

MAI

Deep Learning

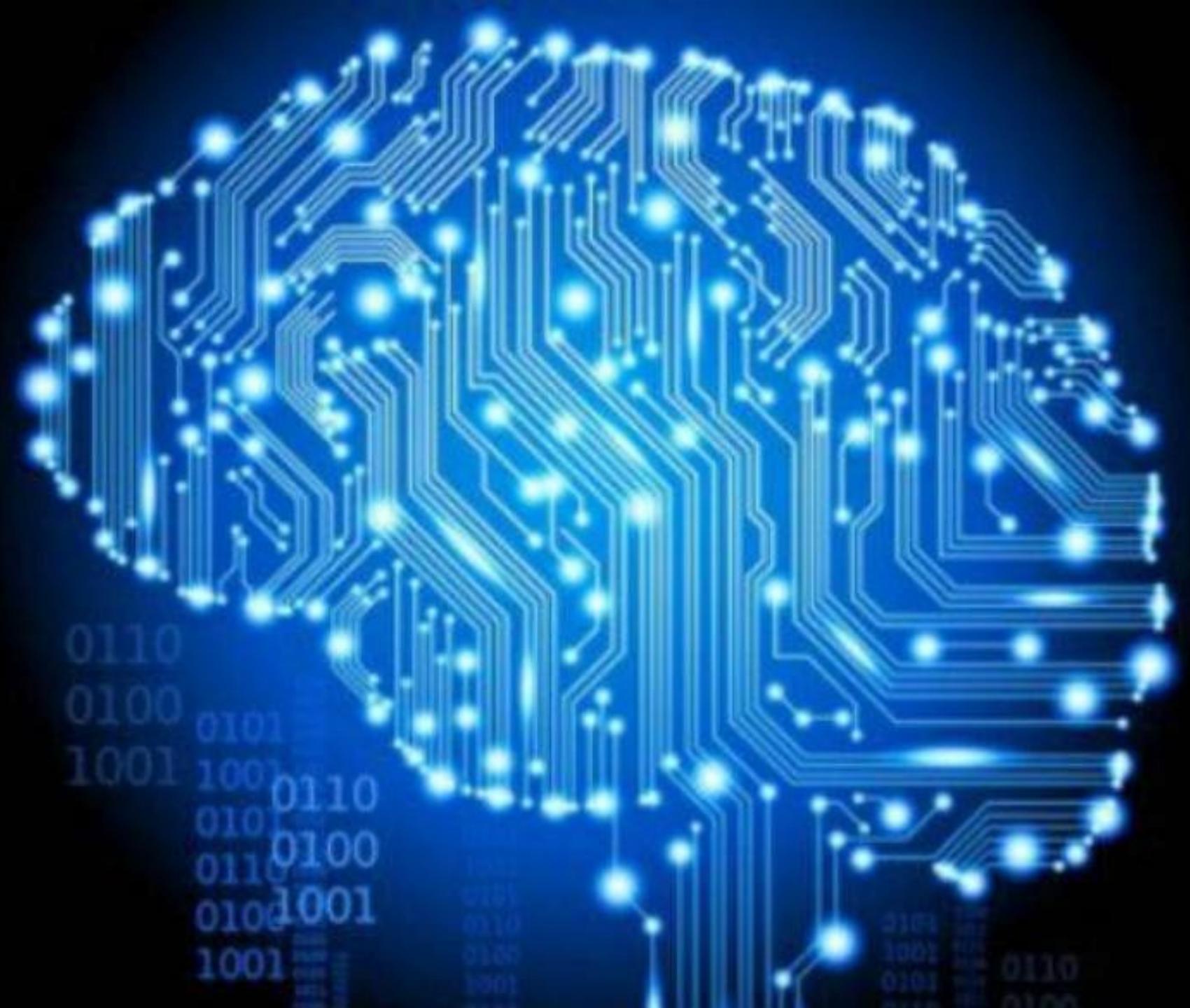


THEORY

Feed-forward and convolutional neural networks

Dario Garcia Gasulla
dario.garcia@bsc.es

A bit of History



Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH

A bit of History

McCulloch & Pitts / Hebb
Rosenblatt's Perceptron
Minsky & Papert - XOR



[McCulloch, 43]

- 1943 Warren McCulloch & Walter Pitts:
 - How To: From neurons to complex thought
 - Binary threshold activations
 -
 -
 -
 -
- 1949 Howard Hebb:
 - Neurons that fire together wire together
 - Weights: Learning and memory

A Logical Calculus of Ideas Immanent in Nervous Activity

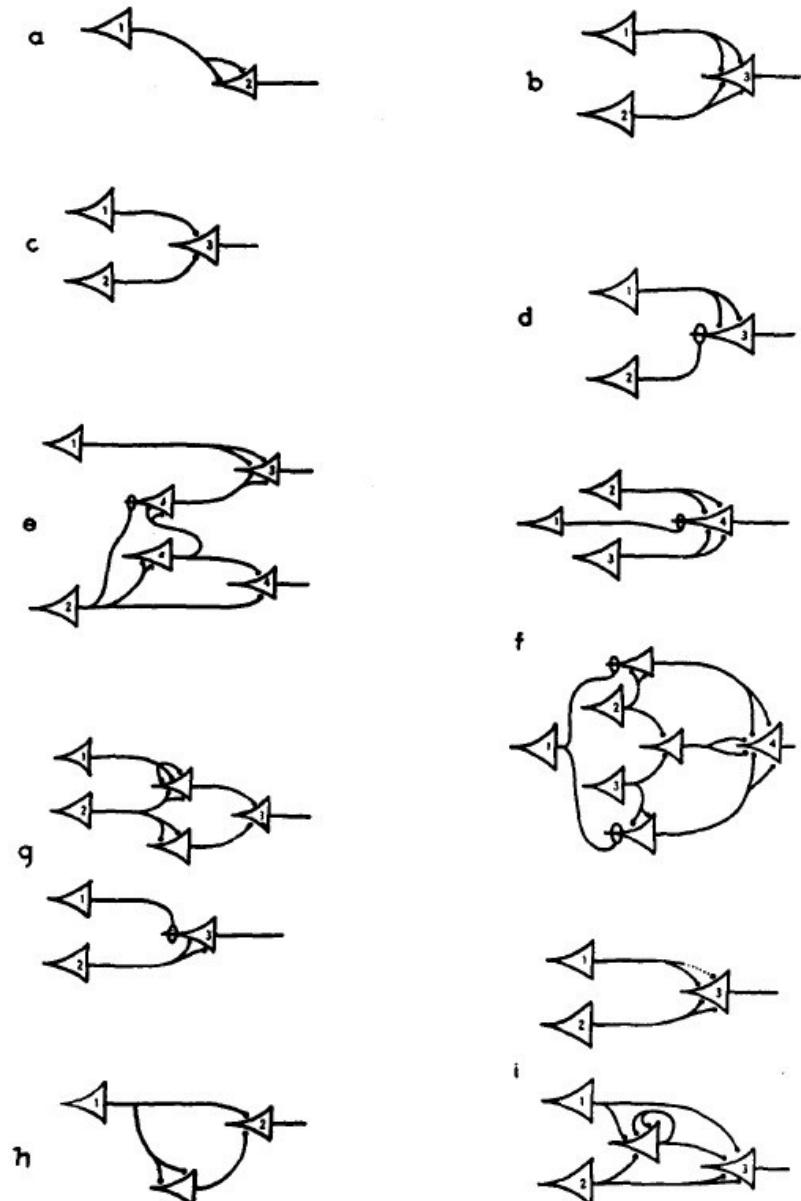


FIGURE 1

A bit of History

McCulloch &
Pitts / Hebb
Rosenblatt's
Perceptron
Minsky & Papert
- XOR
Backpropagation
Algorithm



HIGH PERFORMANCE
ARTIFICIAL INTELLIGENCE

[Rosenblatt, 58]
[Mark I Perceptron]
[Perceptrons]
[1]

1948, Rosenblatt applied
Hebb's learning to
McCulloch & Pitts design

$$f(x) = \begin{cases} 1 & \text{if } w \cdot x + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

- w real-valued weights
- dot product
- b real scalar constant

The Mark I Perceptron. A visual classifier with:

- 400 photosensitive receptors (sensory units)
- 512 stepping motors (association units, trainable)
- 8 output neurons (response units)

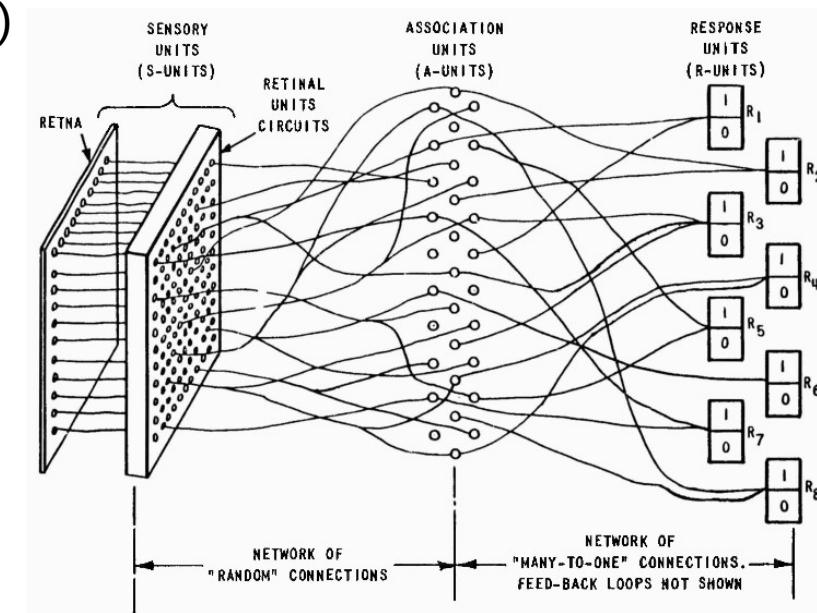


Figure 1 ORGANIZATION OF THE MARK I PERCEPTRON

A bit of History

McCulloch &
Pitts / Hebb
Rosenblatt's
Perceptron
Minsky & Papert
- XOR
Backpropagation
Algorithm

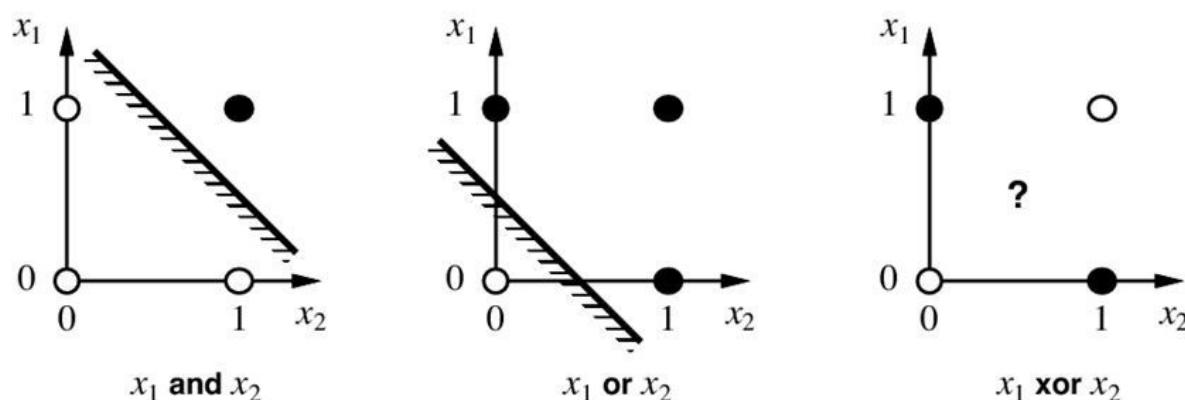


[Minsky, 69]

Rosenblatt acknowledged a set of limitations in the Perceptron machine.

Minsky & Papert did too in “*Perceptrons: an introduction to computational geometry*”, including:

- A multilayer perceptron (MLP) is needed for learning basic functions like XOR
- MLP cannot be trained.



This had a huge impact on the public, resulting in a drastic cut in funding of NNs until the mid 80s

1st AI WINTER

A bit of History

McCulloch &
Pitts / Hebb
Rosenblatt's
Perceptron
Minsky & Papert
- XOR

Backpropagation

Algorithm



[Werbos, 74]
[Rumelhard, 85]

How can we optimize neuron weights which are not directly connected to the error measure?

Backpropagation algorithm:

Use the chain rule to find the derivative of cost with respect to any variable.

In other words, find the contribution of each weight to the overall error.

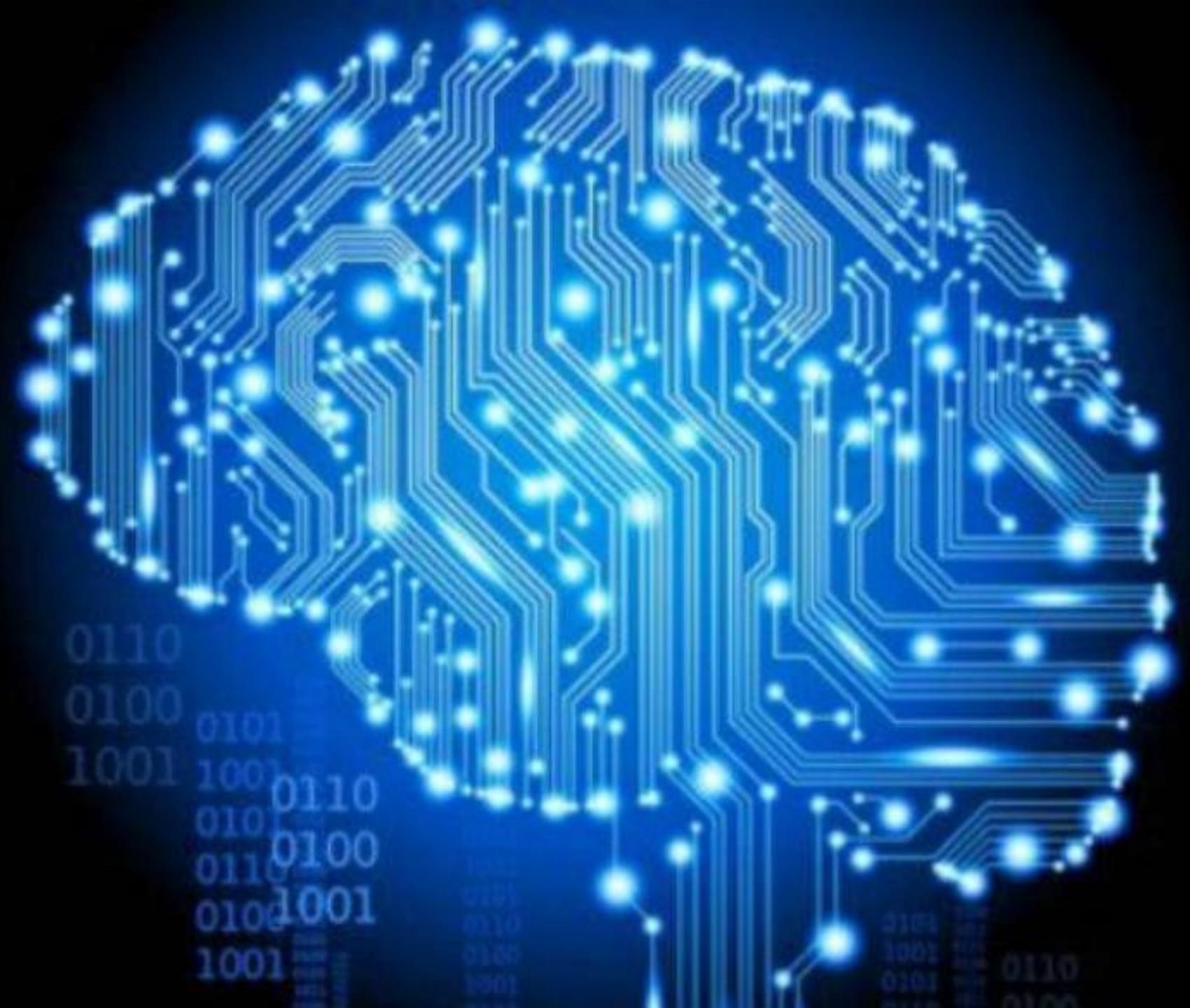
First proposed for training MLPs by Werbos in '74.
Rediscovered by Rumelhart, Hinton and Williams in '85.

End of NNs Winter

Training with backprop

1. Forward pass from input to output
2. Error measurement (loss function)
3. Find gradients towards minimizing error layer by layer
(backward pass)

Feedforward Neural Networks



Feedforward Neural Networks

**SGD, Epochs,
Batches and Steps**
Activation functions
SGD learning rate
**Other optimization
methods**
Regularization
Normalizing inputs
**Vanishing/Exploding
Gradients**
Weights initialization



Computing the gradients using all available training data would require huge amounts of memory.

Stochastic Gradient Descent: Iteratively update weights using random samples (hence, *stochastic*)

-
-

Each feedforward/backward cycle (a **step**) processes a random **batch** of images.

- Typical batch sizes: Powers of 2.
- Batch size = 1 --> Full stochastic (slower)
- Batch size = dataset_size --> Deterministic (bad generalization)

An **epoch** is the processing of the whole dataset once. It corresponds to processing as many batches as:

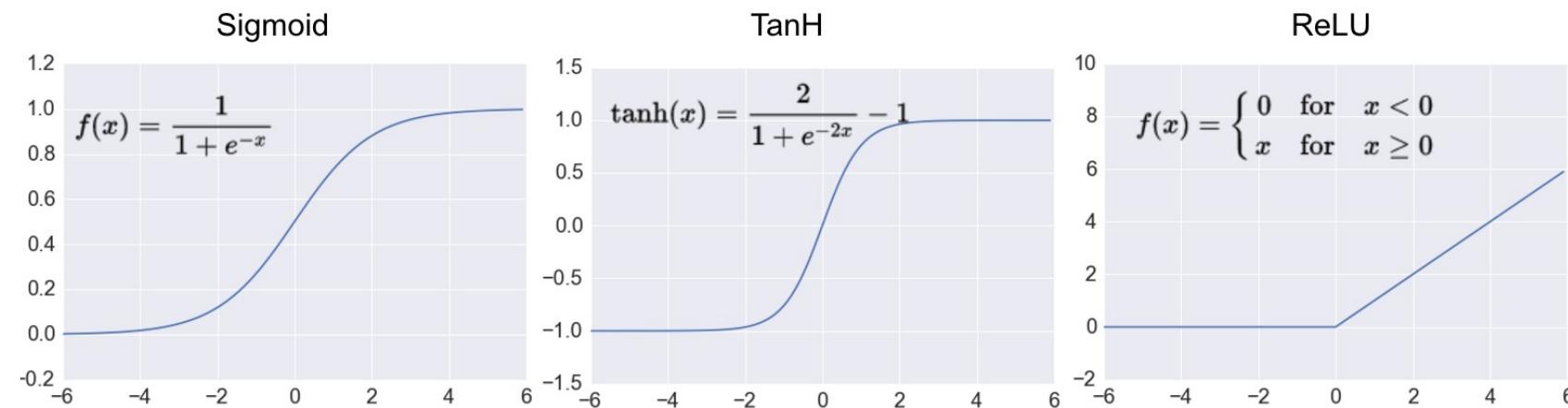
$$\text{dataset_size} / \text{batch_size}$$

Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization methods
Regularization
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization



Activation functions transform the output of a layer to a given range. If the function is non-linear, the net can learn non-linear patterns (e.g., XOR).



- Zero gradient in most of $f(x)$. Saturates!
- Max gradient is 0.25 or 1. Vanishing!
- Does not saturate
- Does not vanish
- Faster
- May die

ReLU is a safe choice in most cases
Undying alternatives: Leaky ReLU, ELU, ...

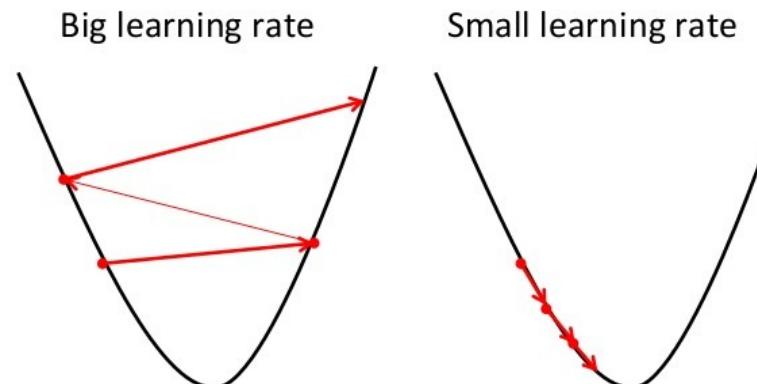
Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization methods
Regularization
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization



[Dauphin, 14]

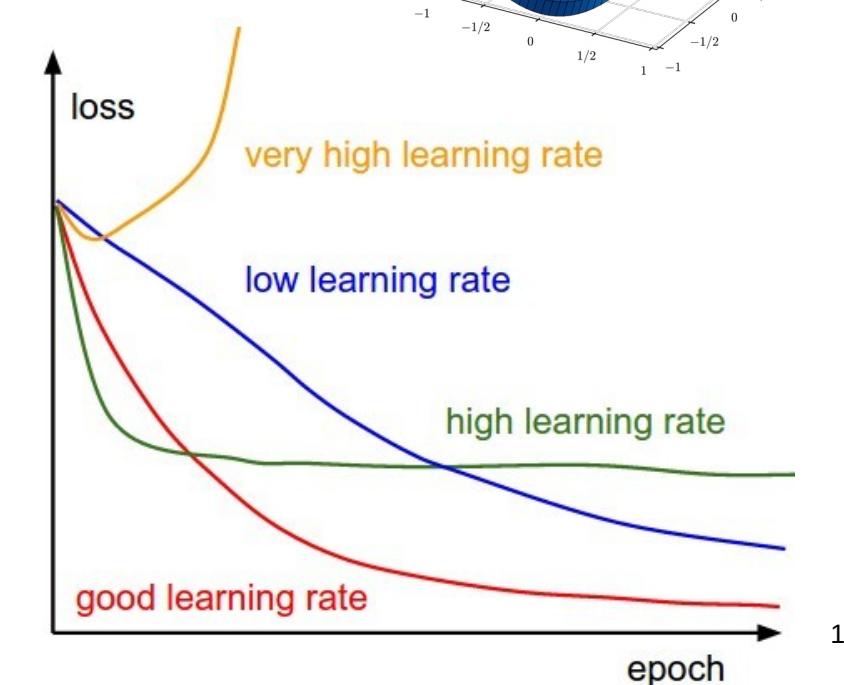
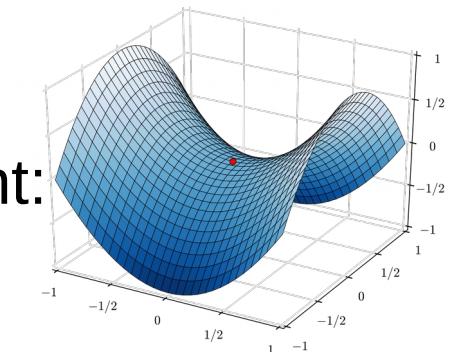
SGD will overshoot unless we keep decreasing the LR.



