \frac{\partial L}{\partial b_{t-1}} \quad \text{bias smoothening}$$

~~SGD~~ this reflects SGD with MOMENTUM

Momentum

~ All combined in weight update

formula and this Vdx_t allow MOMENTUM \rightarrow smoothening weight update

also RMSPROP with smoothening is applied on Adam rely on EWA

EWA updates in weights/bias update

Adam optimizer helps in weight convergence quickly, since it mixes the best properties

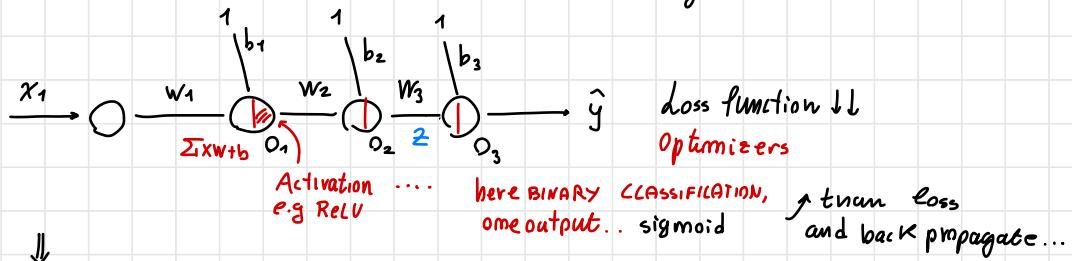


quick convergence + batch size + learning rate dyn + smoothening

best in MOST of use cases !!

EXPLODING GRADIENT PROBLEM

It is related to the weights and can be "bad"...



$$w_{1\text{ new}} = w_{1\text{ old}} - \eta \left| \frac{\partial L}{\partial w_{1\text{ old}}} \right|$$

when solving the $\frac{\partial L}{\partial w}$
using chain rule

$$\frac{\partial L}{\partial w_{1\text{ old}}} = \frac{\partial L}{\partial o_3} \cdot \boxed{\frac{\partial o_3}{\partial o_2}} \cdot \frac{\partial o_2}{\partial o_1} \cdot \frac{\partial o_1}{\partial w_{1\text{ old}}}$$

As next step...

$$z = o_2 w_3 + b_3$$

then Activation function applied on z b(.) sigmoid

$$\frac{\partial o_3}{\partial o_2} = \underbrace{\frac{\partial \delta(z)}{\partial z}}_{\in [0, 0.25]} \cdot \boxed{\frac{\partial z}{\partial o_2}} = [0 \div 0.25] \circled{w_3}$$

\leftarrow Vanishing previously was related to $\frac{\partial \delta}{\partial z} \dots$ now not issue
of VANISHING ✓ :

What if w_3 is BIG

$w_3 = 100 \div 1000 \dots$ in any part of chain Rule because weight initialization \Rightarrow high value

$$\frac{\partial L}{\partial w_{1\text{ old}}} = \frac{\partial L}{\partial o_3} \cdot \frac{\partial o_3}{\partial o_2} \cdot \frac{\partial o_2}{\partial o_1} \cdot \frac{\partial o_1}{\partial w_{1\text{ old}}}$$

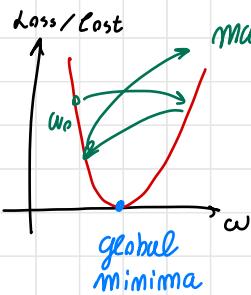
↑↑ with big w_i initialization...
high high high high \leftarrow big numbers..

big value! $w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}$

When this update by big step happen

$$\begin{cases} w_{\text{new}} \gg w_{\text{old}} \\ w_{\text{new}} \ll w_{\text{old}} \end{cases}$$

← big value



Maybe it diverges! and never reach minimum

EXPLODING GRADIENT DESCENT

due to weight initialization problem!

IF w_i are too big...



different w_i initialization techniques
allow to avoid the EXPLODING GRADIENT
PROBLEM!

weight initialization

techniques → allow to guarantee convergence of weights
(scientifically proven)

WEIGHT INITIALIZATION TECHNIQUES

- 1) Uniform Distribution
- 2) Xavier / Glorot Initialization
- 3) Kaiming He Initialization

Key Points

- Weights should be small ← if big you face gradient explosion

- Weights should Not be same ←

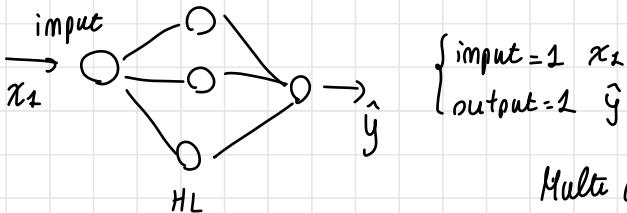
- Weights should have good variance ↗

$w_1, w_2 \dots$ variance changes
as different with big diff

$w_i \neq w_j$ and small!
otherwise they "think"
in the same way and does
the same processing