Gradient descent is an optimization algorithm to update weights towards a min.

Learning rate determines how much we move in that direction. With the wrong LR you may end up in local minima, a saddle point, or be too slow.

Saddle point:



Feedforward Neural Networks

SGD, Epochs,
Batches and Steps

Activation functions

SGD learning rate

Other optimization
methods

Regularization

Normalizing inputs

Vanishing/Exploding
Gradients

Weights initialization



Momentum: Include a fraction of the previous gradient.
Keeps the general direction so far.

Nesterov: Compute current gradient considering where
the previous gradient took you. (RNNs?)

Adagrad: Parameter-wise LR considering past updates.
Good for infrequent patterns. Vanishing LR due to growing
history.

Adadelta: Adagrad with a decaying average over history.
Typically set around 0.9.

Adam: Adadelta + Momentum

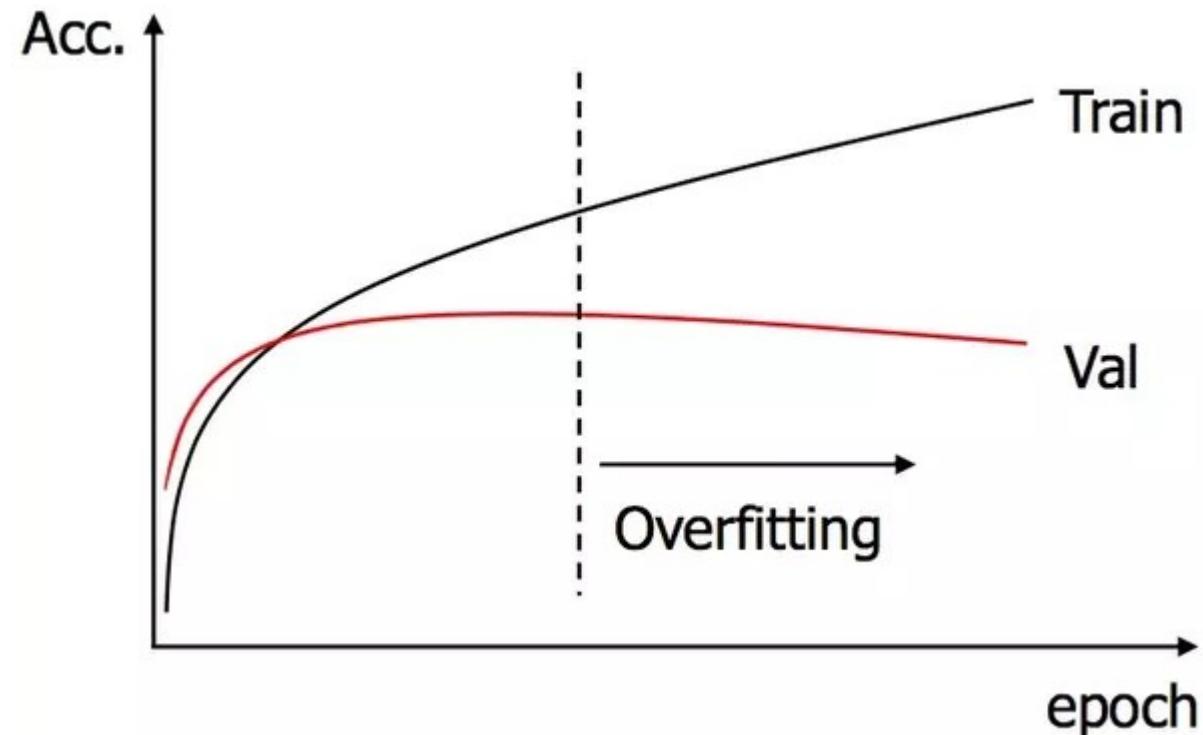
Feedforward Neural Networks

- SGD, Epochs,
- Batches and Steps
- Activation functions
- SGD learning rate
- Other optimization methods
- Regularization
- Normalizing inputs
- Vanishing/Exploding Gradients
- Weights initialization



Why do we need regularization?

Because the difference between **Machine Learning** and **Optimization** is called **Generalization**



Feedforward Neural Networks

Generalization

Polynomial regression

1 $h(x) = w_1x + b$

2 $h(x) = w_3x^3 + w_2x^2 + w_1x + b$

3 $h(x) = w_{14}x^{14} + w_{13}x^{13} + \dots + w_1x + b$

SGD, Epochs, Batches and Steps

Activation functions

SGD learning rate

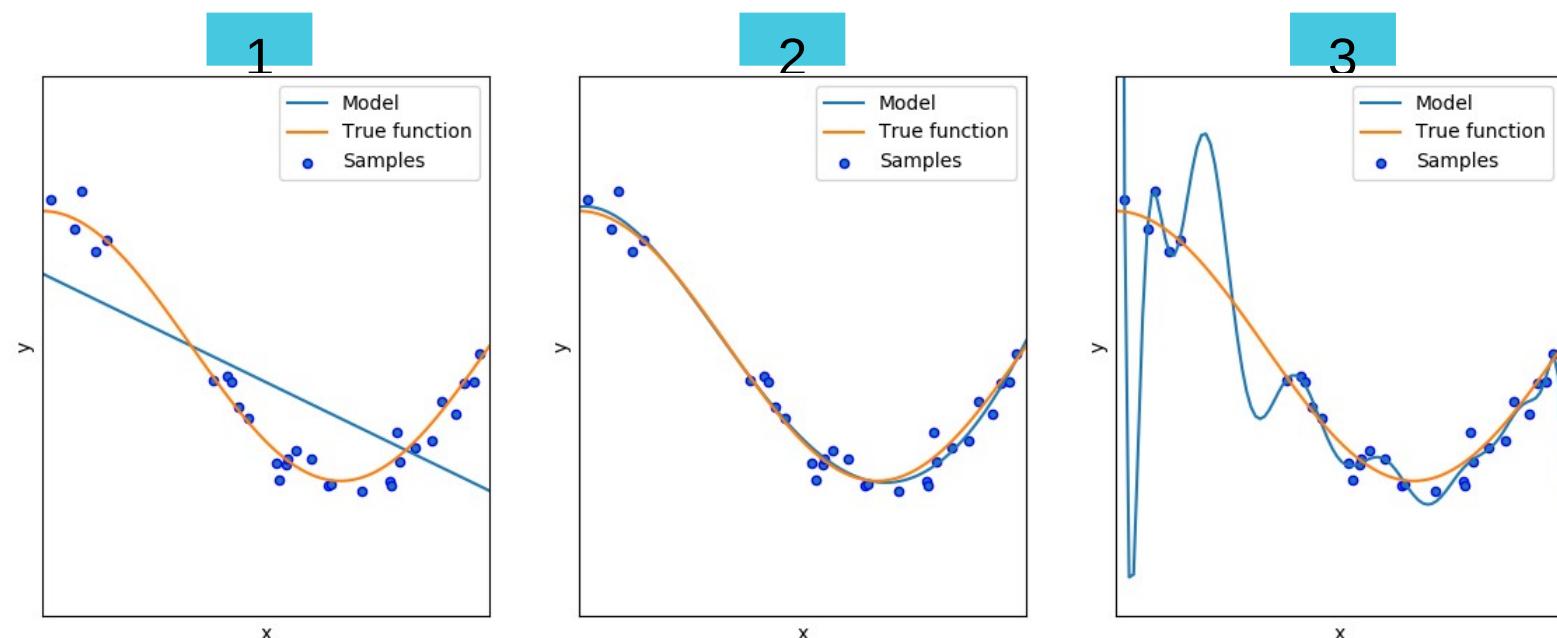
Other optimization methods

Regularization

Normalizing inputs

Vanishing/Exploding Gradients

Weights initialization



Feedforward Neural Networks

SGD: Epochs, Batches and Steps
SGD: Gradient Descent
SGD: Learning rate
Other optimization methods
Adam
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

Generalization

Polynomial regression

Huge

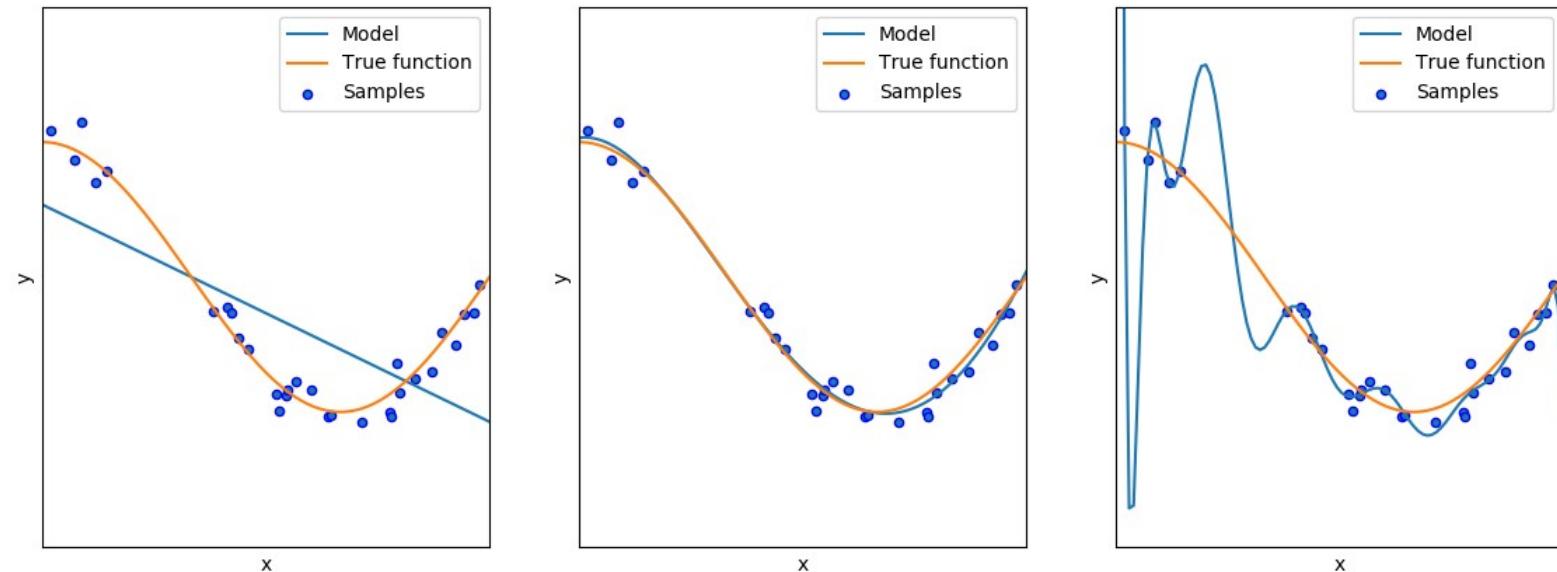
Small

Almost Zero

Bad

Good

Horrible



Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Learning Curves
SGD Learning Rate
Other optimization methods
Adam
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

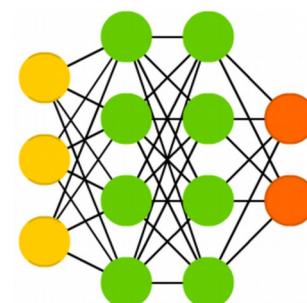
Generalization

What **policy** can we use to **improve** model generalization?



Occam's Razor

when you have **two competing hypotheses** that make the **same predictions**, the **simpler one is the better**



Machine Learning

given **two models** that have a **similar performance**, It's better to **choose the simpler one**

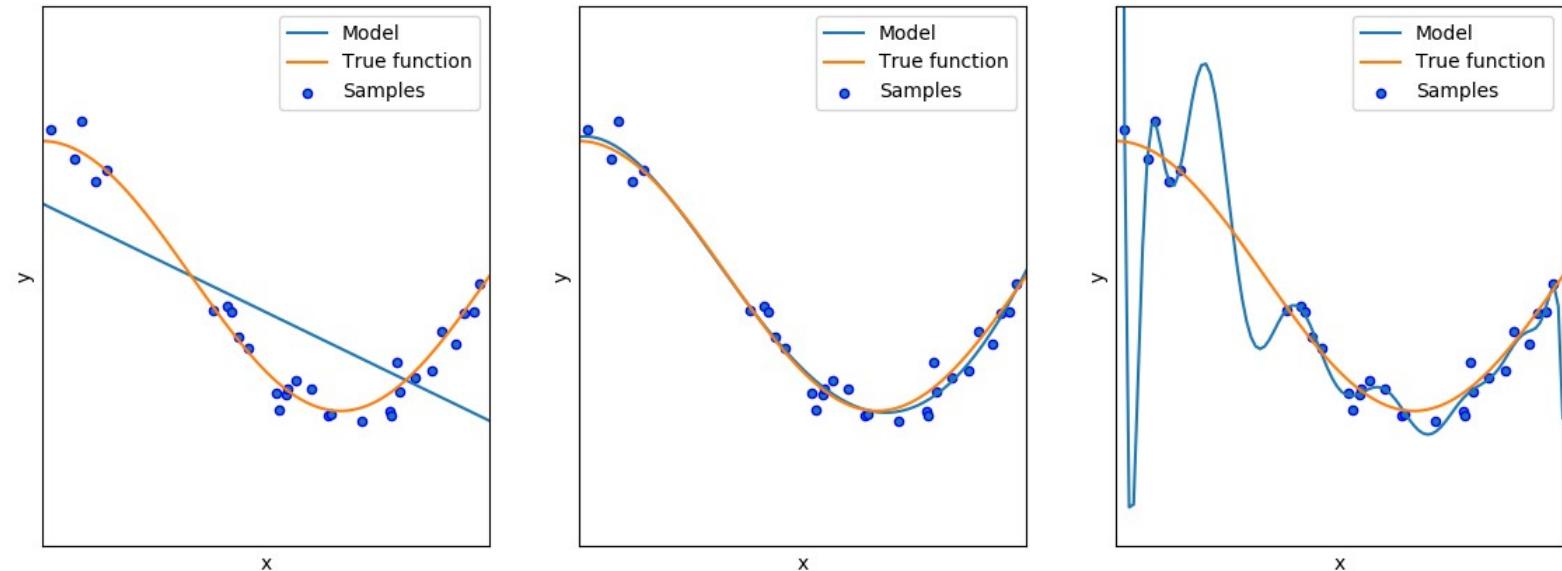
Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Activation functions
ReLU and Softmax
Other optimization methods
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

Model Complexity

What **policy** can we use to **improve** model generalization?

Cost function = Training Error + Model Complexity



Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Activation functions
ReLU implementation
Other optimization methods
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

Model Complexity

$$h(x) = w_3x^3 + w_2x^2 + w_1x + w_0 \quad \text{vs} \quad h(x) = 0x^3 + 0x^2 + w_1x + w_0$$



Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Activation functions
ReLU and Softmax
Other optimization methods
Adam
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

Model Complexity

$$h(x) = w_3x^3 + w_2x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + 0x^2 + w_1x + w_0$$

$$h(x) = w_3x^3 + w_2x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + w_2x^2 + 0x + w_0$$

Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Activation functions
ReLU implementation
Other optimization methods
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

Model Complexity

$$h(x) = w_3x^3 + w_2x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + 0x^2 + w_1x + w_0$$

$$h(x) = w_3x^3 + w_2x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + w_2x^2 + 0x + w_0$$

$$h(x) = 0x^3 + 0x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + w_2x^2 + 0x + w_0$$

Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Activation functions
ReLU and Softmax
Other optimization methods
Adam
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

Model Complexity

$$h(x) = w_3x^3 + w_2x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + 0x^2 + w_1x + w_0$$

$$h(x) = w_3x^3 + w_2x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + w_2x^2 + 0x + w_0$$

$$h(x) = 0x^3 + 0x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + w_2x^2 + 0x + w_0$$

$$h(x) = 0x^3 + 0x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + w_2x^2 + 0x + 0$$

Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Activation functions
ReLU and Softmax
Other optimization methods
Adam
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

Model Complexity

$$h(x) = w_3x^3 + w_2x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + 0x^2 + w_1x + w_0$$

$$h(x) = w_3x^3 + w_2x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + w_2x^2 + 0x + w_0$$

$$h(x) = 0x^3 + 0x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + w_2x^2 + 0x + w_0$$

$$h(x) = 0x^3 + 0x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + w_2x^2 + 0x + 0$$



?

?

?

Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Activation functions
ReLU and Softmax
Other optimization methods
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

Model Complexity

$$h(x) = 0x^3 + 0x^2 + w_1x + w_0$$



$$h(x) = 0x^3 + w_2x^2 + 0x + 0$$

ℓ_0 complexity: Number of non-zero coefficients

ℓ_1 "lasso" complexity: $\sum_{i=0}^d |w_i|$, for coefficients w_0, \dots, w_d

ℓ_2 "ridge" complexity: $\sum_{i=0}^d w_i^2$, for coefficients w_0, \dots, w_d

Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Activation functions
ReLU and Softmax
Other optimization methods
Adam
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

Model Complexity

$$h(x) = 0x^3 + 0x^2 + w_1x + w_0 \quad \text{vs} \quad h(x) = 0x^3 + w_2x^2 + 0x + 0$$

$$w_0 = 1.3 \quad w_1 = -1.2 \quad w_2 = 2.2$$

ℓ_0 complexity

$$|\{w_1, w_0\}| = 2$$



$$|\{w_2\}| = 1$$

ℓ_1 complexity

$$|1.3| + |-1.2| = 2.5$$



$$|2.2| = 2.2$$

ℓ_2 complexity

$$1.3^2 + (-1.2)^2 = 3.13$$



$$2.2^2 = 4.84$$

Feedforward Neural Networks



L1 / L2 Regularization

$$\text{Cost function} = \text{Loss} + \frac{\lambda}{m} \sum_{i=0}^m |w_i|$$

$$\text{Cost function} = \text{Loss} + \frac{\lambda}{2m} \sum_{i=0}^m w_i^2$$

Regularization parameter $\rightarrow \lambda$

What **complexities** do these methods use?

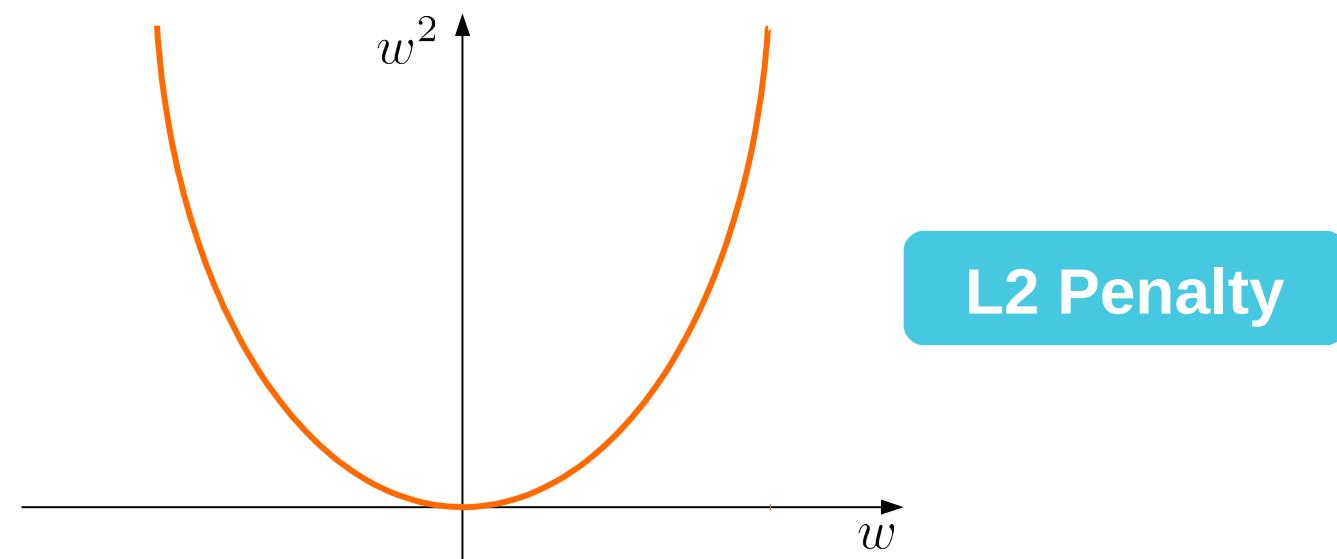
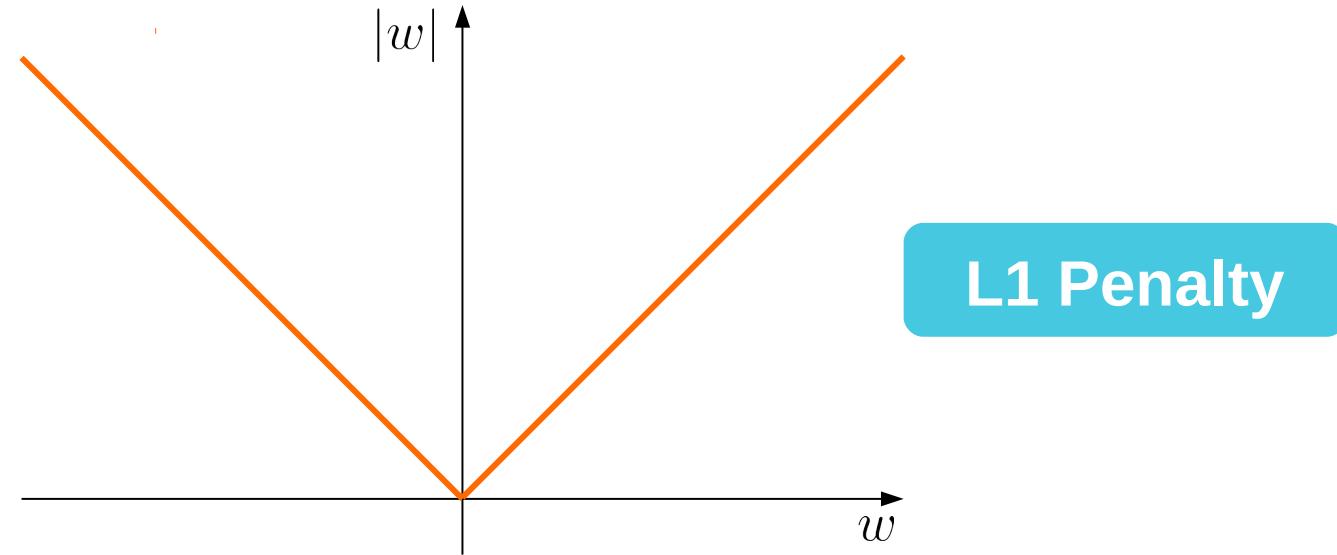
ℓ_1 "lasso" complexity: $\sum_{i=0}^d |w_i|$, for coefficients w_0, \dots, w_d