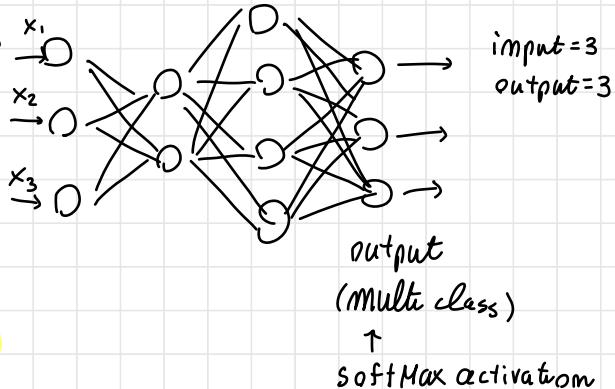
↓
you should give
different processing
with different weights

Important from one input $x_2 \rightarrow$ 3 output \rightarrow then one output \hat{y}



Multi Layer

If multiple input features



In weight initialization w_1, w_2, w_3, \dots

We initialize it with UNIFORM DISTRIBUTION

$$w_{ij} \sim \text{Uniform Distribution } (a, b) = \text{Uniform Distribution } \left[-\frac{1}{\sqrt{\text{input}}}, \frac{1}{\sqrt{\text{input}}} \right]$$

Input layer hidden layer

inside range a, b

$$a = -1/\sqrt{\text{input}}$$
$$b = 1/\sqrt{\text{input}}$$

input = # input features

e.g. if input = 3 $w_{ij} \sim \text{Uniform } \left[-\frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}} \right]$

XAVIER/Glorot INITIALIZATION

(avoid exploding gradient and proof the validity...)

Invented by Researcher → Xavier Glorot

1) Xavier Normal Initialization

$$W_{ij} \sim \mathcal{N}(0, \sigma)$$

Normal
distribution

$$\text{with } \sigma = \sqrt{\frac{2}{\text{input} + \text{output}}}$$

2) Xavier Uniform Initialization

$$W_{ij} \sim \text{Uniform Distribution } (a, b)$$

$$a = -\frac{\sqrt{6}}{\sqrt{\text{input} + \text{output}}} \quad b = \frac{\sqrt{6}}{\sqrt{\text{input} + \text{output}}}$$

KALINING HE INITIALIZATION

1) He Normal initialize.

$$W_{ij} \sim \mathcal{N}(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{\text{input}}}$$

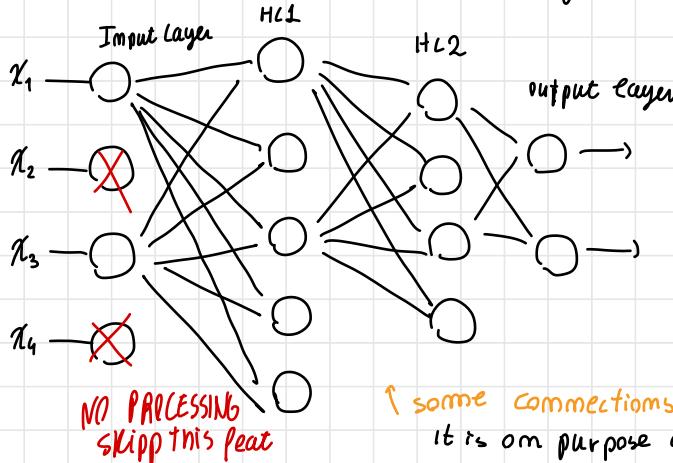
2) He Uniform initialize

$$W_{ij} \sim \text{Uniform Distribution } (a, b)$$

$$a = -\sqrt{\frac{6}{\text{input}}} \quad b = +\sqrt{\frac{6}{\text{input}}}$$

DROPOUT LAYERS

Multi layered NN



↑ some connections are missing
it is on purpose as dropout

it is related to **OVERRFITTING** ⇒ When creating a deep learning model, accurate during training...
but then in testing data it performs very BAD!

Model) Training accuracy ~90%.

Test accuracy ~60% ← training

→ the model has overfitt! we want a generalized model

We should focus ← that performs well
on generalizing the also in test data
model

⇒ In machine learning algorithms such as
Decision tree

created to its depth...

D.T. ... leads to **OVERRFITTING**..



we do pruning technique OR **Random Forest**



In Random Forest \rightarrow for every decision tree we

do Feature sampling

↑ from all dataset I take only a sub-set of features

and Raw sampling we get only some datapoints...

so In Random Forest we create subset of many decision trees built on its depth with different data subset...

↓

during construction, the output of DT is compared and the max occurrence is output... create generalized model with Random Forest

↳ similarly in a MLNN with MANY WEIGHTS... many w_{ij} , b_{ij} initialized on hidden layers / neurons...

↓

the ANN tends to overfit, if fully connected!

To address OVERFITTING you rely on DROPOUT LAYER

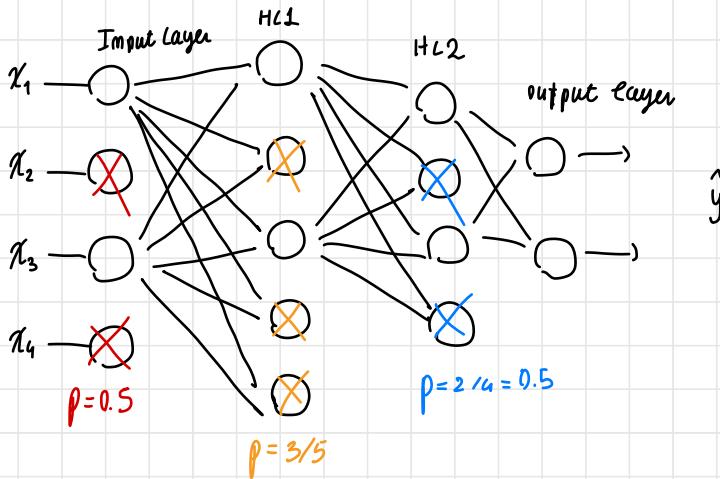
We define a DROPOUT RATE $p \approx 0.5$ (e.g.) ($0 \leq p \leq 1$)

We define the dropout probability of a specific layer (input / hidden)

e.g. $p_{\text{input}} = 0.5 \leftarrow$ Feature sampling (subset of features (input) are used)
with 50% prob. you ignore the mode..

↓

with $4 \text{ input} \sim 2$ mode are discarded during first epoch



$3/5$ of neurons are skipped \rightarrow this p value must be tuned!

during FORWARD PROPAGATION based on p some neurons are Deactivated than in Back propagation this steps happen again \rightarrow neurons may be updated

In next iteration/Epoch (Next Forward propagation) the new dropout neurons change \rightarrow randomly selected

↓
this continue until loss satisfactory and weight converges

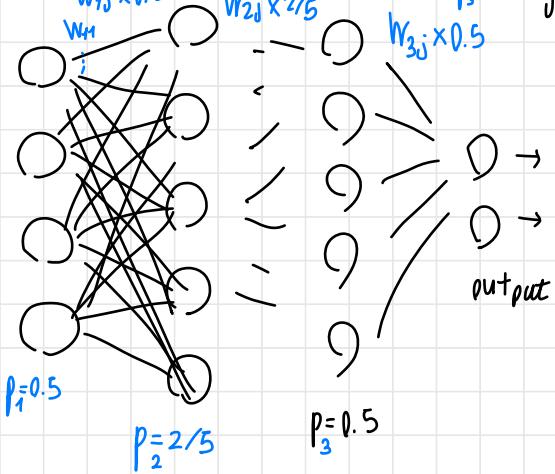
purpose: deactivate some neurons / features s.t. model does not overfit

([↑] **dropping training**) \downarrow after fixing all weights/bias

than for testing you CONNECT all (Fully connected)

(the probability of dropout $p \in [0,1]$ depends on hyperparameter)

Test data (do the prediction)



the dropout doesn't hold
anymore for testing

im test data we
don't apply dropout,
we need to predict..

instead we multiply
all weights by prob.

No deactivation...

to connect the NN
you multiply all
weights by probability
initialized before

CNN INTRODUCTION

covered ANN, Optimizers, Activation Functions... → used in CNN theory

Remember...

- ANN → supervised Learning
 - ↓

Dataset: Inputs features Output feature
 { numerical / categorical }
 { values }

↳ In

CNN

Input: Images → you can solve many different problems

◦ Image classification

↙ ◦ Object detection

◦ Segmentation

...

How does a CNN work?

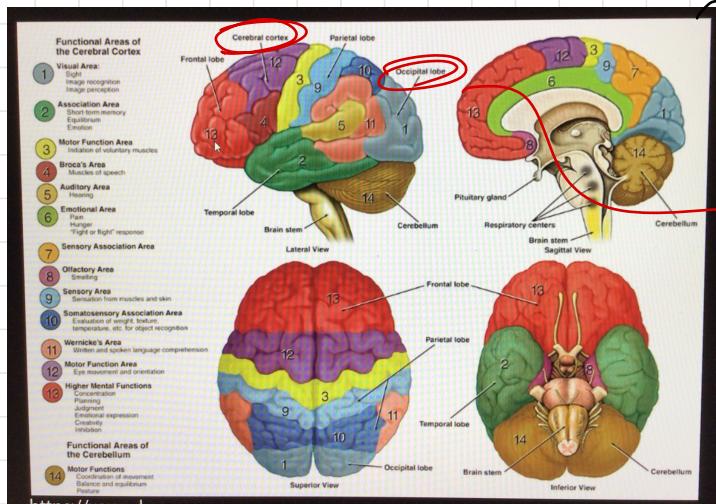
How the CNN process the

data? take images and solve complex problem statements

CNN take images as INPUTS → it's working principle is influenced
on how human brain works w.r.t images

↓

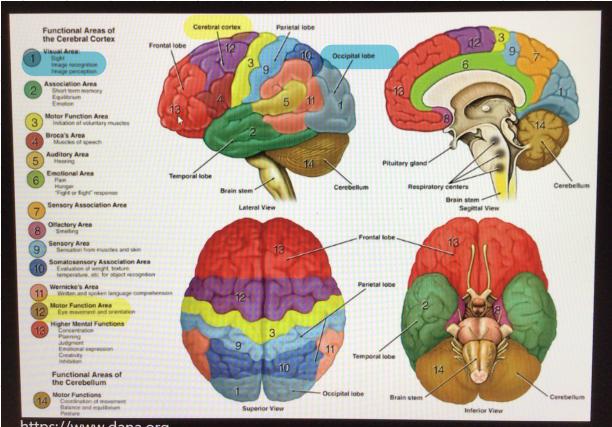
at the high level:



→ Important parts of the brain that NN try to emulate

→ for CNN we focus on

- { • cerebral cortex
- occipital lobe



← these two parts

- CEREBRAL CORTEX

- OCCIPITAL LOBE

have sub-tasks responsible
for different stimuli

1) Occipital lobe ⇒ Visual Area { sight
image recognition
image perception

12) Cerebral Cortex ⇒ Motor Function Area { Eye movement and orientation

↑
process in different layers to overall recognize pattern in images (12) take the image, (1) process it

these two parts are replicated in CNN

(12) pass signals to visual cortex → subdivide layers of image
than in (1) the image is analyzed in details.

↑

Also other brain functionalities are applied in Robotics idea

(12) the eye see something, pass signals to multiple layers in the VISUAL CORTEX

↓

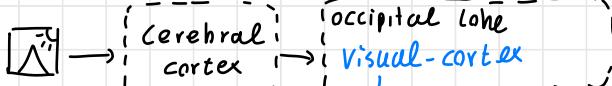
each signal is than processed and passed on (1) to extract and visualize details

here the image is visualized + processed etc!

HUMAN BRAIN and CNN

(Inside occipital lobe) ↓ this help to understand the Cerebral Cortex and Visual Cortex → Working principle of CNN

- How a human brain process a visual image?



VISUAL CORTEX (V1-V5): Region of the brain that receive, integrates and processes visual information relayed from the retinas
divided into 5 layers

↓
signal (retina)

in V1 the edges/lines are processed/analyzed

V1: primary Visual Cortex (orientation of Edges and Lines)

↑ red, green?

V2: Differences in Color, complex patterns, Object orientation

V3 V4 V5 Object Recognition

Detection
....

↖
- round
- squared
- shape..?

Visualize the image

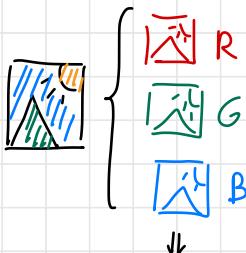
signal goes from the retinas to various layers and get processed... finally after V3-V4-V5 we visualize the image characteristics

↳ In CNN this functionality is replicated! pre-processing and perform this task in a computational way

↑
each layer of network CNN process signal using FILTERS to extract features (edges/lines...)

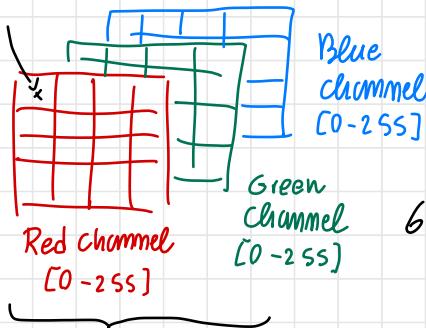
ALL YOU NEED KNOW ABOUT IMAGES

RGB Images and Grayscale Images



You can extract a gray channel (B&W) of each channel

(8 b,t)
each value
in $0 \div 255$

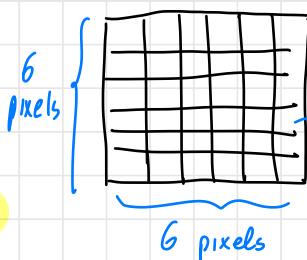


to understand how CNN works...

Images are divided in 3 channels, combined to get any colored image.

In the image representation you have

GRAYSCALE IMAGE



pixels division

6×6 with values in
 $0 \div 255$

+ you specify the # channels

$\begin{cases} 6 \times 6 \times 1 := \text{Grayscale} \\ 6 \times 6 \times 3 := \text{Colored Image} \end{cases}$

$6 \times 6 \times 3$

overall when combined, it forms the colored image

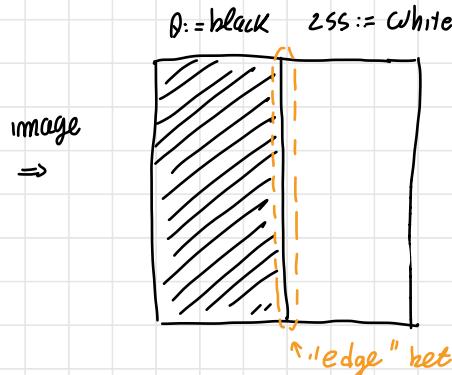
↑ different treatment by a CNN
if the image is Grayscale or Colored

CONVOLUTION OPERATION IN CNN

First type of operation in CNN

example:

6 x 6 x 1 image
(grayscale)



convolution operation
process the image
and find useful
information

"edge" between B&W

as visual cortex has many layers... every layer process some information
(V1 := edges orientation..)

The first step of convolution is pixel normalization.