ℓ_2 "ridge" complexity: $\sum_{i=0}^d w_i^2$, for coefficients w_0, \dots, w_d

Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Optimization techniques
SGD learning rate
Other optimization methods
Regularization
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

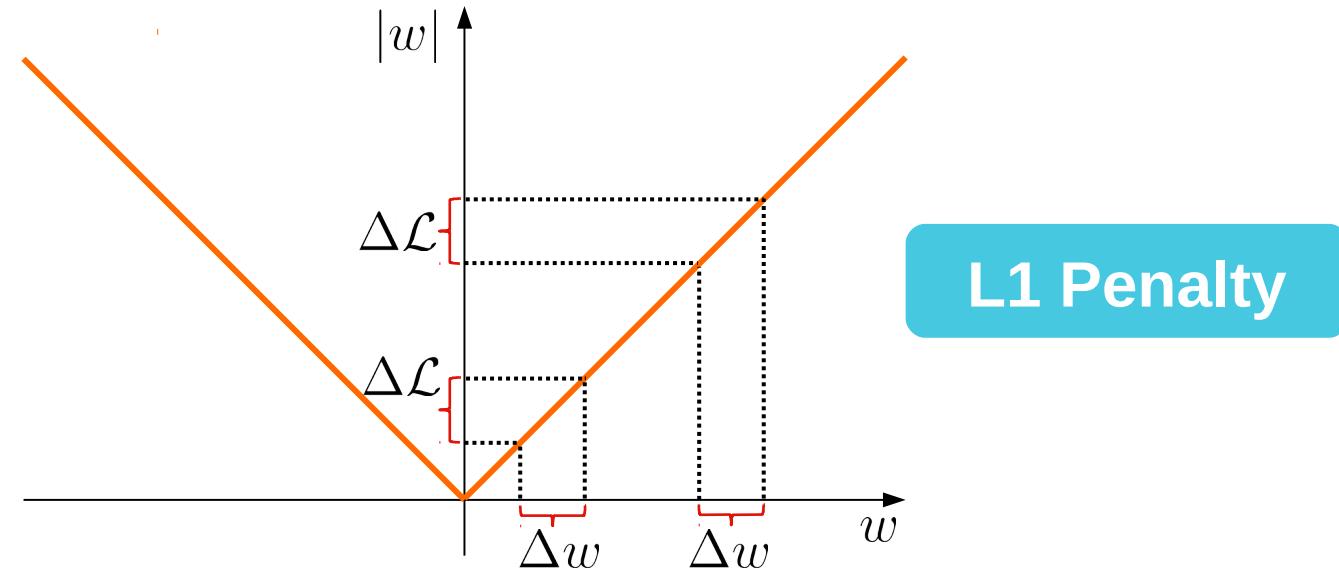
L1 / L2 Regularization



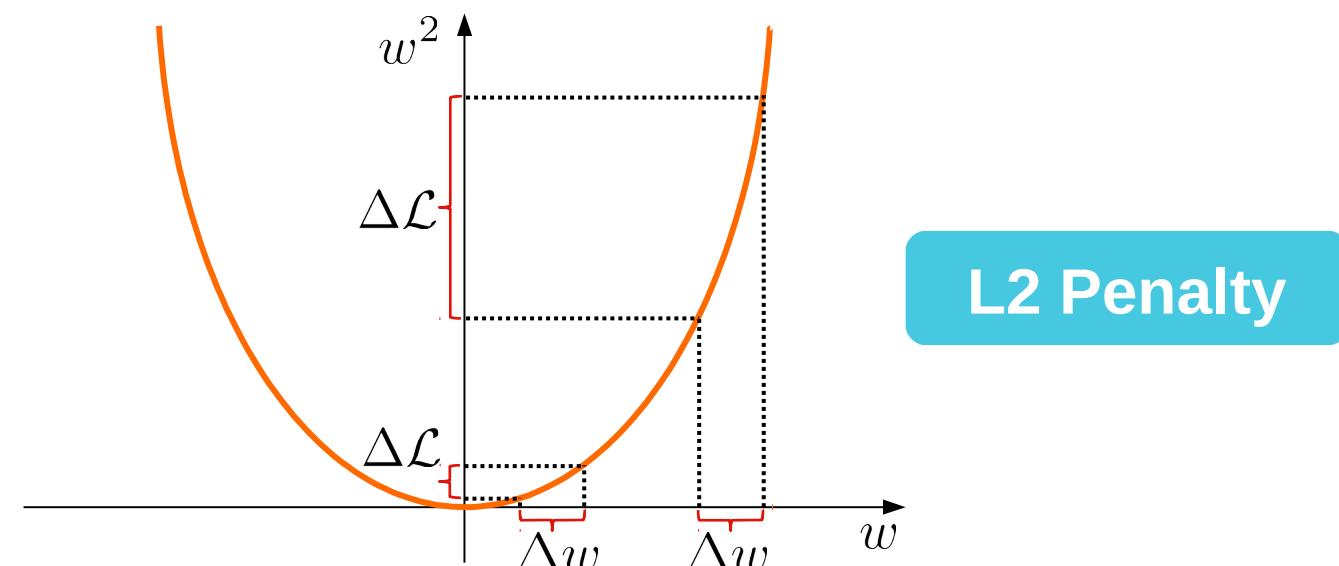
Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Optimization techniques
SGD learning rate
Other optimization methods
Regularization
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

L1 / L2 Regularization



L1 Penalty



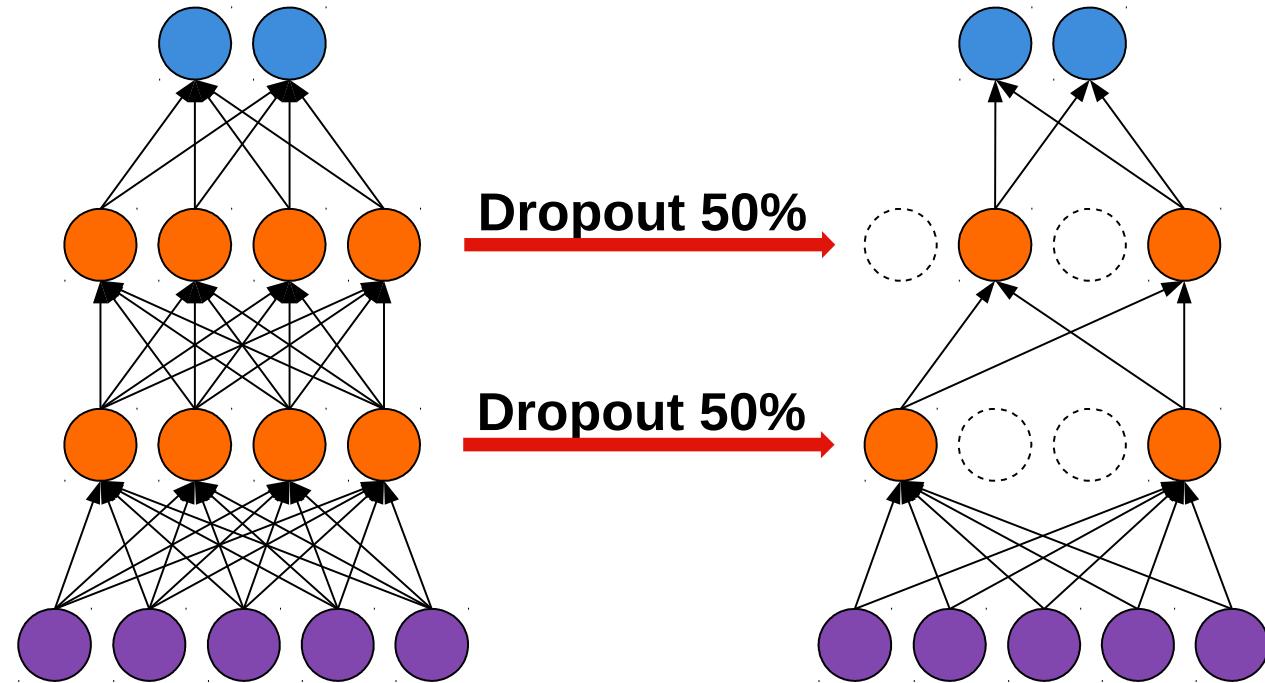
L2 Penalty

Feedforward Neural Networks

Dropout

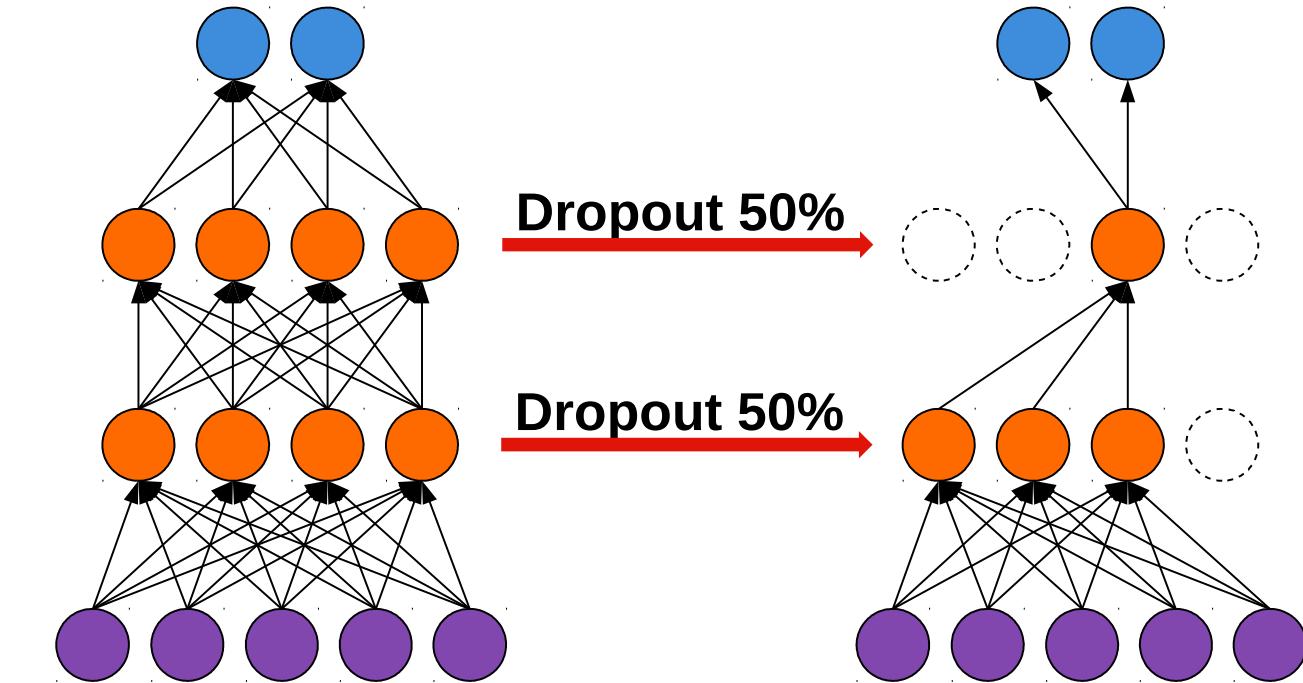
SGD, Epochs, Batches and Steps
Activation Functions
Cross-Entropy Loss
Other optimizers methods
Regularization
Normalizing Inputs
Vanishing/Exploding Gradients
Weights Initialization

Step n



Feedforward Neural Networks

Dropout

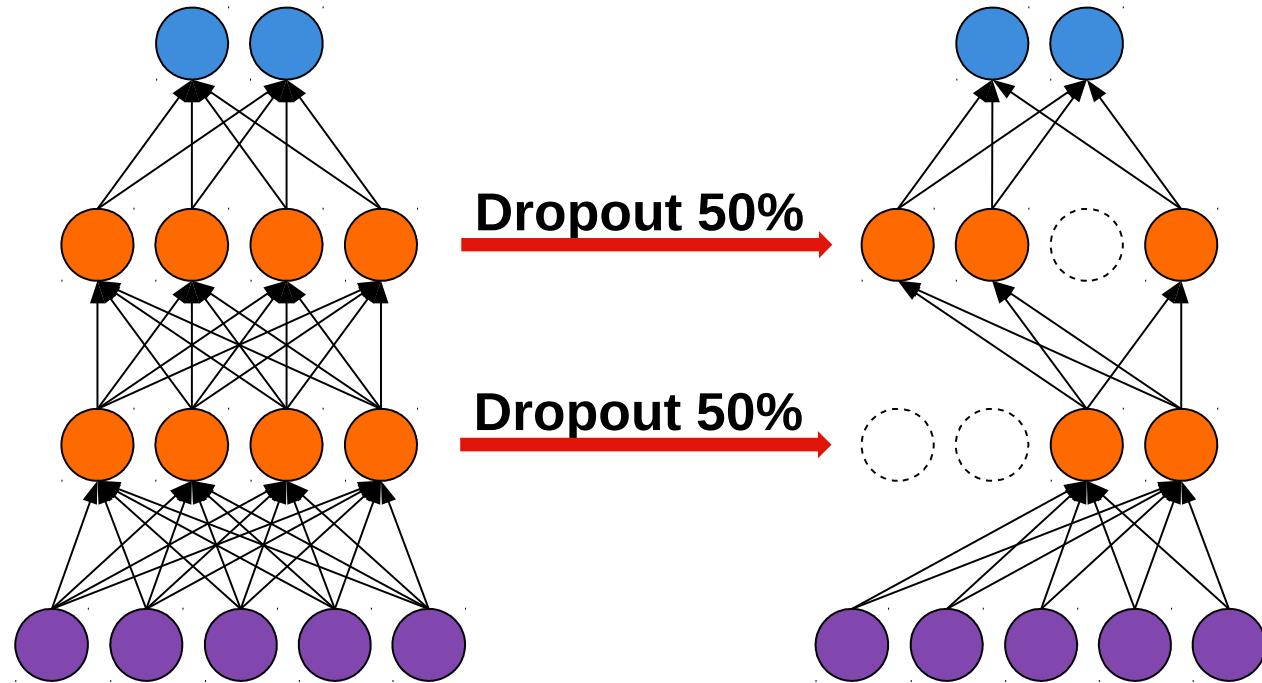


Feedforward Neural Networks

Dropout

SGD, Epochs, Batches and Steps
Activation Functions
Cross-Entropy Loss
SGD Learning Rate
Other Optimization Methods
Regularization
Normalizing Inputs
Vanishing/Exploding Gradients
Weights Initialization

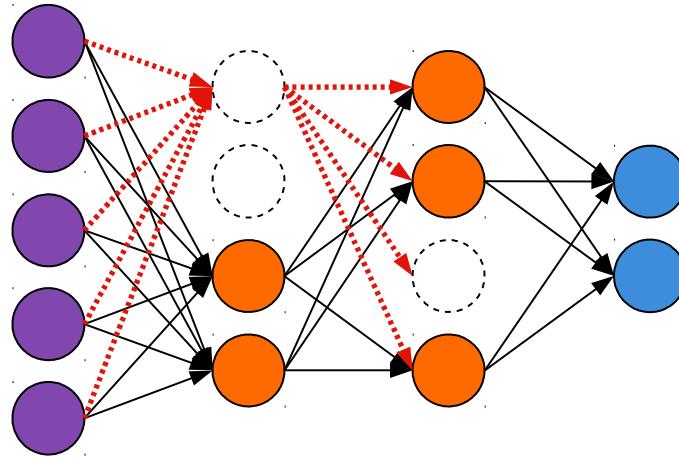
Step n+2



Feedforward Neural Networks

Dropout

SGD, Epochs, Batches and Steps
Activation Functions
ReLU Activation
SGD Learning Rate
Other Optimization Methods
Regularization
Normalizing Inputs
Vanishing/Exploding Gradients
Weights Initialization



Before drop-out:

$$a_0^{[0]} = g \left(w_{00}^{[0]}x_0 + w_{10}^{[0]}x_1 + w_{20}^{[0]}x_2 + w_{30}^{[0]}x_3 + w_{40}^{[0]}x_4 + b_0^{[0]} \right)$$

After drop-out: $a_0^{[0]} = 0$

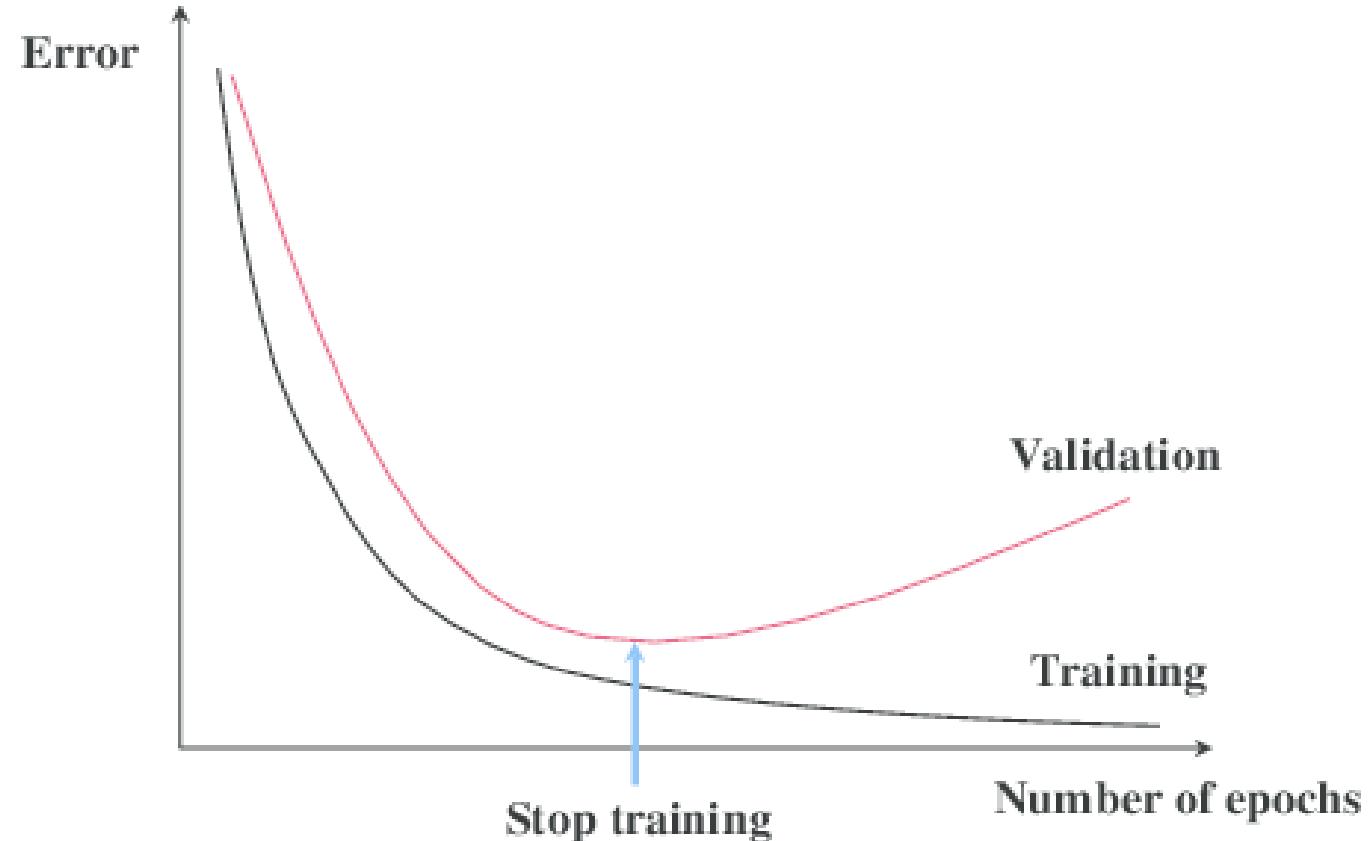
What **complexity** does this method use?