STEPS

1) Normalize: convert all pixels in [0,1] range which help in the next processing steps

↓

Normalize as pixel / 255 (as maximum(um))

as dividing by 255
max pixel intensity)

0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255
0	0	0	255	255	255

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

$$\begin{array}{c} * \\ \begin{array}{|c|c|c|} \hline +1 & 0 & -1 \\ \hline +2 & 0 & -2 \\ \hline +1 & 0 & -1 \\ \hline \end{array} \end{array} = \begin{array}{|c|c|c|} \hline 0 & -4 & -4 & Q \\ \hline 0 & -4 & -4 & Q \\ \hline P & -4 & -4 & Q \\ \hline Q & -4 & -4 & Q \\ \hline \end{array}$$

FILTER
 3×3

Step 2: convolution operation

as the visual cortex processes information in each layer (v1, v2...) in the same way CNN layered process occurs through FILTERS

Its main objective is ← it can be of different sizes
to take out information from the image
↓ (edges horiz./vert., vertical plane)

$6 \times 6 \times 1 \times 3 \times 3 \times 1$ → 4×4 output by convolution operation

stride during convolution we take the 3×3 filter and slide it over top of $6 \times 6 \times 1$ image

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

$$* \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 2 & 0 & -2 \\ \hline 1 & 0 & -1 \\ \hline \end{array} \rightarrow$$

convolution 1
centered on
first 3×3 block

$$1 \cdot 0 + 0 \cdot 0 - 1 \cdot 0 \\ + 2 \cdot 0 + 0 \cdot 0 - 2 \cdot 0 \\ + 1 \cdot 0 + 0 \cdot 0 - 1 \cdot 0 = 0$$

element wise
product and
summation..

0	-4	-4	0
0	-4	-4	0
Q	-4	-4	0
Q	-4	-4	0

next convolution 2 step you jump by one on the right

and multiply for the same filter

$$1 \cdot 0 + 0 \cdot 0 - 1 \cdot 1 + 0 \cdot 0 + 0 \cdot 2 - 2 \cdot 1 \\ + 0 \cdot 1 + 0 \cdot 0 - 1 \cdot 1 = -4$$

"stride" is a jump to right/left/bottom here STRIDE=1

another stride 2 on step 3

then final stride on first "row" of the output and for next step I go to next row (No less right stride)

What does this output represents??

normalized image

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

*

+1	0	-1
+2	0	-2
+1	0	-1

FILTER
3x3

=

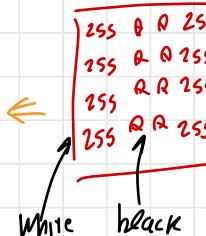
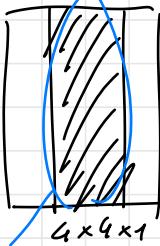
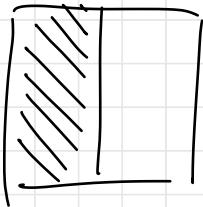
0	-4	-4	Q
0	-4	-4	Q
Q	-4	-4	Q
Q	-4	-4	Q

↓ de-normalization
based on min/max
here, $0 := 255$ max
 $-4 := 0$ min

original format

(NOT computed in CNN,
Here just for visual
purpose)

6x6x1



this Filter 3x3 has found out
the VERTICAL EDGE in original image "edge detection"

+1	0	-1
+2	0	-2
+1	0	-1

vertical
edge
filter

→ in VISUAL CORTEX V1 layer has
the responsibility on edges orientation
↓

in the same way in CONVOLUTION
OPERATION edges are extracted
by FILTERS.

→ this behave as the V1 layer
similarly I can have any # of filters

- Round edge
- Horizontal edge
- ...

e.g Horiz Edge Filter

+1	+2	+1
0	0	0
-1	-2	-1

↖ extract horiz.

Here these filters are pre-defined

Usually in a CNN we don't define filter values... we randomly initialize these values

Then through Forward and Backward propagation, the FILTERS values are updated

With different filter size and values... computed by different CNN techniques, e.g. in transfer learning)

So... a convolution operation pass a filter through an image and you get an output that highlights image informations
(vertical/horiz/round edges etc.)

→ Notice that from 6×6 image size $\{m=6\}$
and filter size f is $3 \times 3 \quad \{f=3\}$
while output size $n \quad \{n=4\}$



Usually a formula is valid on this input/output size

$$m - f + 1 = n \quad (6 - 3 + 1 = 4)$$

Be careful from 6×6 image I get a 4×4 ...

Loss of information! (less pixels/resolution)



It is possible to prevent this loss of information during filtering by applying PADDING.

Convolution is used to process the image and find patterns, that will help the model in its task

PADDING IN CNN

NORMALIZED IMAGE

$$\begin{bmatrix} 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix} \quad (6 \times 6) \quad n=6$$

(3x3)

FILTER

$$\begin{bmatrix} +1 & 0 & -1 \\ +2 & 0 & -2 \\ +1 & 0 & -1 \end{bmatrix} \quad f=3$$

=

(4x4)

FILTERED IMAGE

$$\begin{bmatrix} 0 & -4 & -4 & 0 \\ 0 & -4 & -4 & 0 \\ 0 & -4 & -4 & 0 \\ 0 & -4 & -4 & 0 \end{bmatrix}$$

$$n - f + 1 = 4 \text{ (output)}$$

here from 6×6 to 4×4 we lose some amount of informations...



How to prevent it? to get as output same size as input we apply padding

padding := adding additional layer

zero padding) here one layer
neighbour padding you avoid loss of data by PADDING in CNN

ZERO padding

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

FILTER

$$\begin{bmatrix} +1 & 0 & -1 \\ +2 & 0 & -2 \\ +1 & 0 & -1 \end{bmatrix} \quad 3 \times 3$$

→ 6x6 output

$$\begin{bmatrix} 0 & -4 & -4 & 0 \\ 0 & -4 & -4 & 0 \\ 0 & -4 & -4 & 0 \\ 0 & -4 & -4 & 0 \end{bmatrix}$$

$$6 \times 6 \rightarrow 8 \times 8$$

how much padding I need to apply?

s.t. at output you have same size
as input ↓

$$n - f + 2p + 1 = \text{output (desired)}$$

$$6 - 3 + 2p + 1 = 6 \rightarrow 2p = 2 \rightarrow p=1$$

padding of 1 layer is necessary

↳ what values you use in the padding??

1) zero padding := apply 0 in all the additional layer

2) neighbour padding := based on the close pixels, you use same value

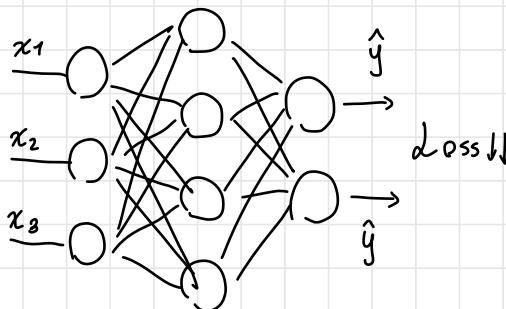
level of padding is computed based on desired output size
+ you need to decide also padding type (zero/ neighbour)

e.g. input $m=7$ filter $f=3$ output desired $q=7$

$(7 \times 7) \quad (3 \times 3) \quad (7 \times 7)$

$$M - f + 2p + 1 = q \Rightarrow 7 - 3 + 2p + 1 = 7 \quad 2p = 2 \rightarrow p=1$$

OPERATION OF CNN vs ANN



In ANN (multi layered NN)
← 2 layered NN
(1 input, 1 HL, 1 output)

$$\begin{aligned} z &= W^T x_i + b \\ \text{Act}(z) &\text{ (e.g. ReLU(z))} \\ \text{then backward, etc...} \end{aligned}$$

Similarly in a CNN

NORMALIZED IMAGE						
0	0	0	1	1	1	1
0	0	0	1	1	1	1
0	0	0	1	1	1	1
0	0	0	1	1	1	1
0	0	0	1	1	1	1
0	0	0	1	1	1	1

(vertical edge)

FILTER f_1

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

f_2



FILTERED IMAGE

$$\begin{bmatrix} 0 & -4 & -4 & 0 \\ 0 & -4 & -4 & 0 \\ 0 & -4 & -4 & 0 \\ 0 & -4 & -4 & 0 \end{bmatrix}$$

This is a
(SINGLE)

CONVOLUTION
LAYER

↓
we can add more
convolution layer
in sequence or
the CNN (horizontally)

→ when training CNN from
scratch this filters are

I can have multiple filters...

Randomly
initialized...

Image pass through the Filter...

as in the VISUAL CORTEX layered, processing of information
coming from the Retina.

(Random)

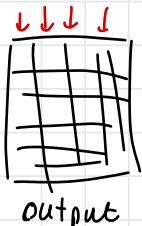
The Filter is applied and my output w.r.t f_1 is computed

multiple filters can be applied to extract patterns (shape/edges...)

each filter plays an important role in understanding the image \Rightarrow the initialized filters (Randomly) during training initialization

then with forward/backward propagation we update the filters values...

$$\text{Image} \times \text{filter} =$$



↑↑↑↑ to all values ReLU is applied
 $\max(0, x) \rightarrow$ to transform the values in the filtered image

\Rightarrow when backpropagating, filter values must be updated ... to update it the derivative should be computed...

(make it differentiable)

↳ easily update filter weights

\leftarrow (+ avoid vanishing gradient with ReLU)

same training / learning approach as ANN

Forward: we apply filters on input image

$$\begin{array}{l} \text{Image } f_1 f_1 \times f_1 \\ \text{mxm} \times f_2 f_2 \times f_3 \\ \vdots \\ f_n f_n \times f_n \end{array} = \begin{array}{ll} \text{out}_1 & m - f_1 + 2p_1 + 1 \\ \text{out}_2 & m - f_2 + 2p_2 + 1 \\ \vdots & \vdots \\ \text{out}_n & m - f_n + 2p_n + 1 \end{array}$$

\uparrow on all values of output image we apply ReLU

because on backpropagation

we should be able to find

the differentiation is required
so that the filters can be
computed (learn parameters)