ℓ_0 complexity: Number of non-zero coefficients

Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Adam optimization
SGD learning rate
Other optimization methods
Regularization
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

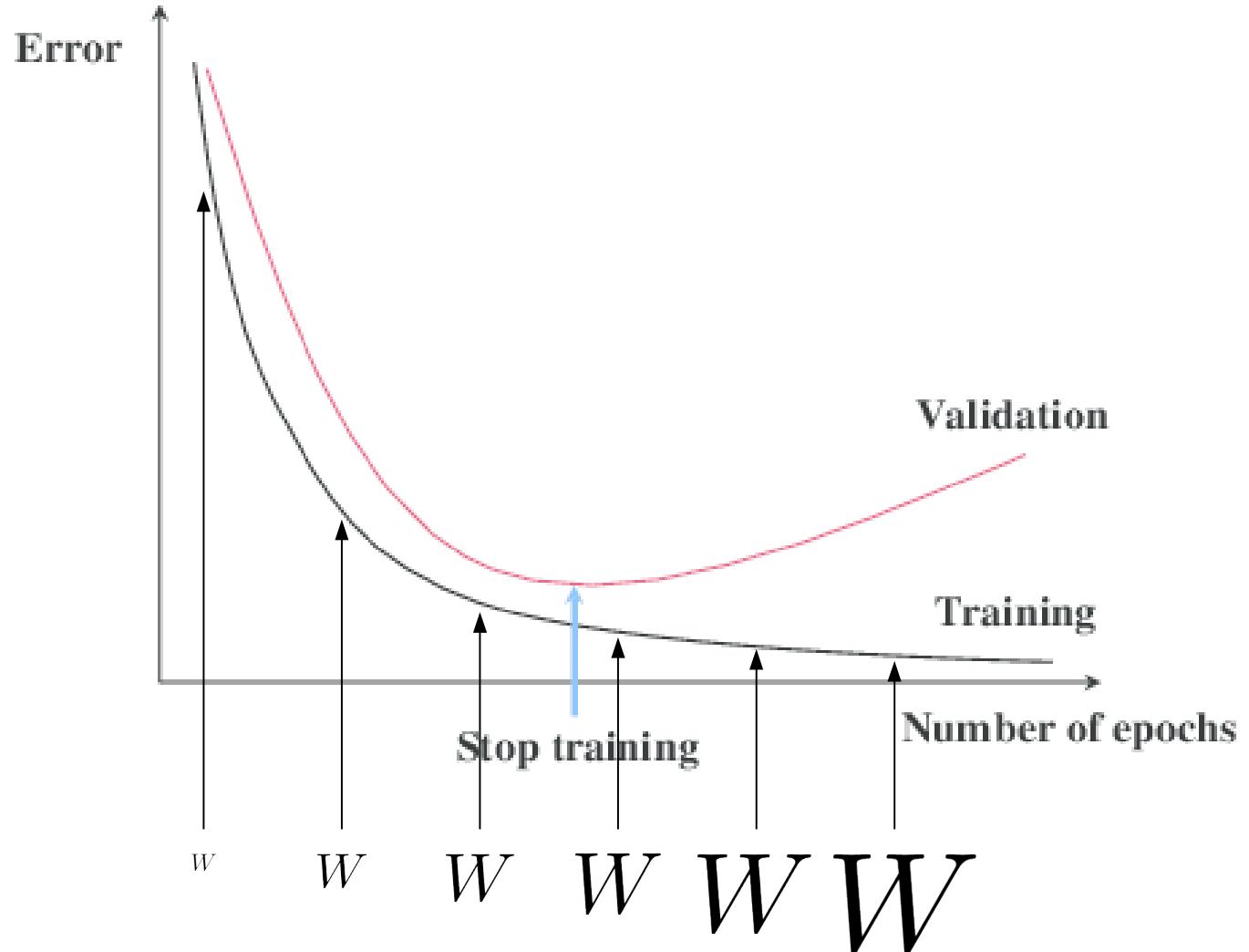
Early Stopping



Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Adam optimization
SGD learning rate
Other optimization methods
Regularization
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

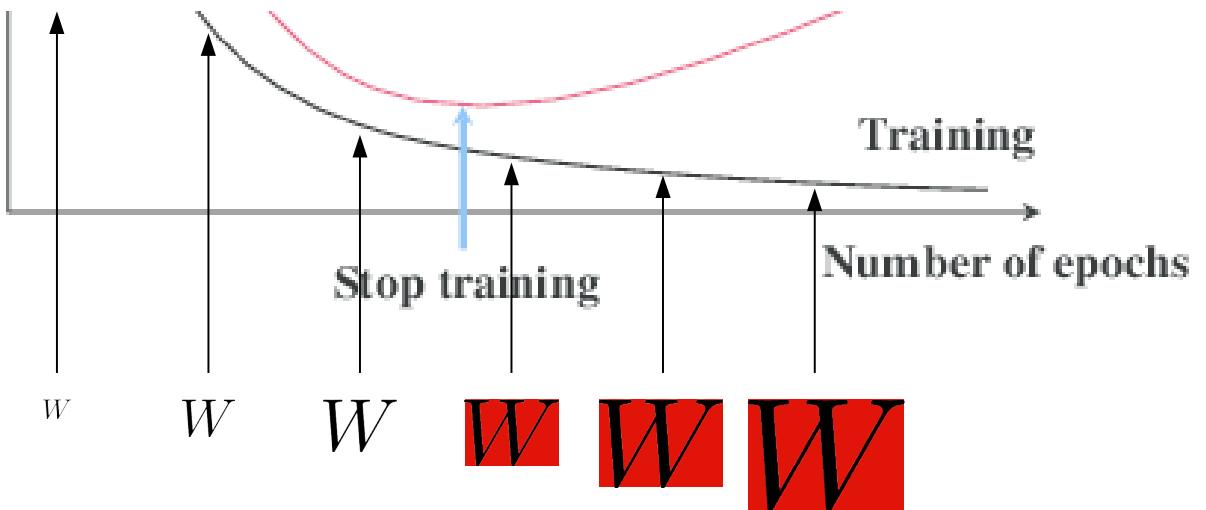
Early Stopping



Feedforward Neural Networks

SGD, Epochs, Batches and Steps
Adam optimization
SGD learning rate
Other optimization methods
Regularization
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

Early Stopping



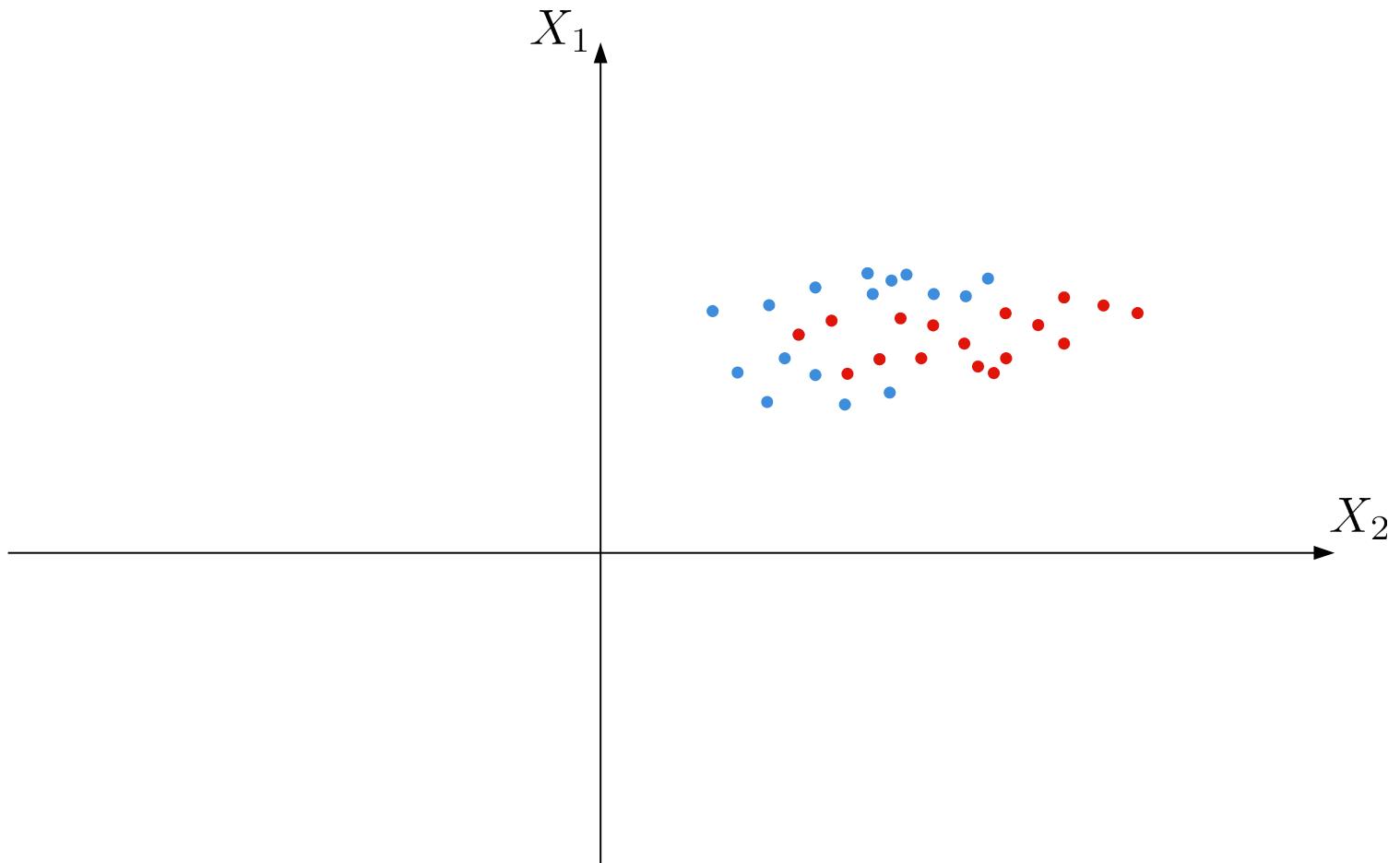
What **complexity** does this method use?

ℓ_1 "lasso" complexity: $\sum_{i=0}^d |w_i|$, for coefficients w_0, \dots, w_d

ℓ_2 "ridge" complexity: $\sum_{i=0}^d w_i^2$, for coefficients w_0, \dots, w_d

Feedforward Neural Networks

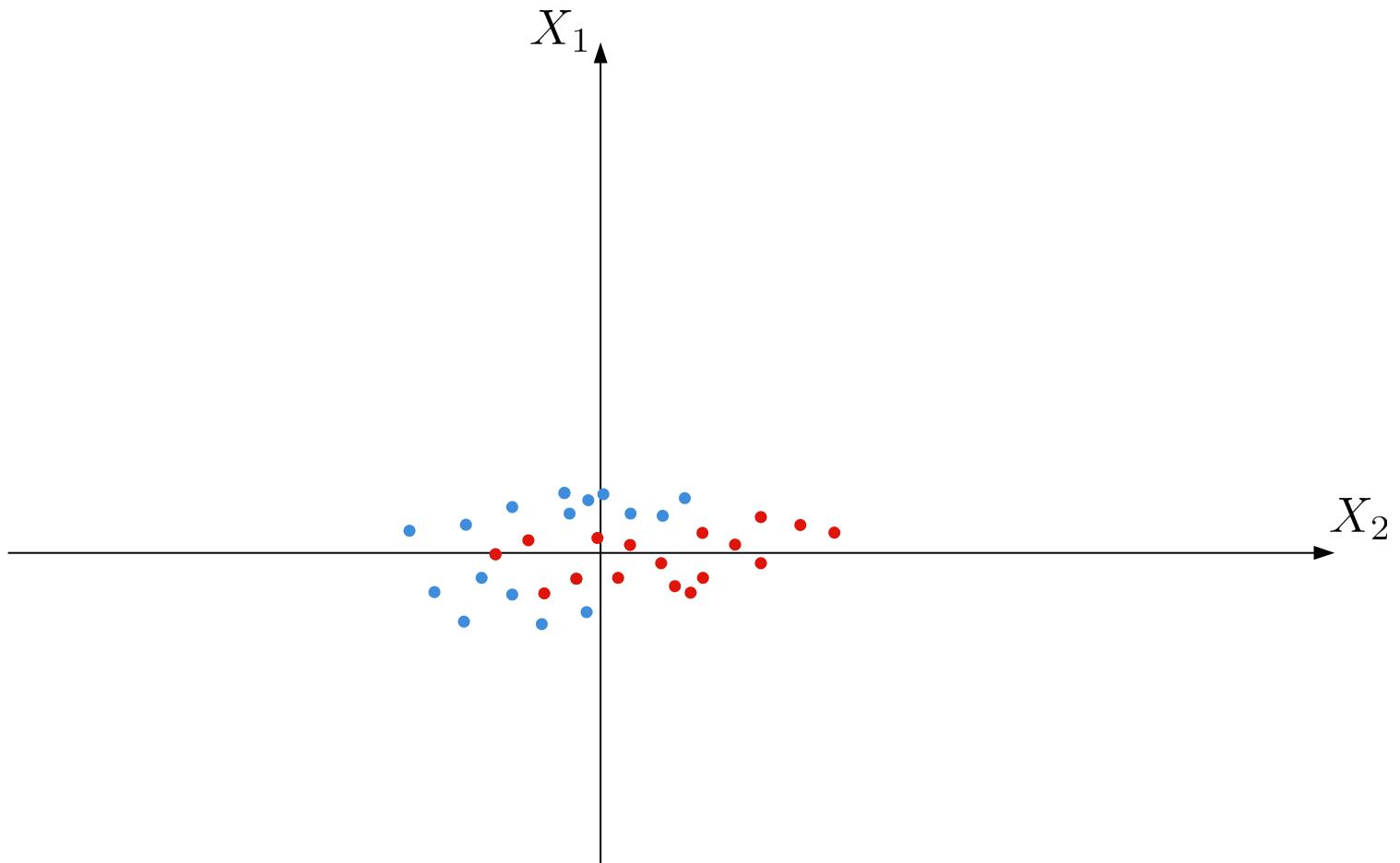
**SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization**



$$x = \frac{x - \mu}{\sigma^2}$$

Feedforward Neural Networks

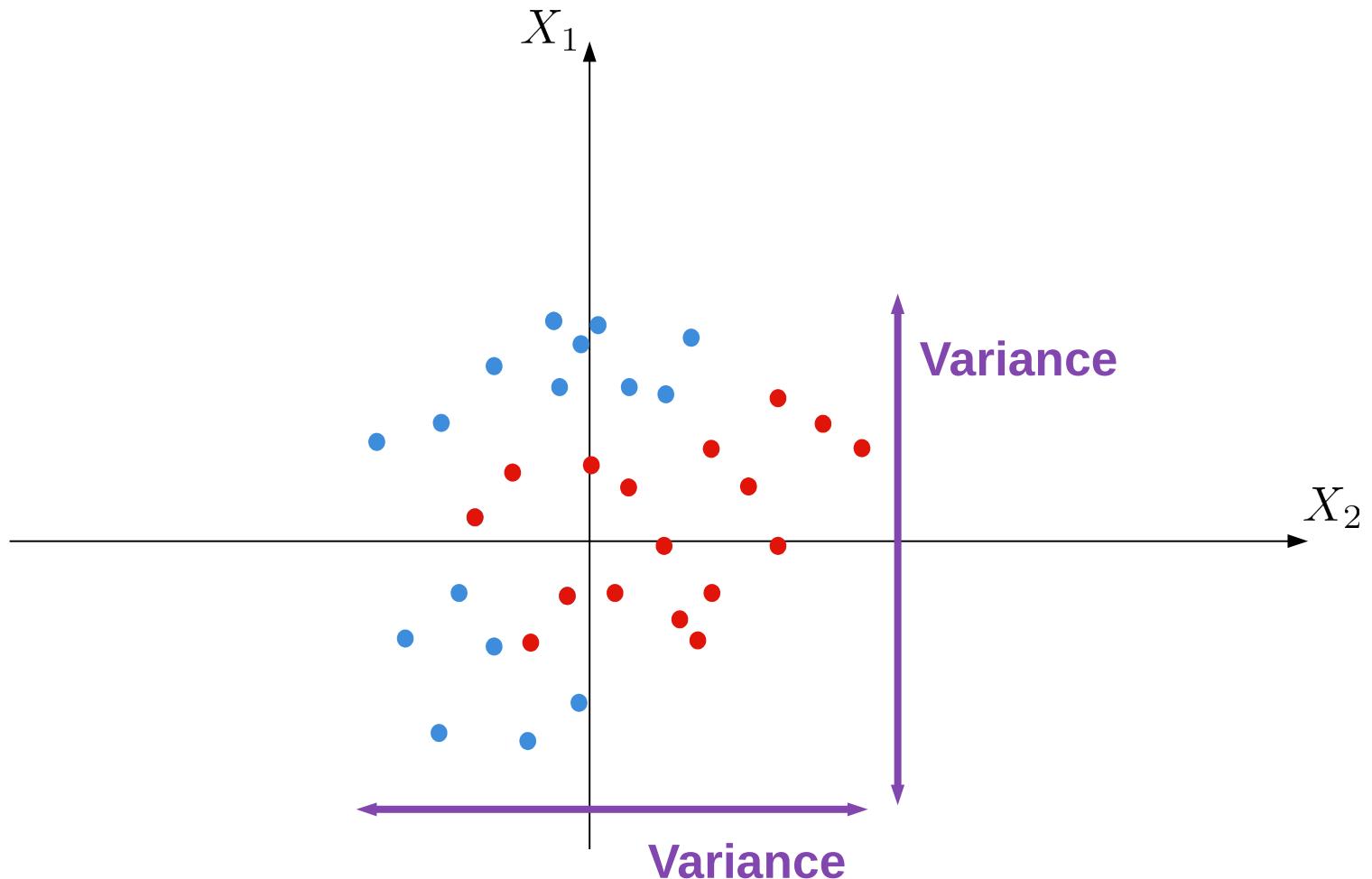
**SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization**



$$x = \frac{x - \mu}{\sigma^2}$$

Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization



$$x = \frac{x - \mu}{\sigma^2}$$

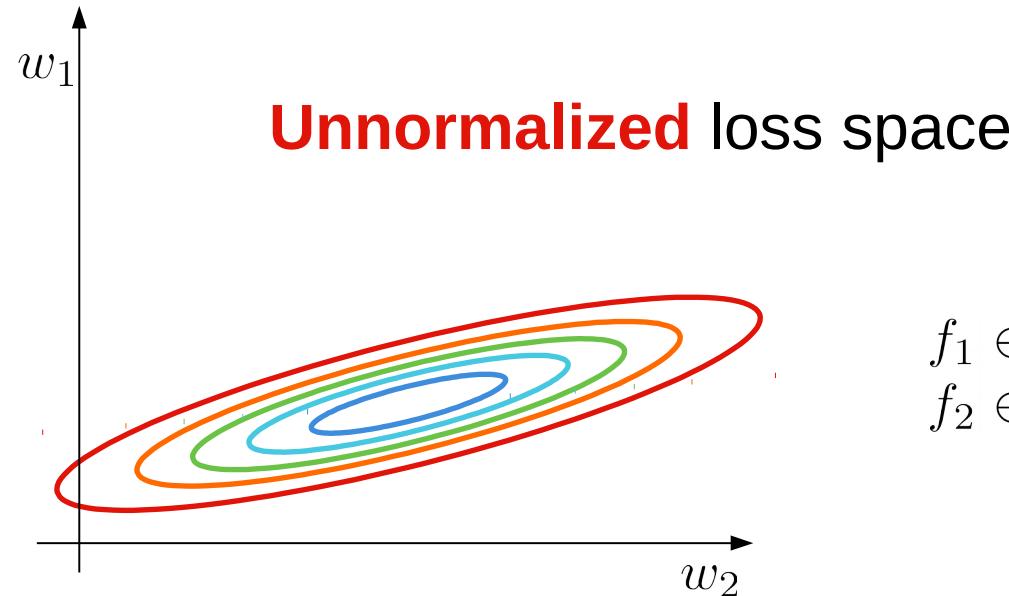
Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization methods
Regularization
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

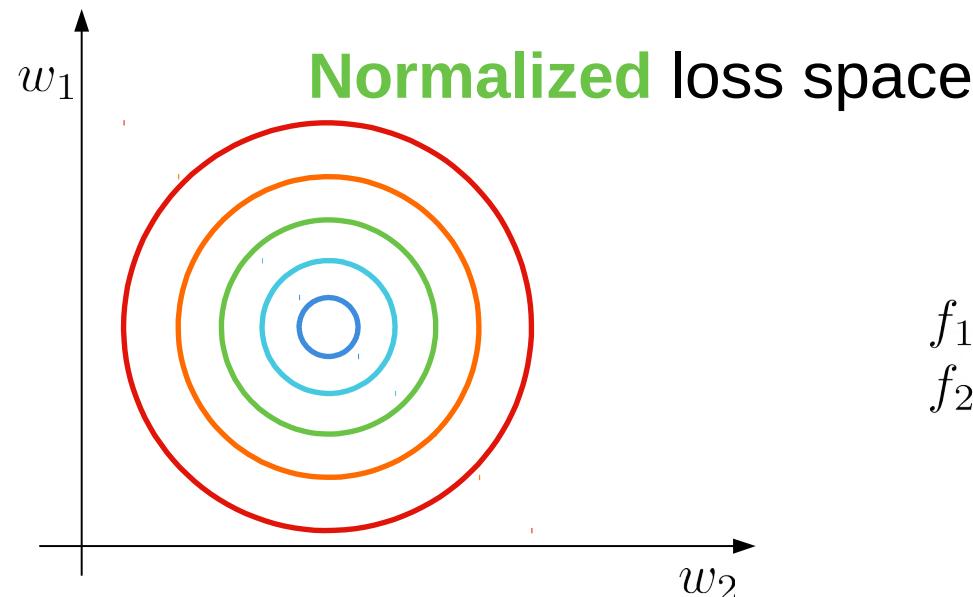


HIGH PERFORMANCE ARTIFICIAL INTELLIGENCE

Why **input normalization** matters?



$$f_1 \in [1, 1000]$$
$$f_2 \in [0, 1]$$



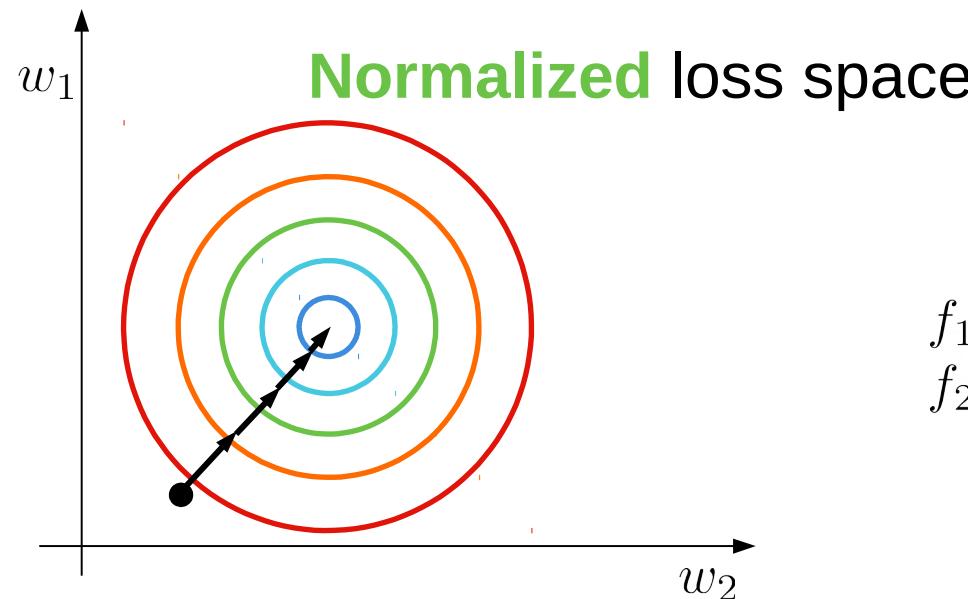
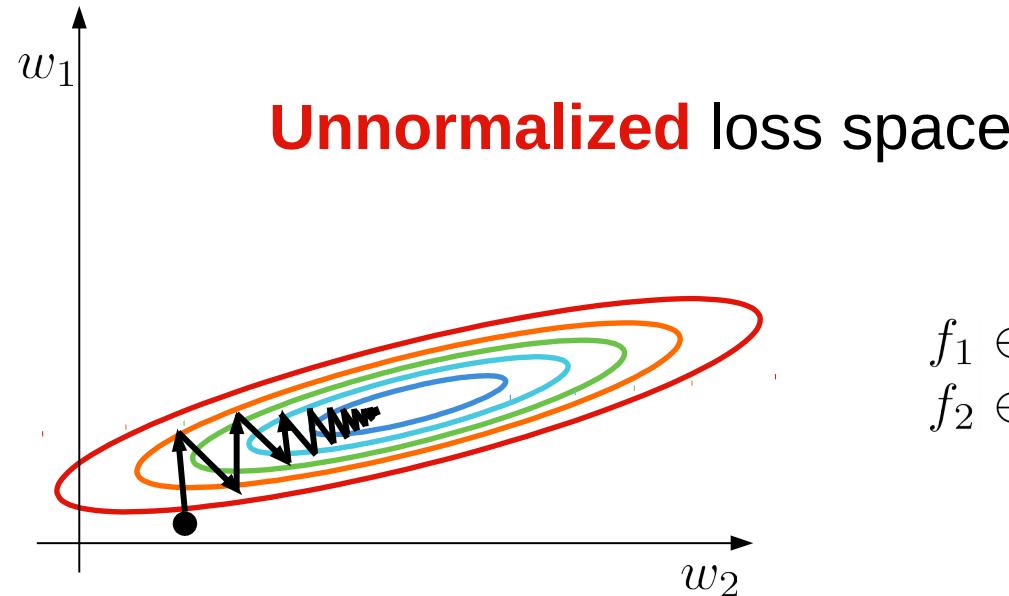
$$f_1 \in [-0.5, 0.5]$$
$$f_2 \in [-0.5, 0.5]$$

Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization methods
Regularization
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization

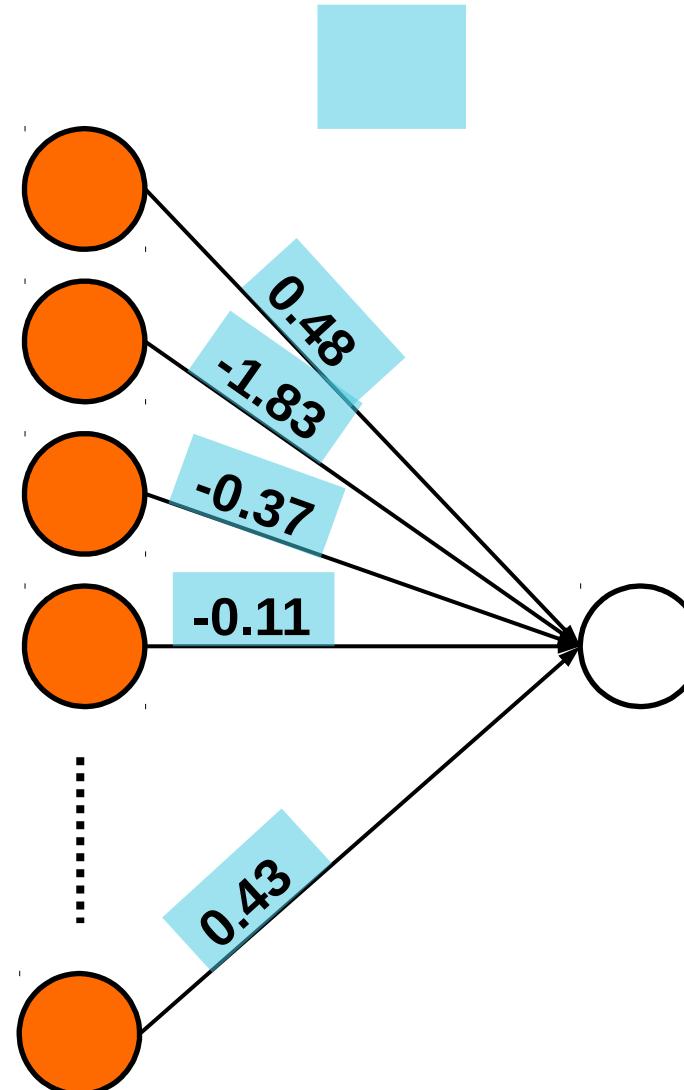


Why **input normalization** matters?

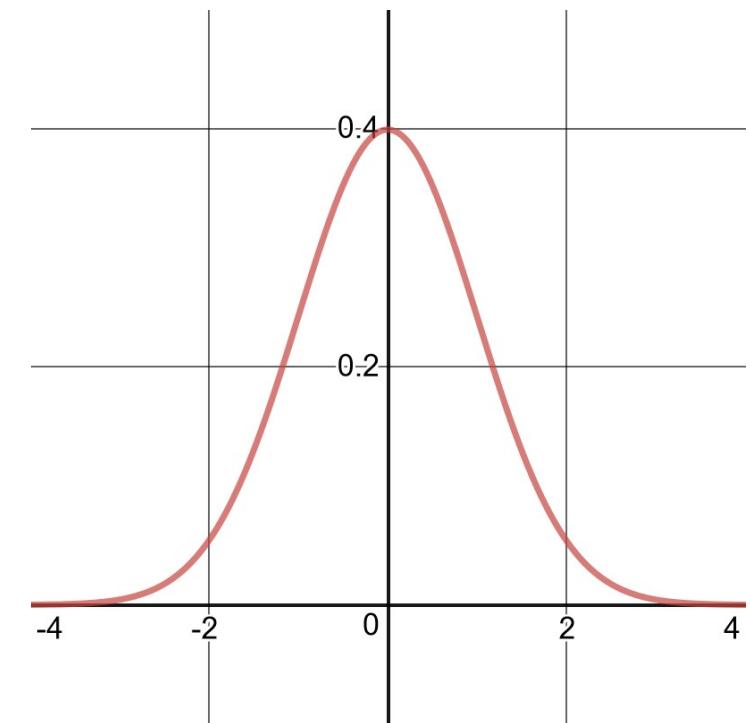


Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization

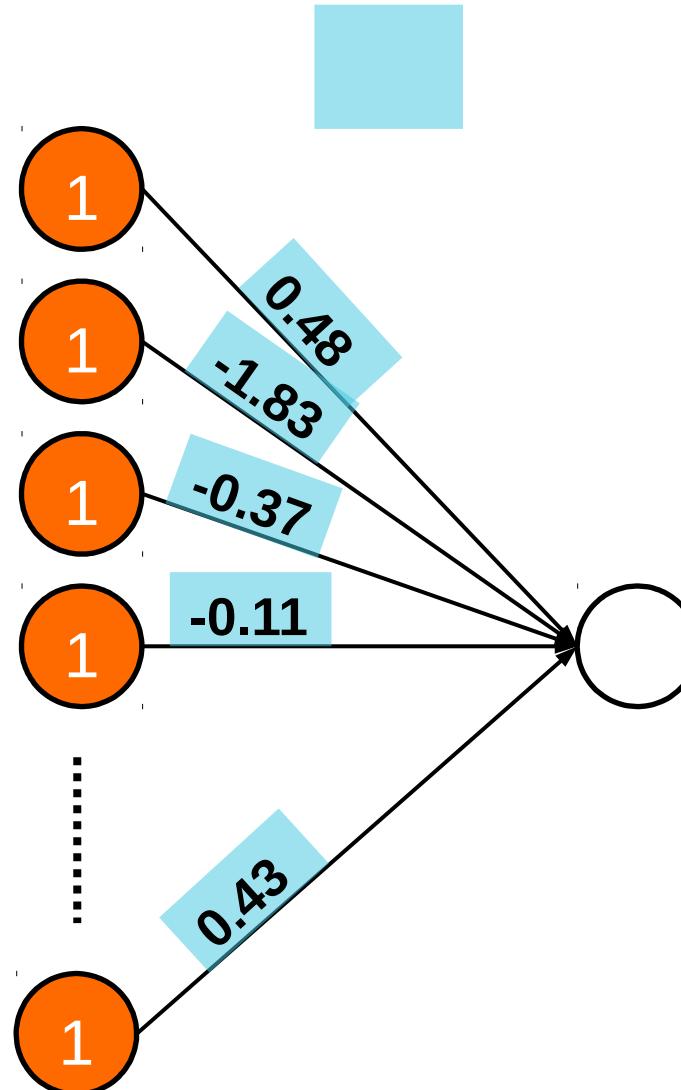


Normal Distribution

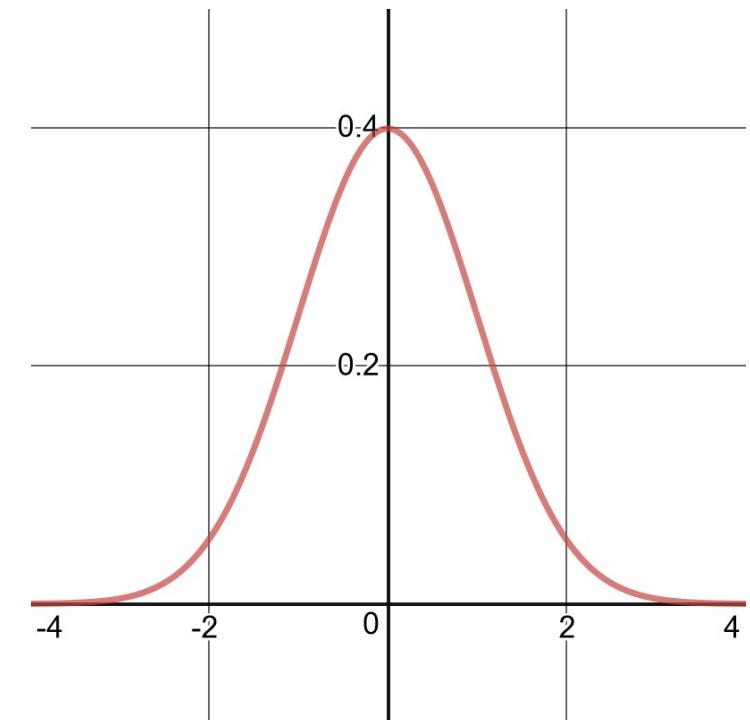


Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization

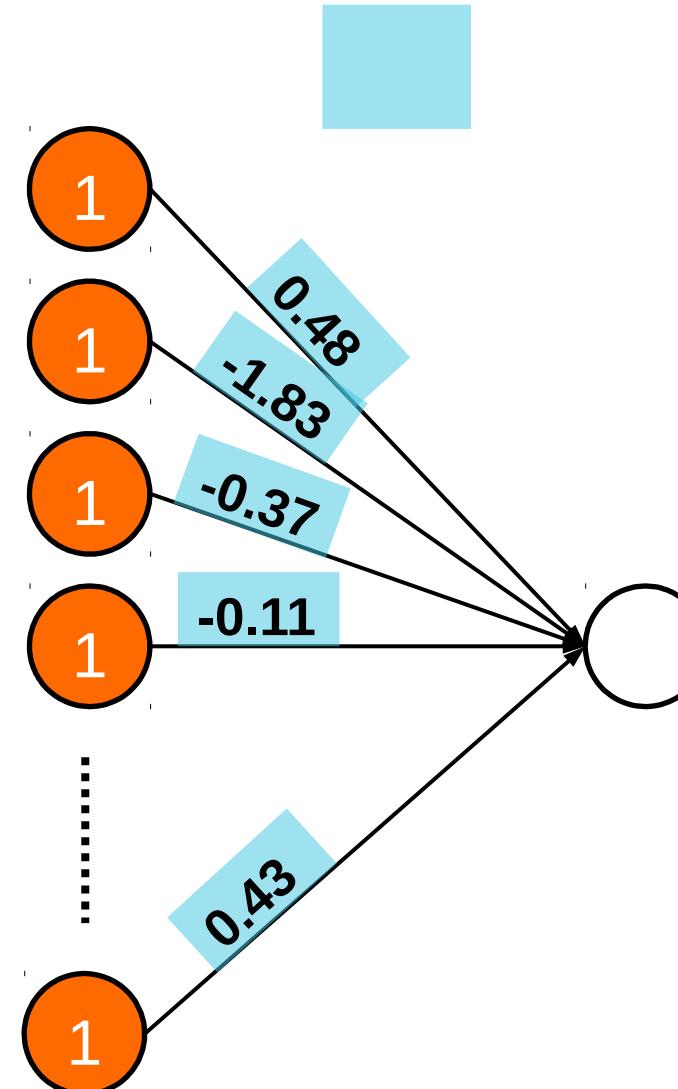


Normal Distribution



Feedforward Neural Networks

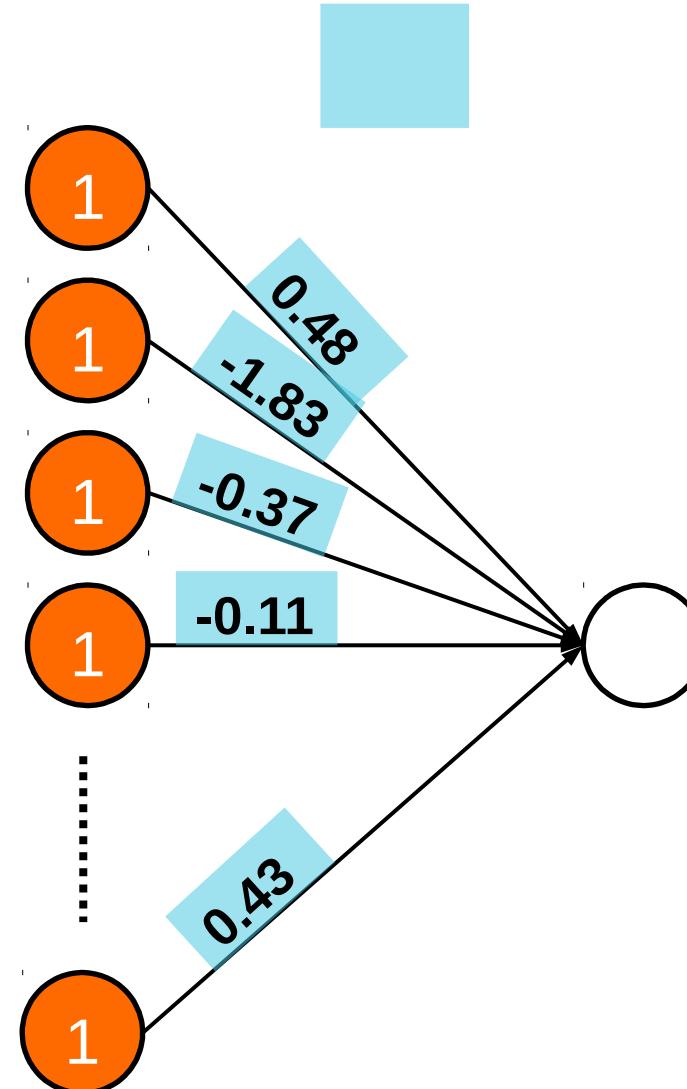
SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization



Sum of **independent random variables** that are **normally distributed**:

Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization

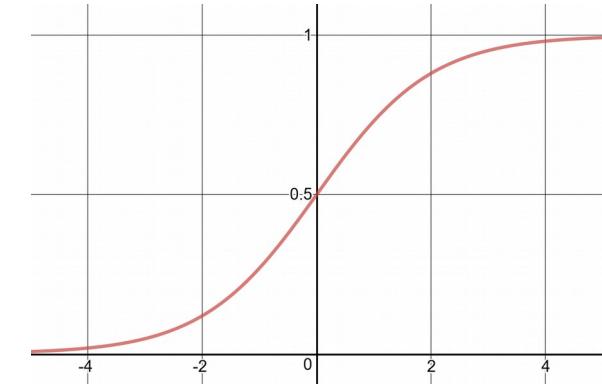
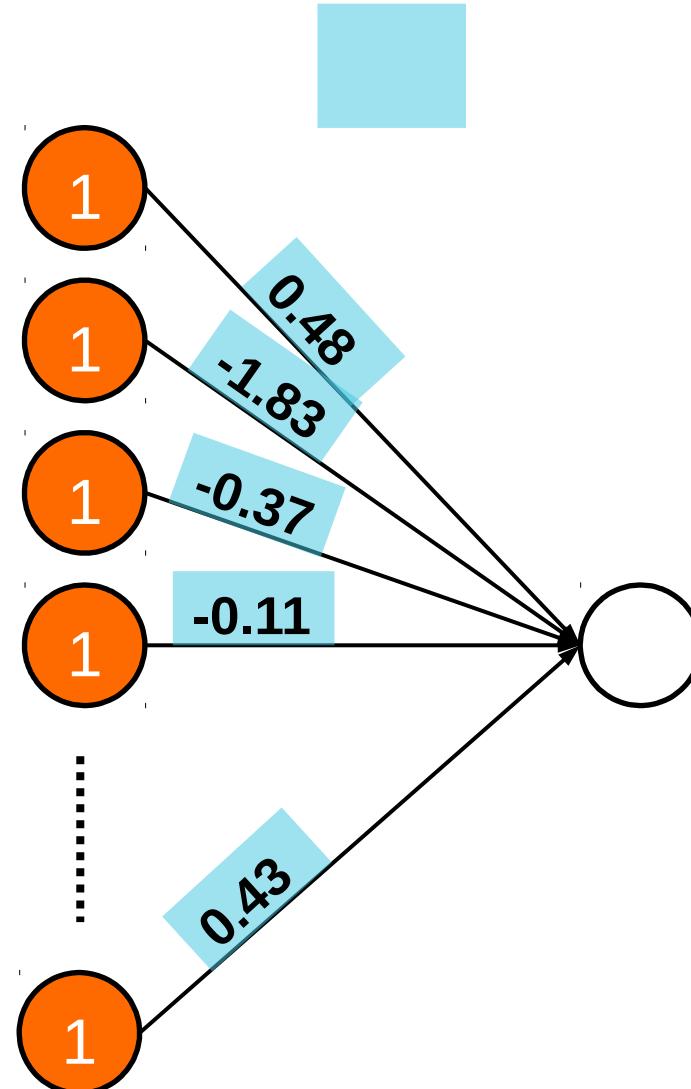


Sum of **independent random variables** that are **normally distributed**:



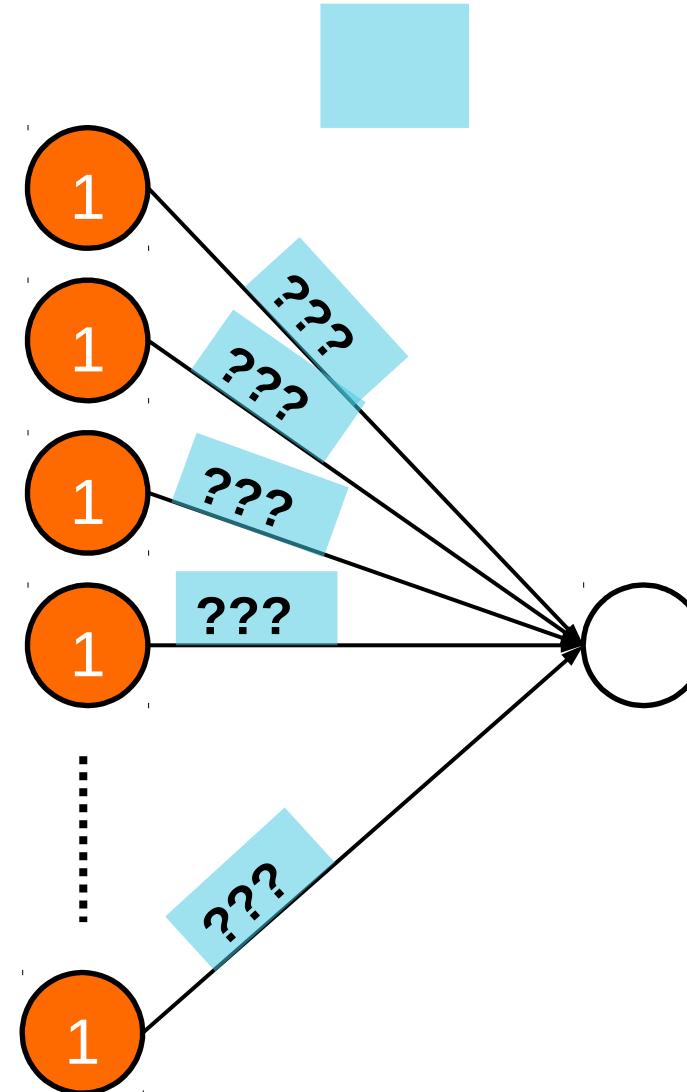
Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization



Feedforward Neural Networks

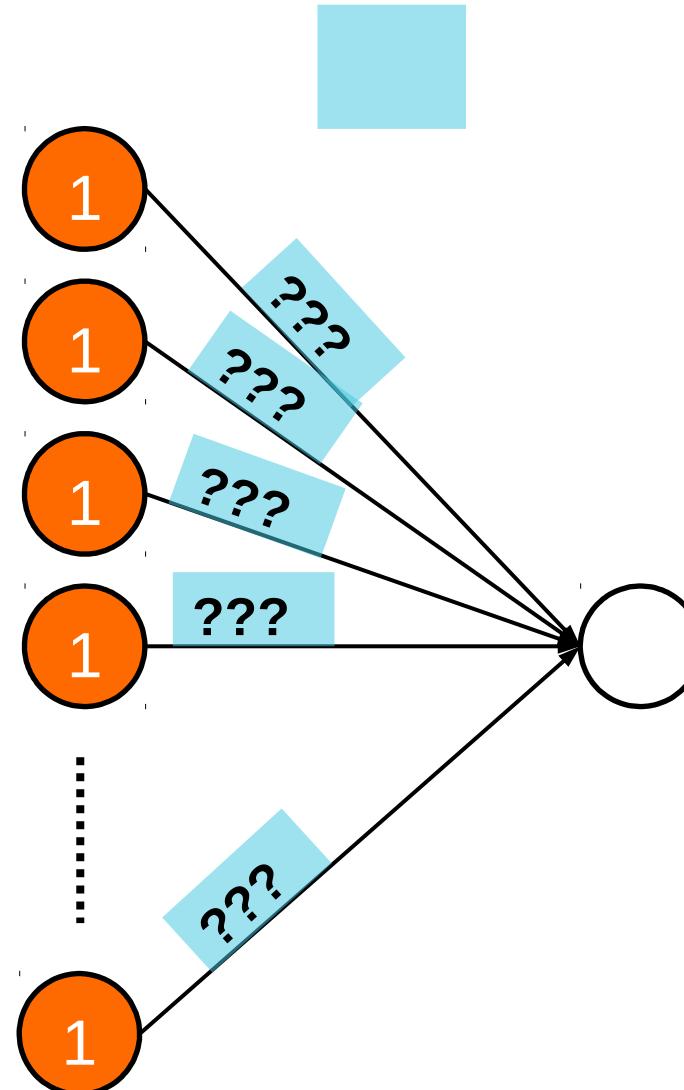
SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization



Alternative weights
initializations?