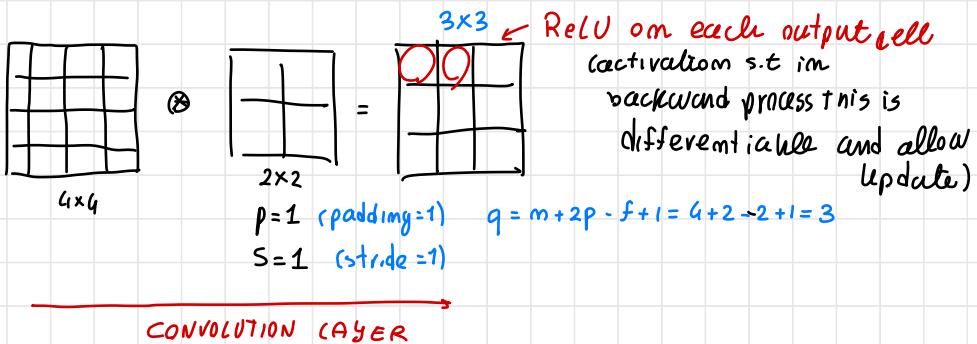
\leftarrow the derivative...

after convolution, ReLU is required (or other variation)
so that it is derivable...
then other operations are applied after ReLU

allow the learning of filter based on the problem task,
profile output for classification/detection...

after convolution, ReLU \rightarrow max pooling is applied
(average)

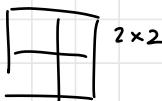
MAX POOLING, MIN POOLING, MEAN POOLING



IF after convolution

max pooling?
choose a filter size

ReLU output		
1	2	3
4	5	G
2	8	4



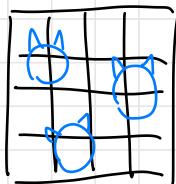
I apply Max pooling
on top of the FILTER

this filter is applied on top of the output

max pooling := applied for LOCATION INVARIANT

LOCATION INVARIANT

Since we use multiple FILTERS to determine different edges/shapes...



↙ I have 3 cat img \Rightarrow location invariant

Any number of cat imgs present...

IF FILTER \Rightarrow responsible to find cat faces
↓

When Apply max pooling on output layer I can extract faces much clearly!

MAX POOLING

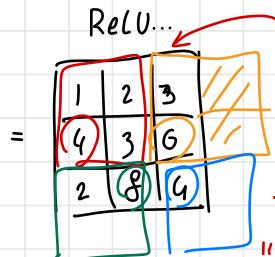
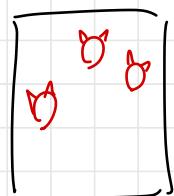
It concentrate completely on the part of interest to get output

allow to focus on the most important informations irrespective of the particular filter

e.g. multiple cars in image... given a FILTER that extract the wheels then max pooling focus on cars wheels to take out the specific information

+ LOCATION INVARIANT \rightarrow whenever those img shapes/features are present, max pooling can detect those parts

overall max pooling act in this way:



max pooling act as in the convolutions...

"MAX" pooling since it extract the max value in the filtersize stride of 2!

then make stride = 2, since it take max intensity (NOT all pixels)

4	6
8	4

1	2	3
4	3	6
2	8	9

max pooling



→

4	6
8	6

highest intensity (most important feature) is extracted... ↗

↓
output is more focused on the part of interest

In CNN architectures: convolutional layer + max pooling
are stacked horizontally multiple times...

Input → convolution → max pool → convolution → max pool ...

applied to all
output of the filtered image

other pooling techniques are MIN / MEAN POOLING

interest in the
minimum
intensity pixel
available

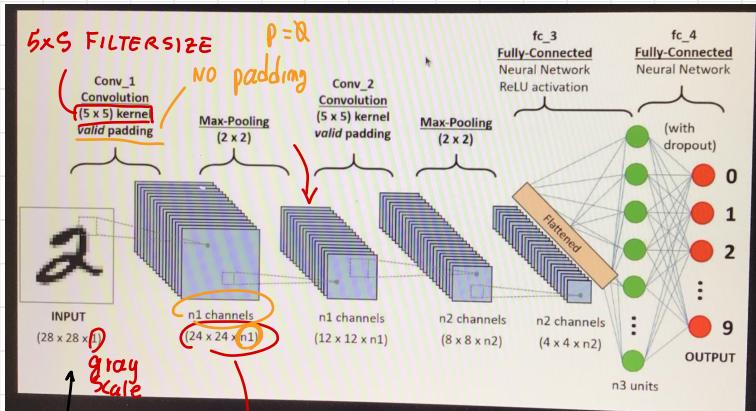
take the
minimum
element

Average of
all the
pixel values

After all this CNN architecture combine convolution + max pooling
to extract information