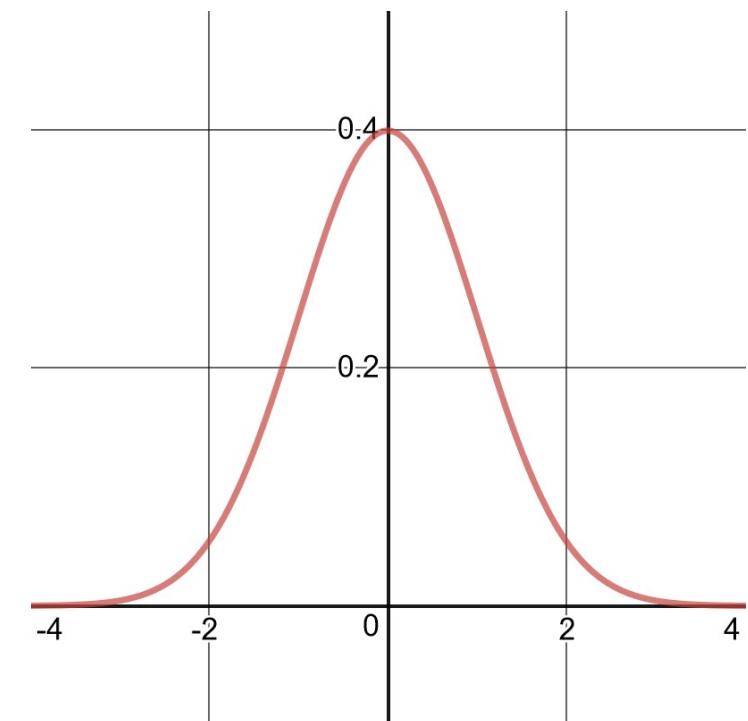
Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization



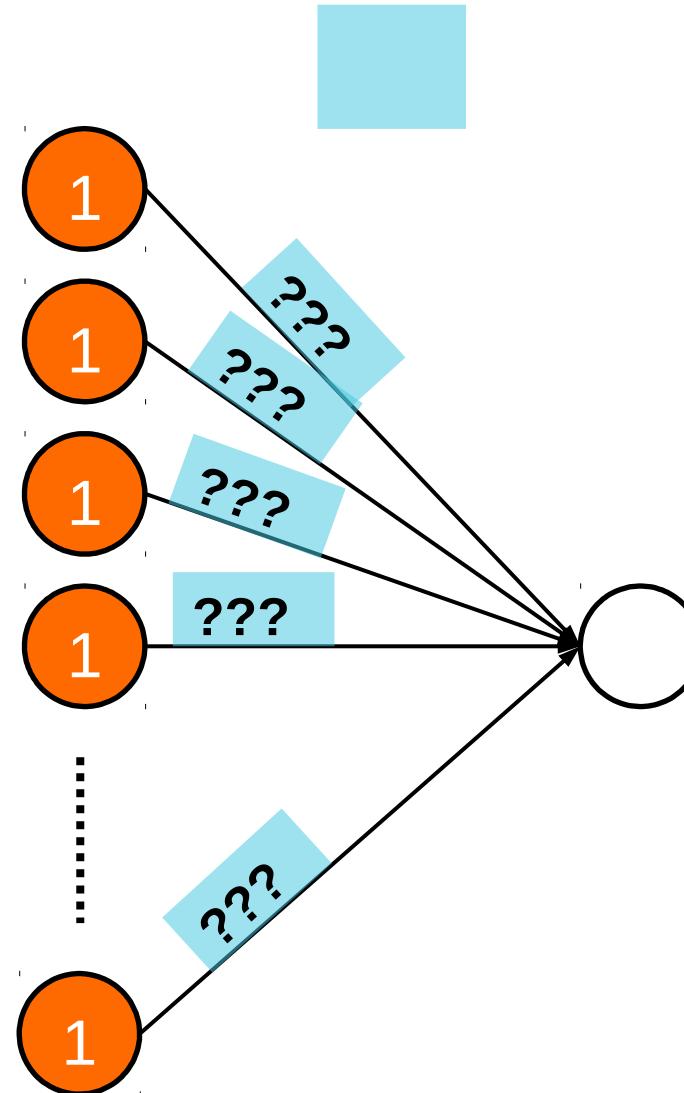
Truncated Normal Initialization

Normal Distribution



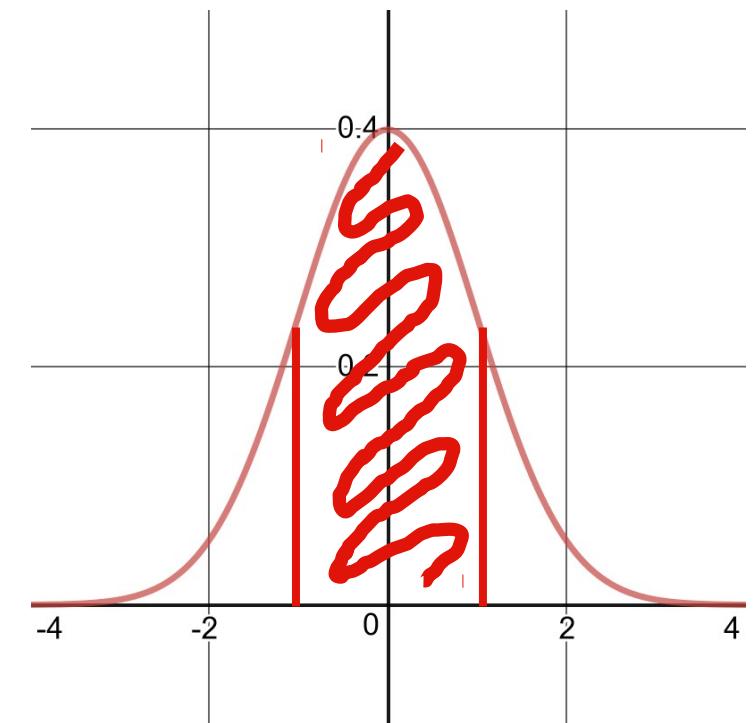
Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization



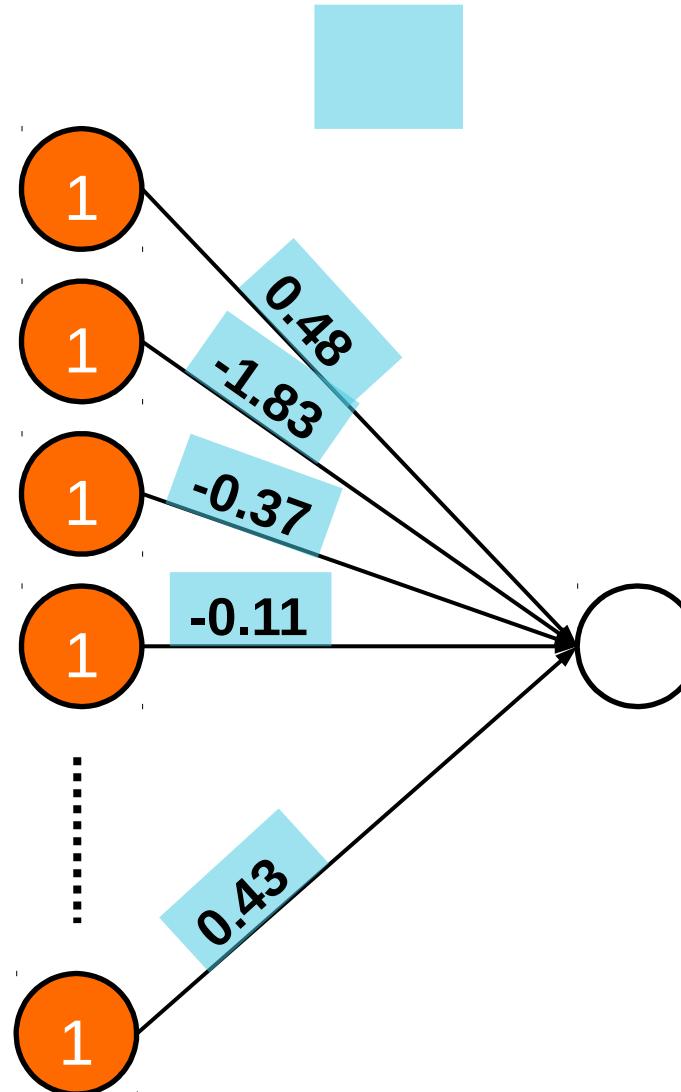
Truncated Normal Initialization

Normal Distribution



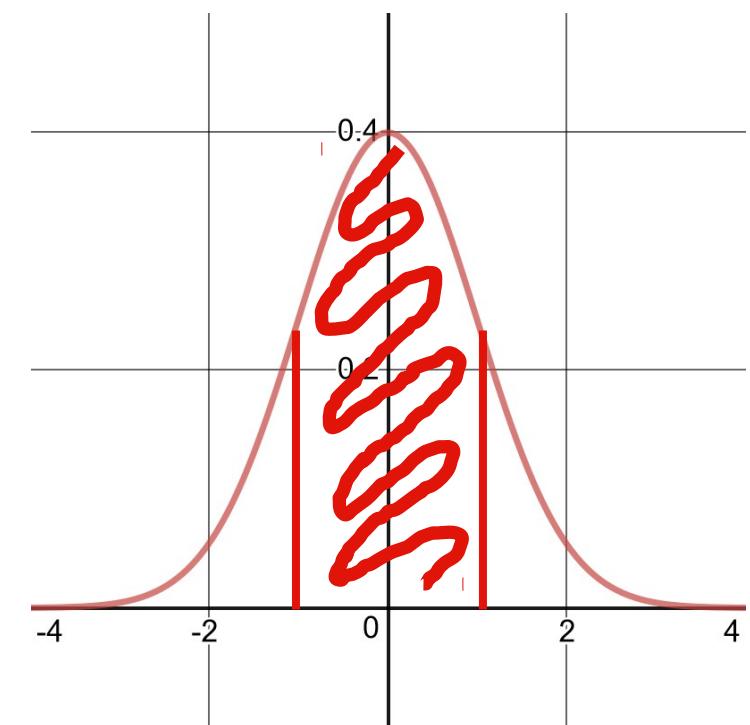
Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization



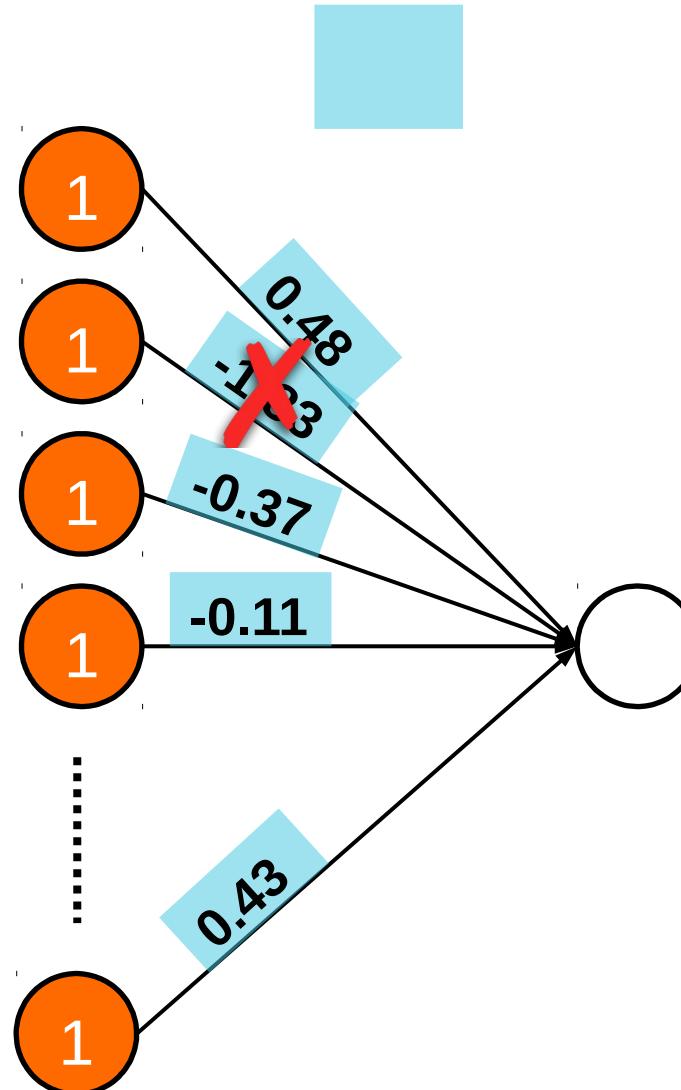
Truncated Normal Initialization

Normal Distribution



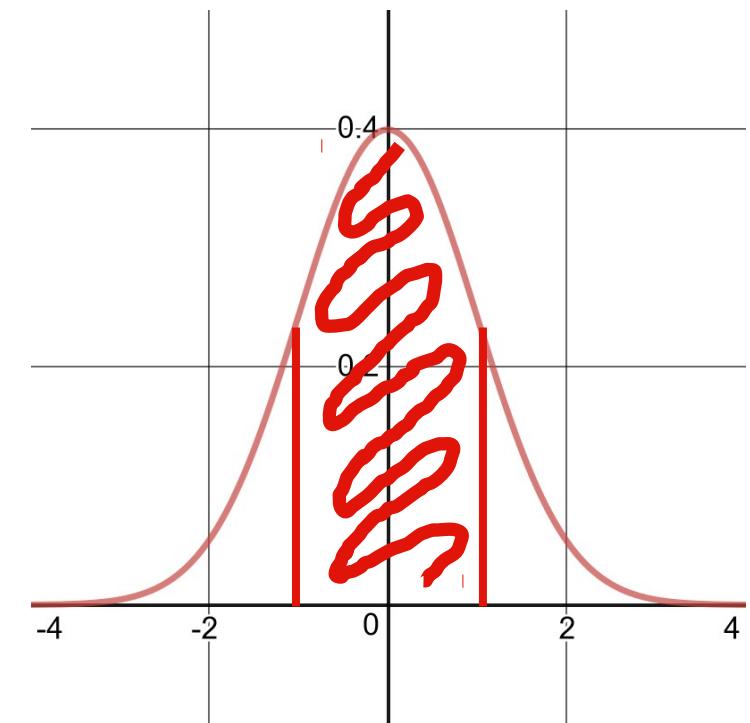
Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization methods
Regularization
Normalizing inputs
Vanishing/Exploding Gradients
Weights initialization



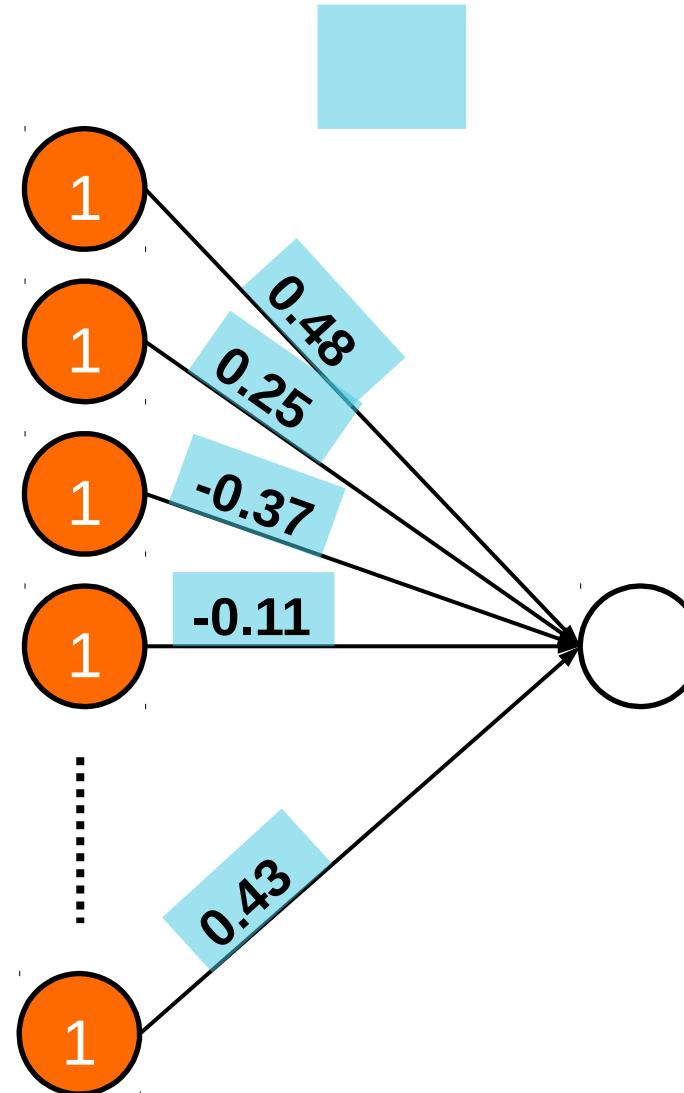
Truncated Normal Initialization

Normal Distribution



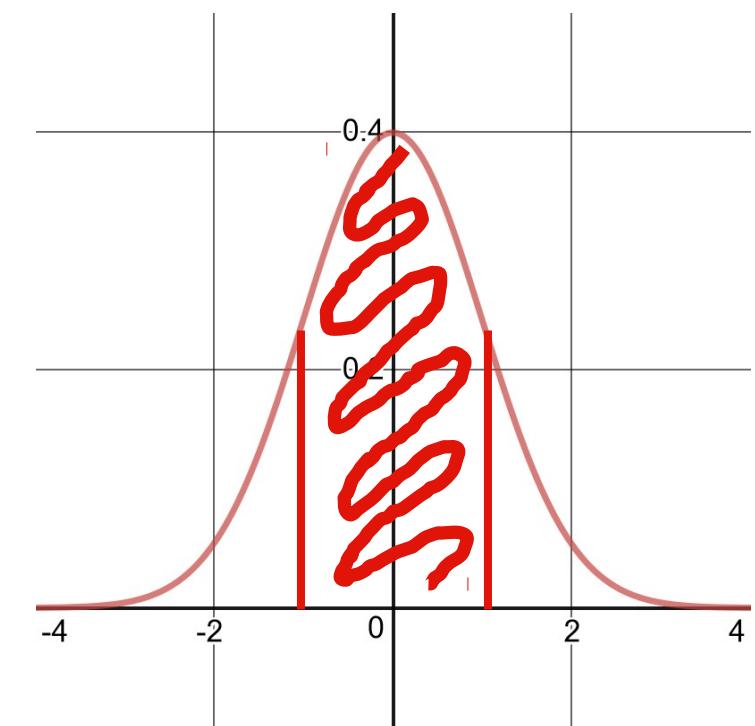
Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization



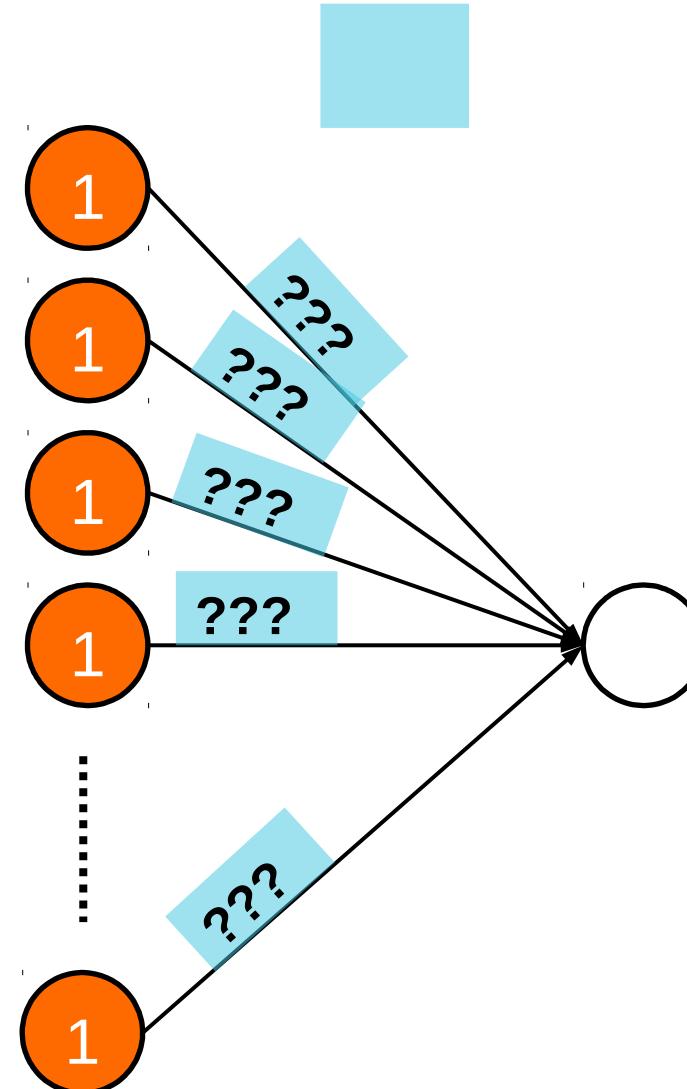
Truncated Normal Initialization

Normal Distribution



Feedforward Neural Networks

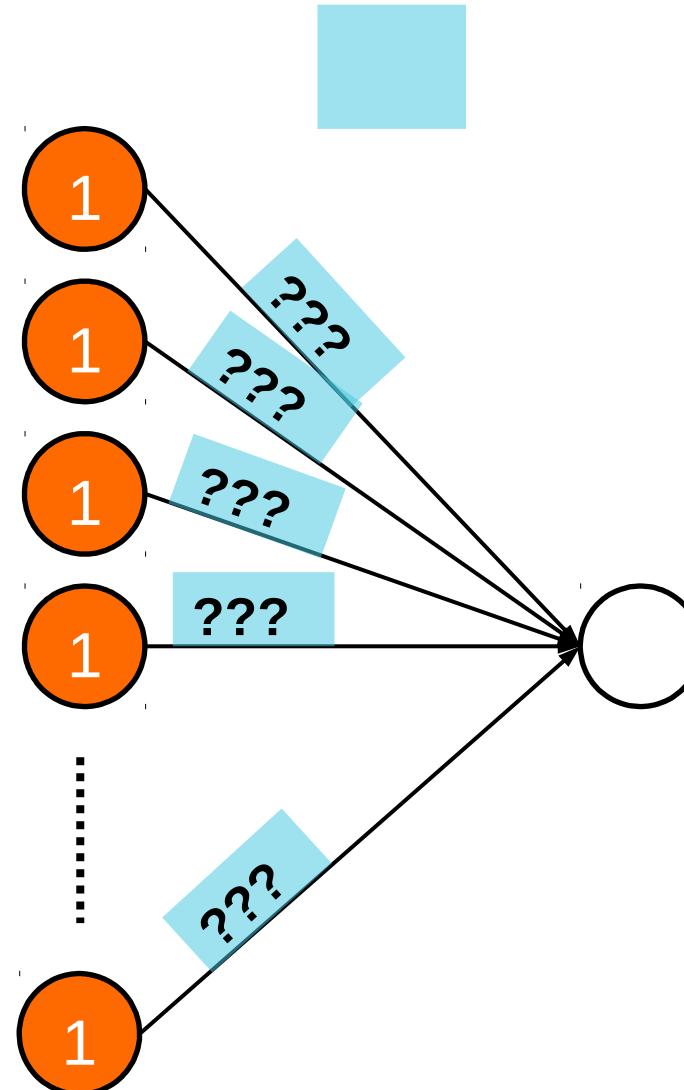
SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization



Xavier / Glorot
Initialization

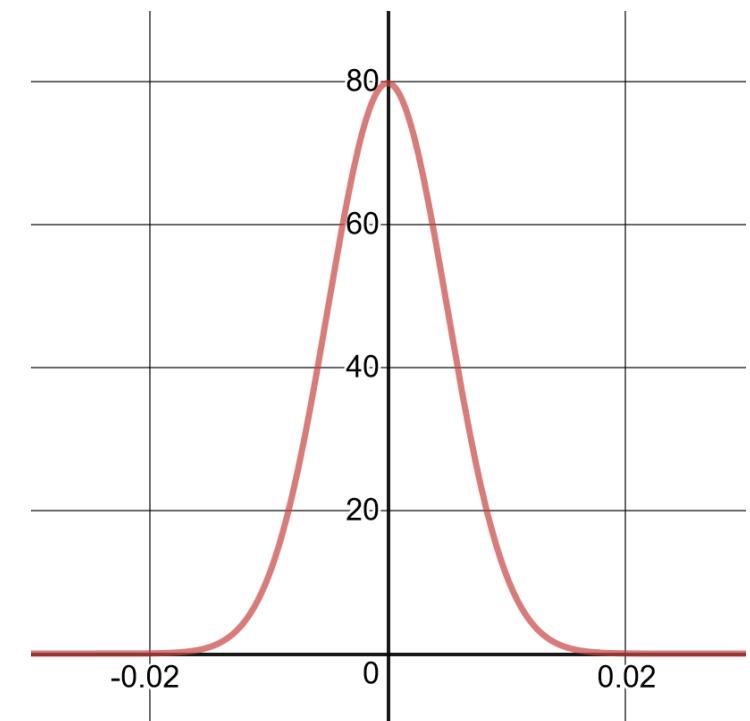
Feedforward Neural Networks

SGD, Epochs,
Batches and Steps
Activation functions
SGD learning rate
Other optimization
methods
Regularization
Normalizing inputs
Vanishing/Exploding
Gradients
Weights initialization



Xavier / Glorot Initialization

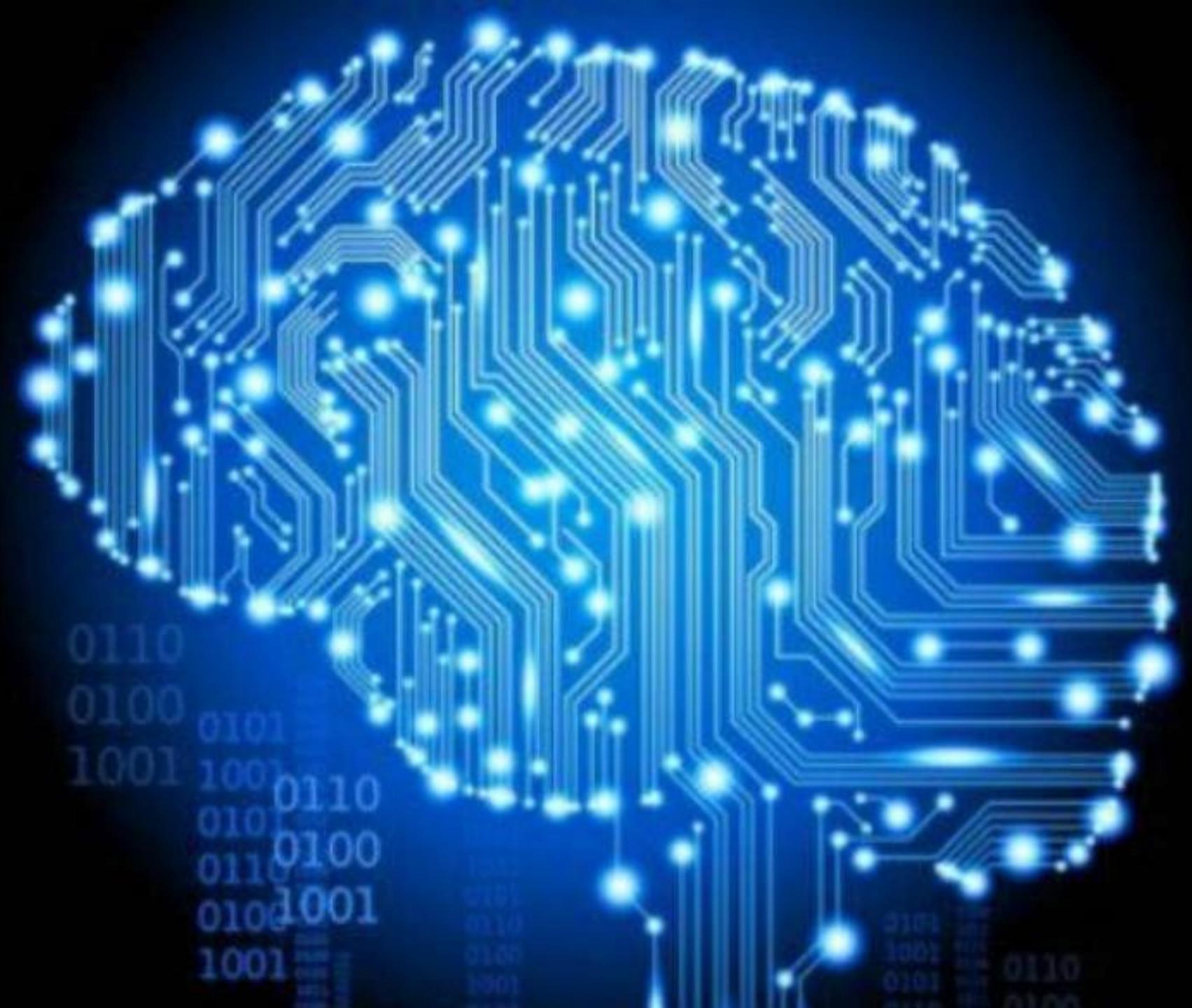
Normal Distribution



Practical Aspects

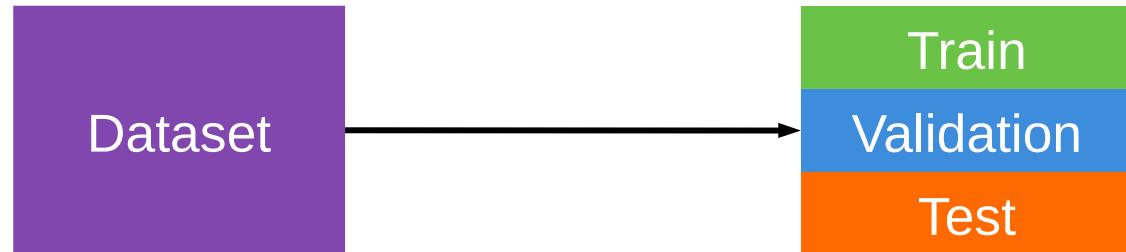


HIGH PERFORMANCE
ARTIFICIAL INTELLIGENCE



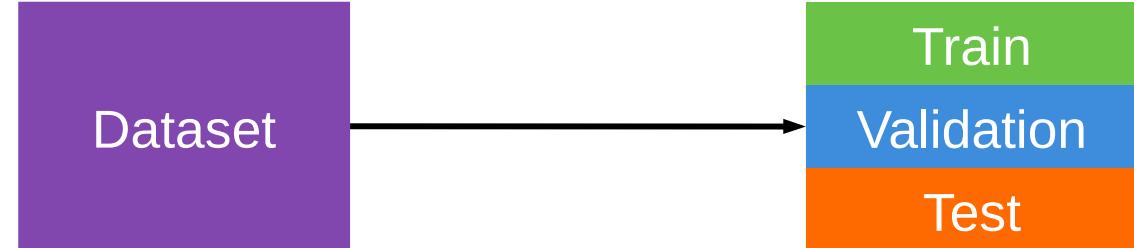
Practical Aspects

Train /
Validation / Test
Workflow



- | | |
|------------|--|
| Train | Set used to train & control bias |
| Validation | Set used to control variance |
| Test | Set used to estimate the generalization error of the final model |

Practical Aspects



**Train /
Validation / Test**

Workflow

Train

Set used to **train & control bias**

Validation

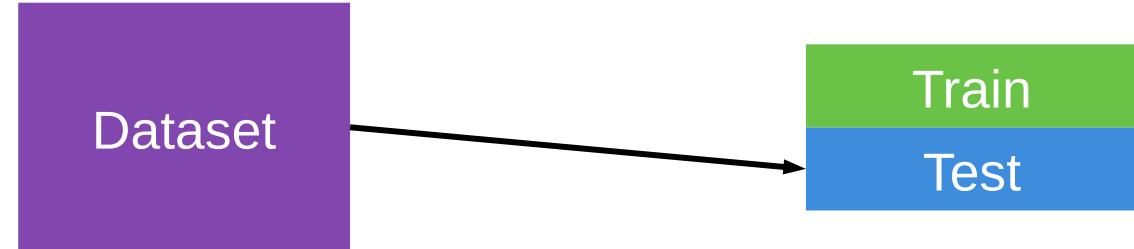
Set used to **control variance**

OPTIONAL

Set used to **estimate** the
generalization error of the final
model

Practical Aspects

Train / Validation / Test Workflow



Train	Set used to train & control bias
Test	Set used to control variance

OPTIONAL

Practical Aspects

Train /
Validation / Test
Workflow



Machine Learning



How many **images** do we **need** for each set?

Train

As many as possible

Validation

The **minimum** amount to appropriately **represent each class**

Test

The **minimum** amount to appropriately **represent each class**

Practical Aspects

Machine Learning



Typical dataset size: 1.000 - 30.000

**Train /
Validation / Test
Workflow**

Deep Learning



Typical dataset size: 30.000 - 10.000.000



Practical Aspects

Train /
Validation / Test
Workflow

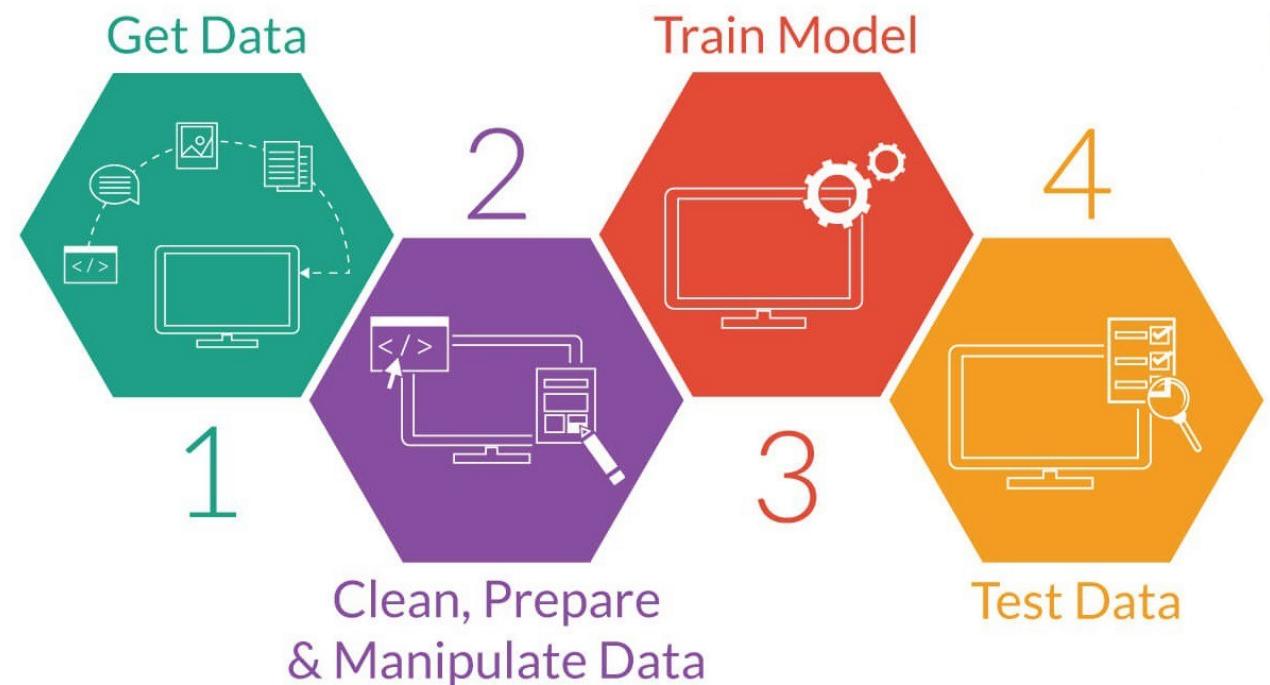


Regardless of what happens, **never never ever**

- Mix correlated data in train/val/test
 - Not straightforward!
- Use different sources for train/val/test
 - Shuffle is your friend
- Perform unrepeatable results
 - Seed is a good friend of Shuffle, which makes it your friend too
- Trust test results only
- Trust validation results only if you have done thorough hyperparameter tuning

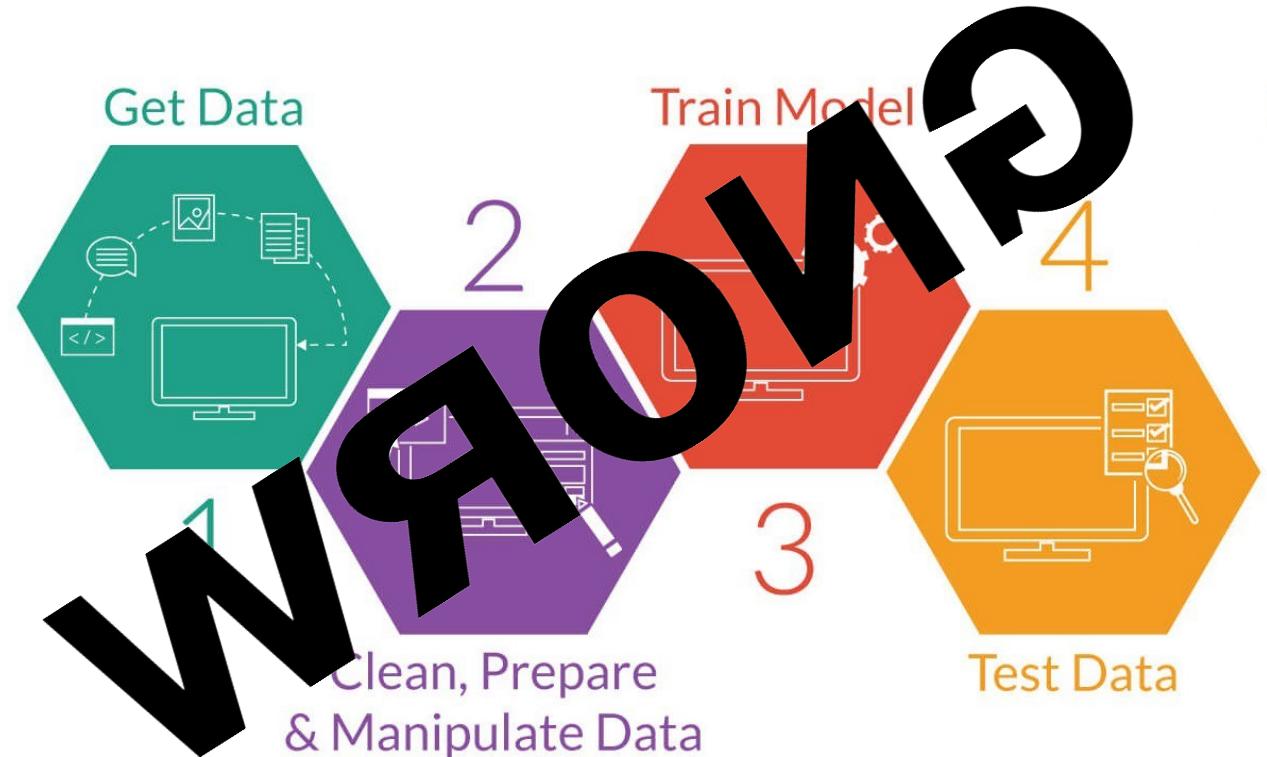
Practical Aspects

**Train /
Validation / Test
Workflow**



Practical Aspects

**Train /
Validation / Test
Workflow**



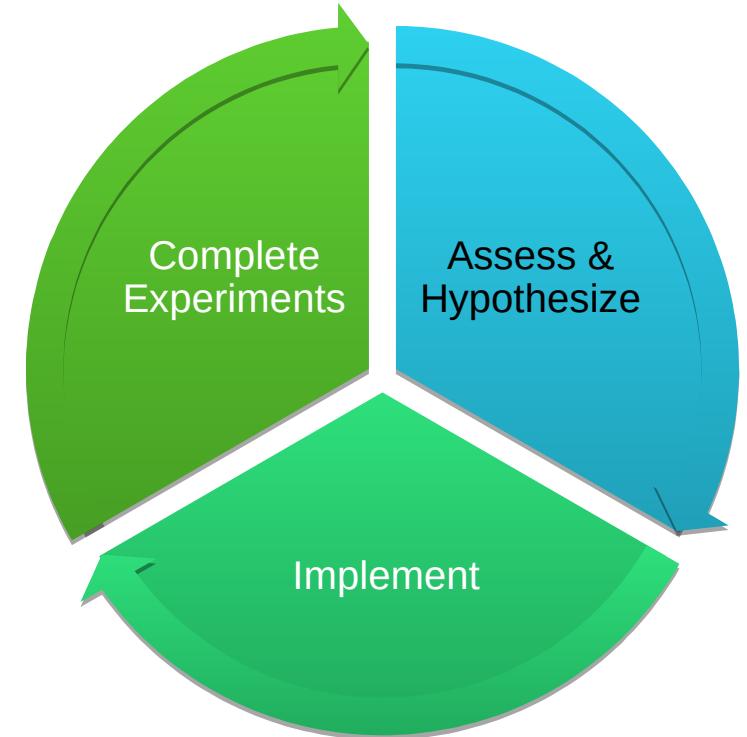
Practical Aspects

Train /
Validation / Test
Workflow



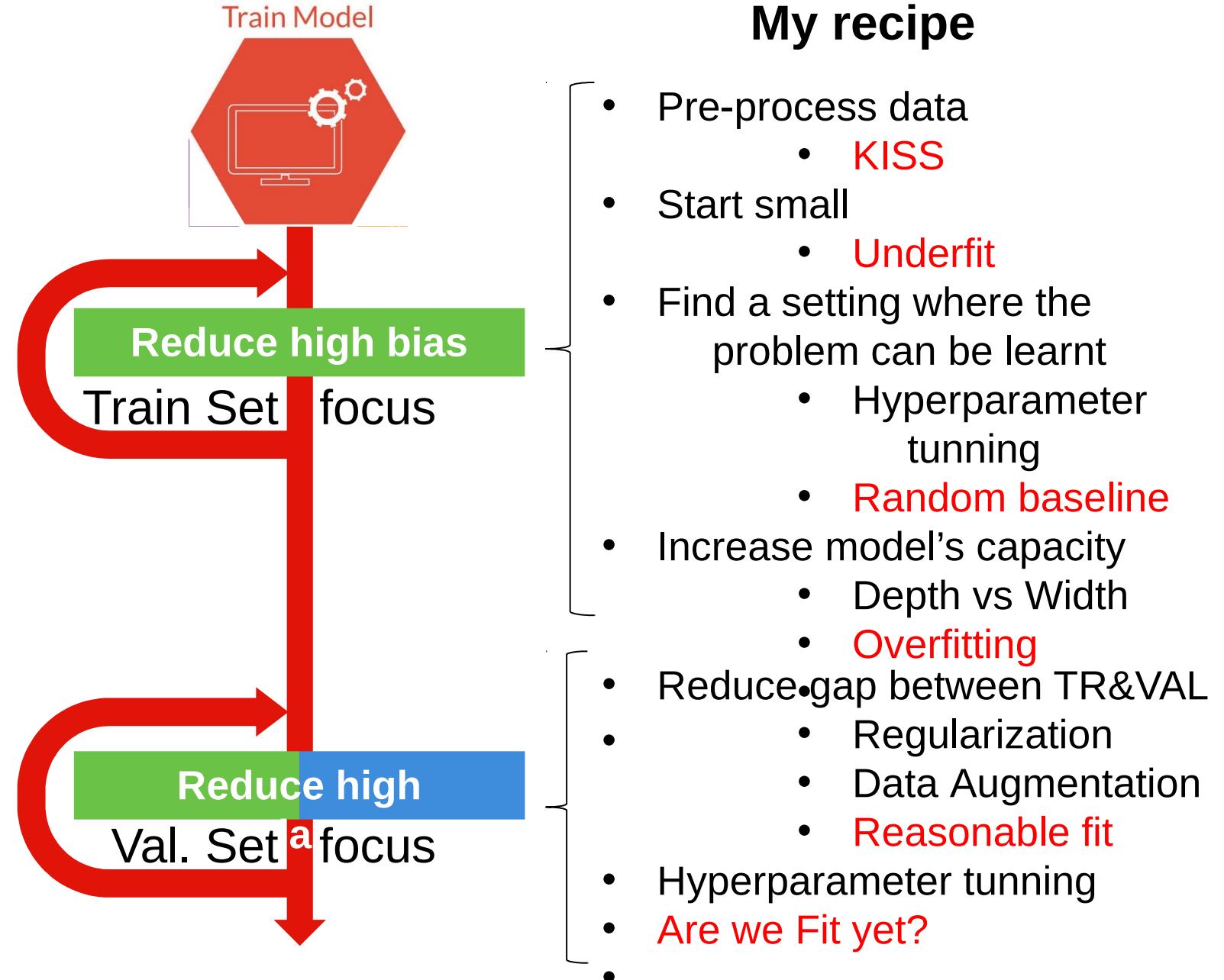
Lots of hyper-parameters:

- Network architecture:
 - # layers
 - # neurons
 - Activation function
 - ...
- Learning rate
- Optimization algorithm
- Generalization methods
- Initialization
- Data pre-processing
- ...