+ place as final layer a fully connected
layer on the CNN output

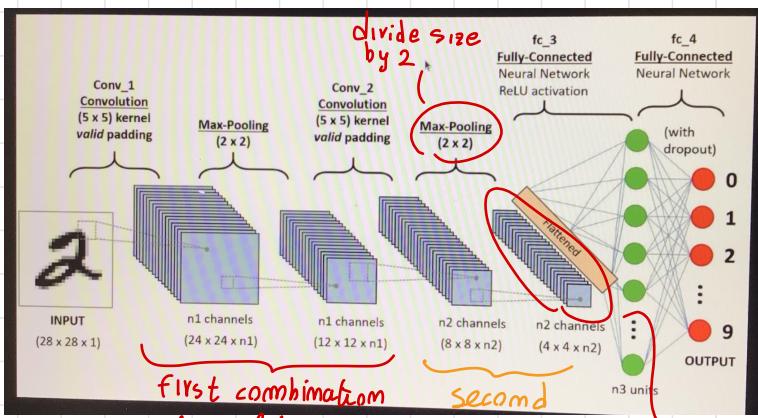
FULLY CONNECTED LAYER IN CNN



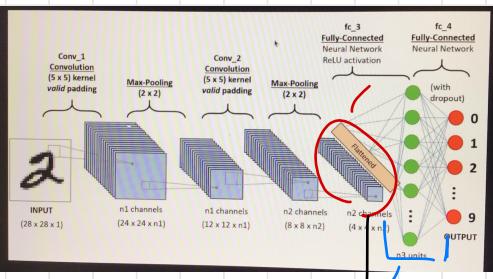
After first convolution layer ...
 With n_1 channels meaning the amount of filters applied on the image
 $28-5+1 = 24$ out of Convolution layer

divide size by 2 } applied to all
 max pooling 2×2 } applied to all
 and reduce img } channels of $1 \dots n_1$
 size → focus on location invariant...
 extract more intensity path

conv := filter + ReLU
 ... applied horizontally in sequence



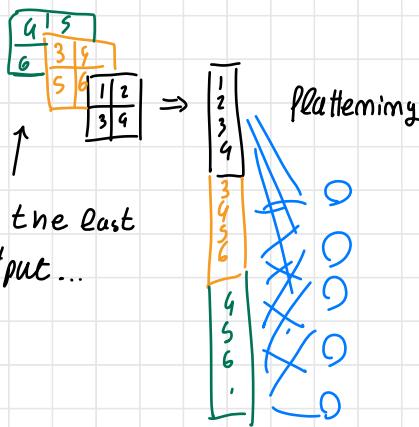
after all conv + max pool, sequence the FLATTENING LAYER is applied



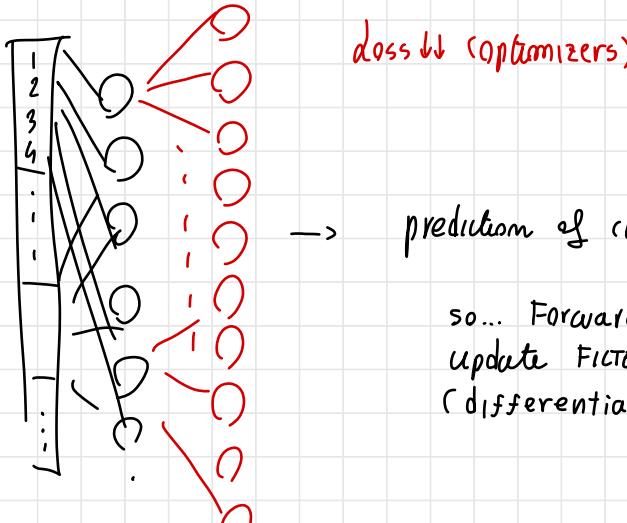
than the last piece is a fully connected layer
(the same as ANN)

than, depend on the problem
e.g. in MNIST you have 0÷9 (10 outputs)
classes

FLATTENED Layer



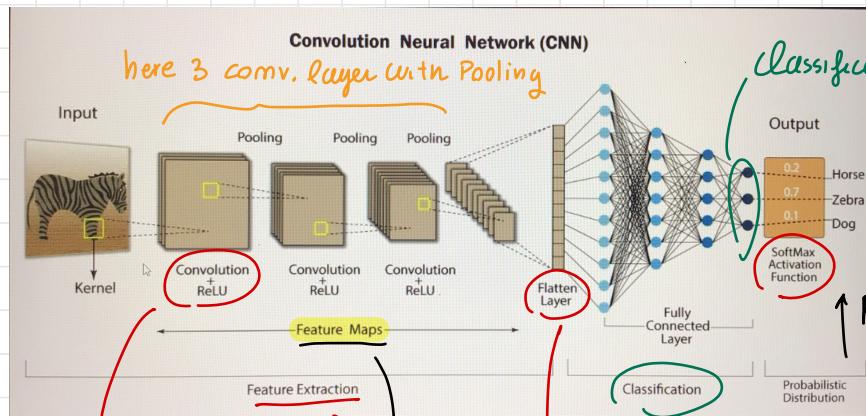
fully connected layer



→ prediction of class

so... Forward + backward
update FILTERS VALUES
(differentiable thanks to ReLU activ.)

CNN With RGB



classification of 3 classes

probability in output

one convolution layer

define the process
of feature

the entire process is the FEATURE EXTRACTION

from flattening
to MLP apply
classification

In Forward propag... we pass through the cell Network...
then in backward you should be able to derivate... why you use ReLU

then loss $(\hat{y} - y)^2$... backpropagate + optimizers...

update FILTERS
weight and ANN weights.