Practical Aspects

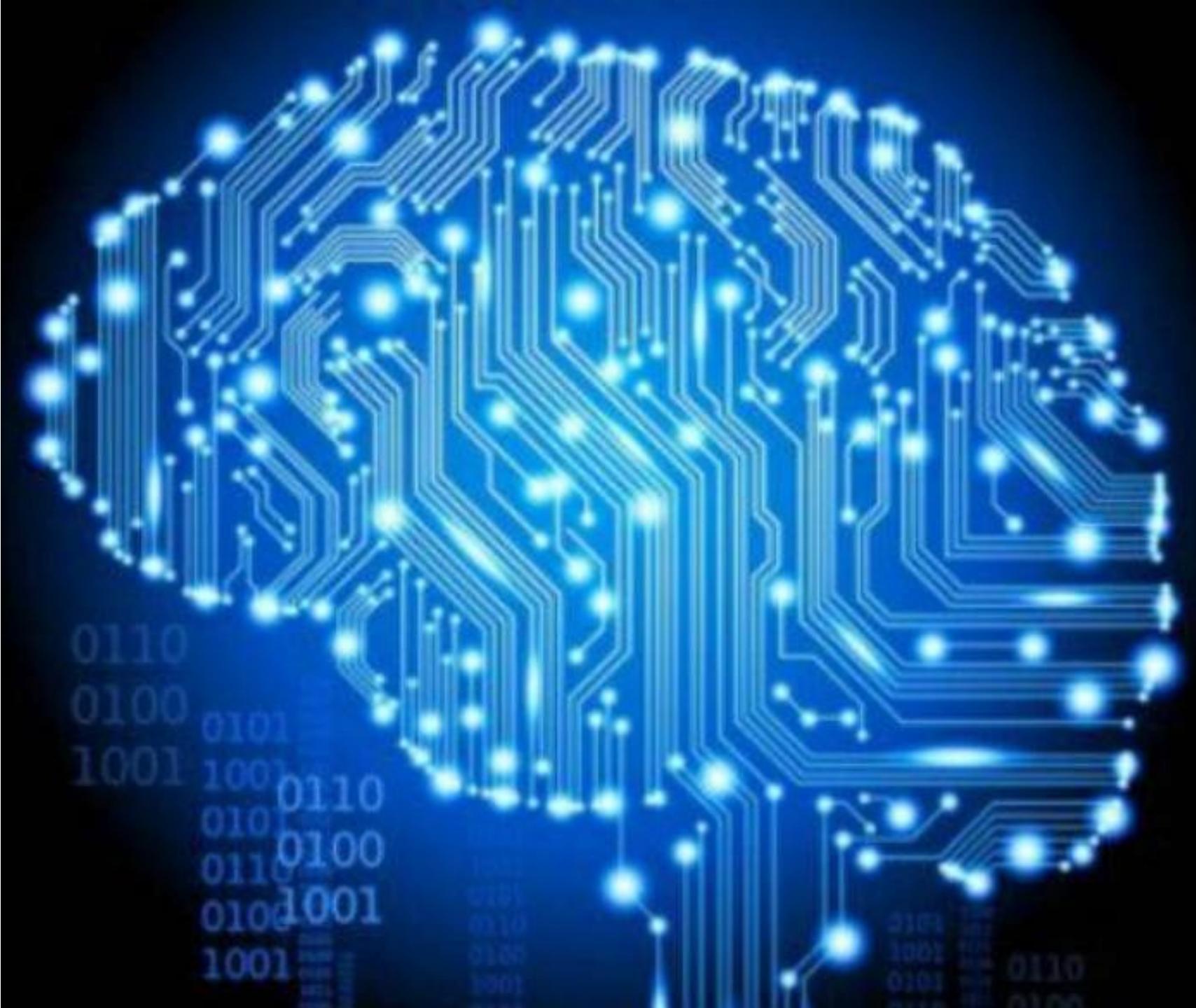
Train / Validation / Test Workflow



Convolutional Neural Networks



HIGH PERFORMANCE
ARTIFICIAL INTELLIGENCE



Convolutional Neural Networks

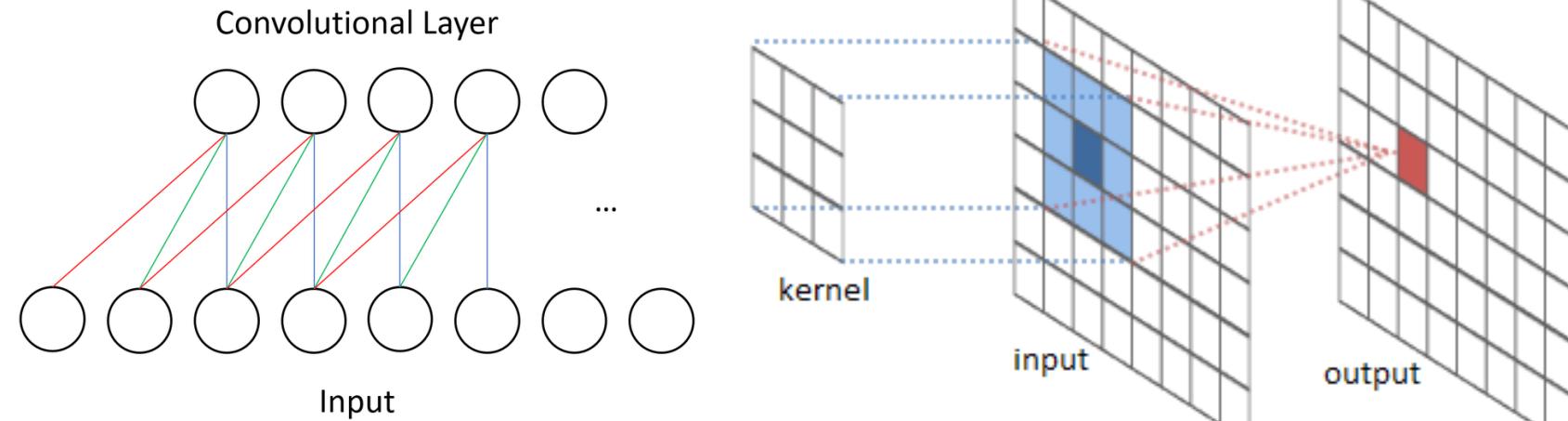
Limited connectivity
Convolution & weight sharing
Filters
Kernel size, stride and padding
Convolutional volumes
Pooling layers
Convolutional architectures
CNNs from the inside
CNN Applications



Some data has spatial correlations that could be exploited (in 1D, 2D, 3D, ...):

- Near-by data points are more relevant than far-away.

If we sparsify connectivity with a consistent purpose, we may **reduce complexity** and ease the learning of **more coherent patterns**



Convolutional Neural Networks

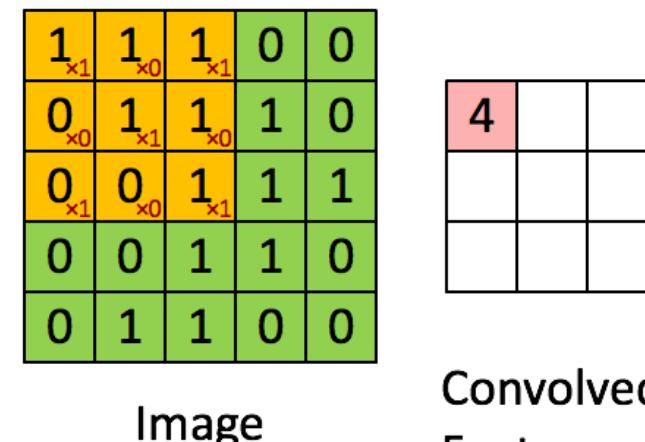
- Limited connectivity
- Convolution & weight sharing
- Filters
- Kernel size, stride and padding
- Convolutional volumes
- Pooling layers
- Convolutional architectures
- CNNs from the inside
- CNN Applications



Sparse connectivity is nice, but we still want to apply filters everywhere.

Each limited connectivity pattern (a **kernel**) will get **convolved** all over the image, generating a number of values.

Notice each kernel generates a 2D matrix of values.



In practice we have sets of neurons **sharing weights**

Convolutional Neural Networks

- Limited connectivity
- Convolution & weight sharing
- Filters**
- Kernel size, stride and padding
- Convolutional volumes
- Pooling layers
- Convolutional architectures
- CNNs from the inside
- CNN Applications



Convolution kernels can do all sorts of things on an image:

Input image



Edge detection

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$



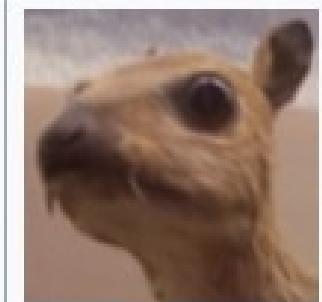
Sharpen

$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$



Gaussian blur
 3×3

$$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$



Let's let the model learn them

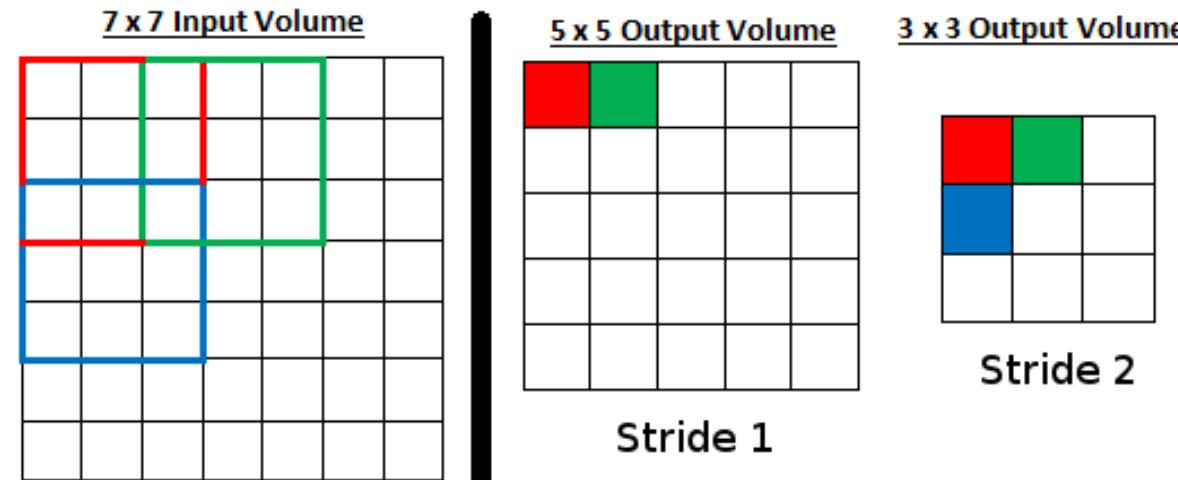
Convolutional Neural Networks

Limited connectivity
Convolution & weight sharing
Filters
Kernel size, stride and padding
Convolutional volumes
Pooling layers
Convolutional architectures
CNNs from the inside
CNN Applications



Kernel size: Size of the receptive field of convolutional neurons. Typically 3x3, 5x5, 7x7

Stride: Number of steps while convolving filter.



Stride 1 is the most common. Larger strides can replace pooling.

Padding: Border added to center conv. everywhere

- No padding: Dimensionality reduced
- Most common, zero equal/same padding

$$\text{OutputSize} = \frac{\text{InputSize} - \text{KernelSize} + 2 * \text{Padding}}{\text{Stride}} + 1$$

Convolutional Neural Networks

**Limited connectivity
Convolution & weight
sharing**

Filters

Kernel size, stride and padding

Convolutional volumes

Pooling layers

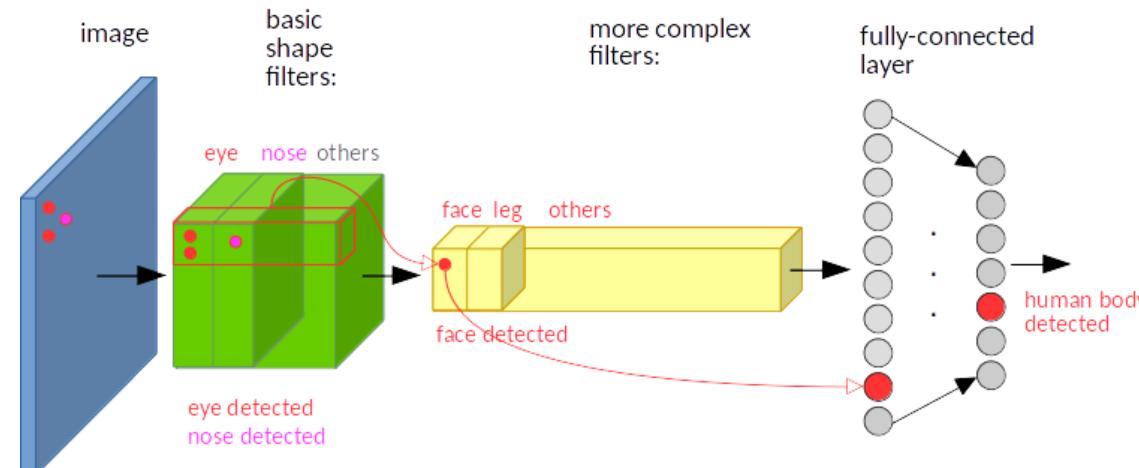
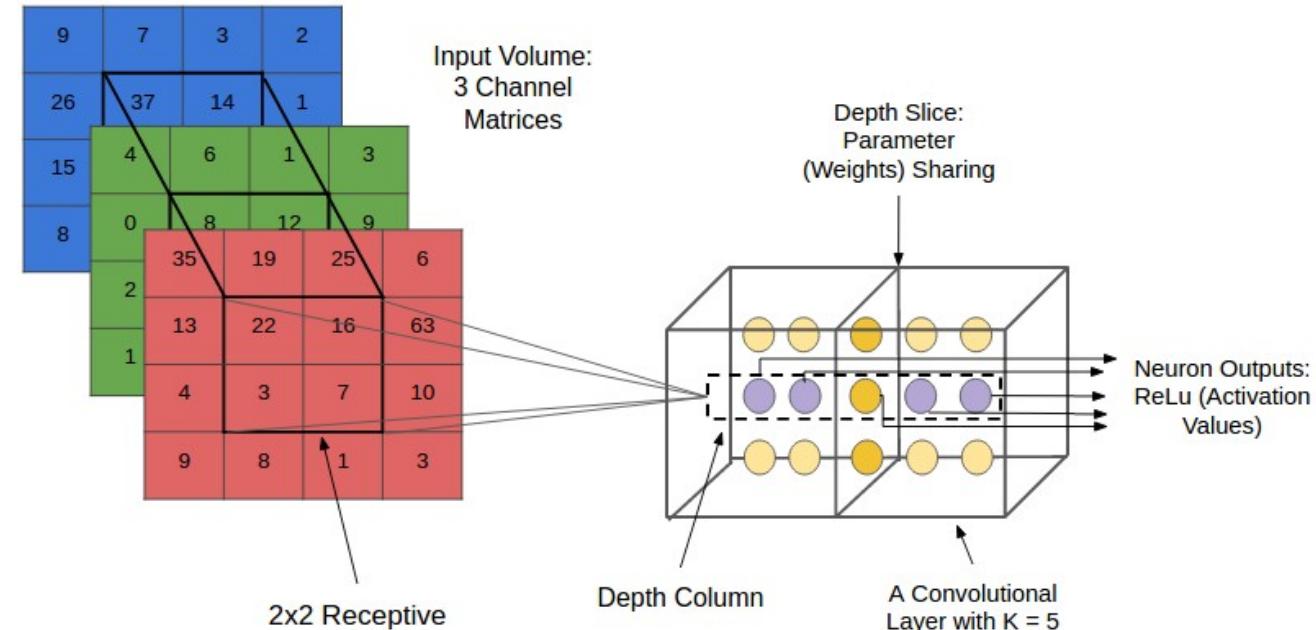
Convolutional architectures

CNNs from the inside

CNN Applications



- In a typical 2D CNN, conv filters are 3D (full depth).
 - Each filter convolved generates a 2D plane of data.
 - Depth provides all the neural views on a part of data



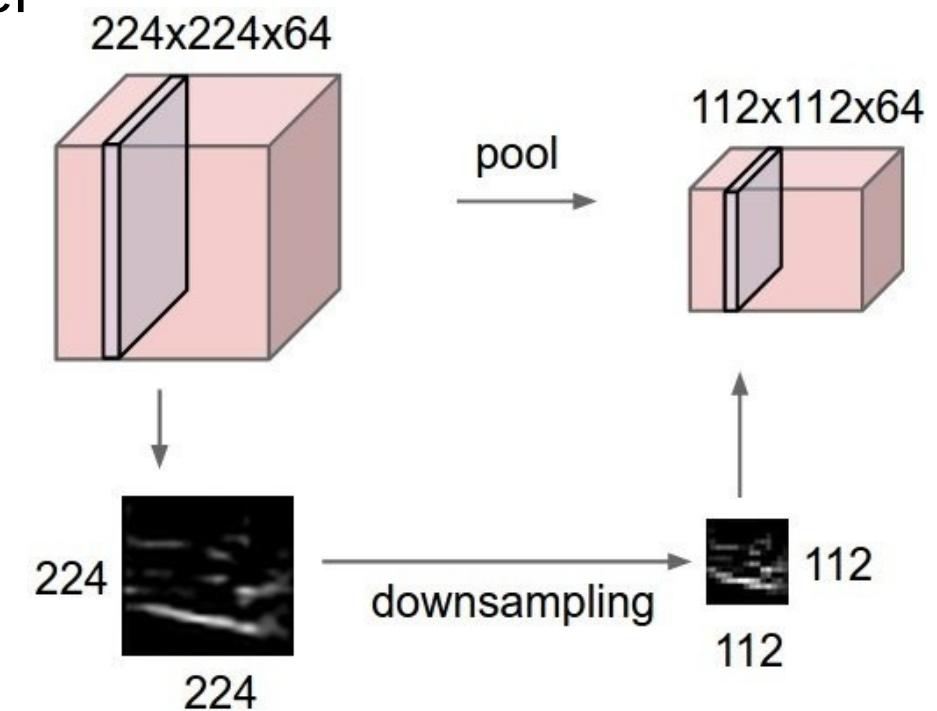
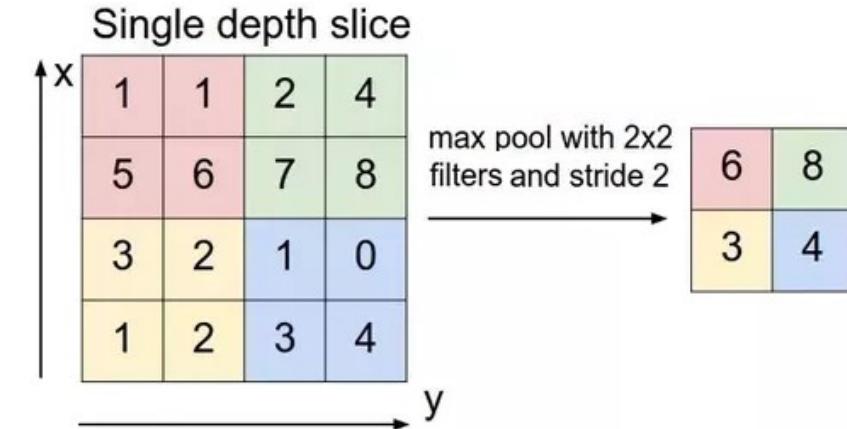
Convolutional Neural Networks

Limited connectivity
Convolution & weight sharing
Filters
Kernel size, stride and padding
Convolutional volumes
Pooling layers
Convolutional architectures
CNNs from the inside
CNN Applications



Pooling:

- Small spatial invariance
- Dimensionality reduction (along x and y only)
- Never applied full depth!
- Parameter free layer
- Hyperparams:
 - Size & Stride
 - Loss in precision
 - Max >> Avg
 -
 -



Convolutional Neural Networks

Limited connectivity
Convolution & weight sharing
Filters
Kernel size, stride and padding [AlexNet,12] [VGG,14]
Convolutional volumes
Pooling layers
Convolutional architectures
CNNs from the inside
CNN Applications



The first influential architecture was **AlexNet**:

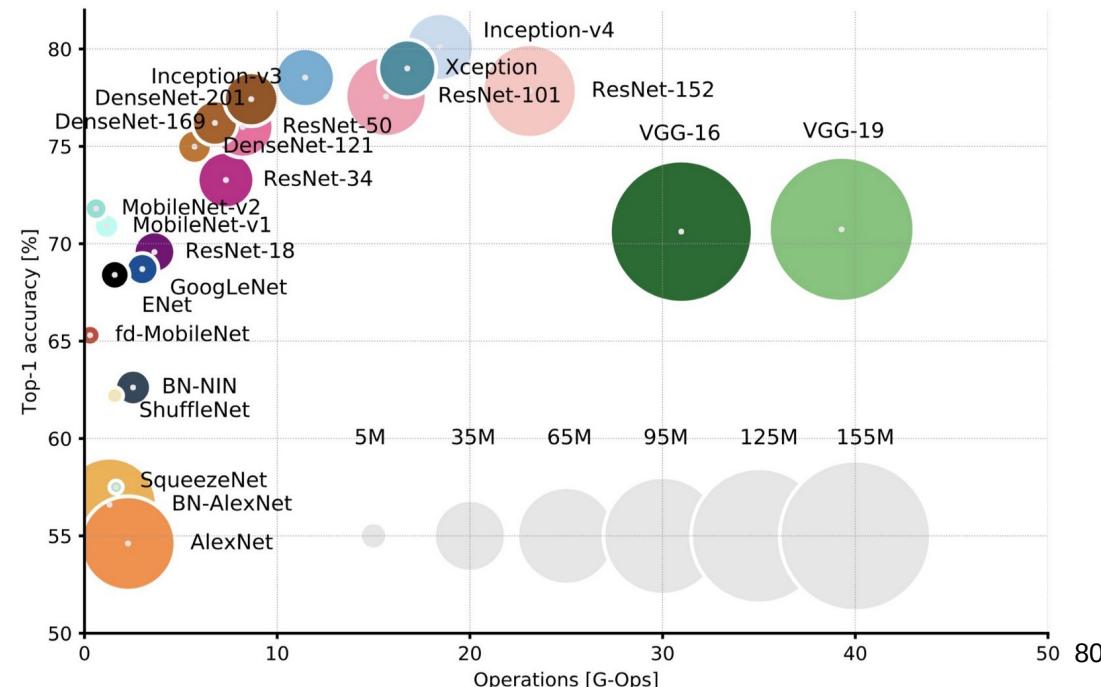
- 5 layers using convs, pools, *ReLU*, 2 dense, and *dropout*.
- 62M parameters

VGG16/19 extends the (conv-pool)*dense design:

- Smaller, 3x3 filters, but more
- 138M parameters
-

Some design principles: **KISS, be repetitive & pyramidal**

Bigger is not better!



Convolutional Neural Networks

Limited connectivity
Convolution & weight sharing
Filters
Kernel size, stride and padding
Convolutional volumes
Pooling layers
Convolutional architectures
CNNs from the inside
CNN Applications

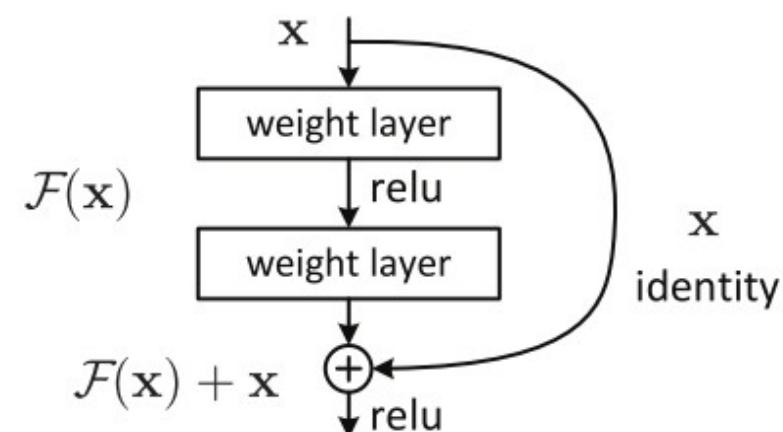
[Inception,15]
[ResNet,16]
[Huang,16]
[Xu,WWW]

But deeper should never be worse!

- In theory, one more layer is innocuous
- In practice, it's often terrible. Why?
 - Learning not alter anything is harder than it seems.
 - Its easier to learn to say always NO.

ResNet: Learning zero is easier than learning the identity function

- We can now train a 1K layer net
- They call them skip connections and residual blocks



Convolutional Neural Networks

Limited connectivity
Convolution & weight sharing
Filters
Kernel size, stride and padding
Convolutional volumes
Pooling layers
Convolutional architectures
CNNs from the inside
CNN Applications

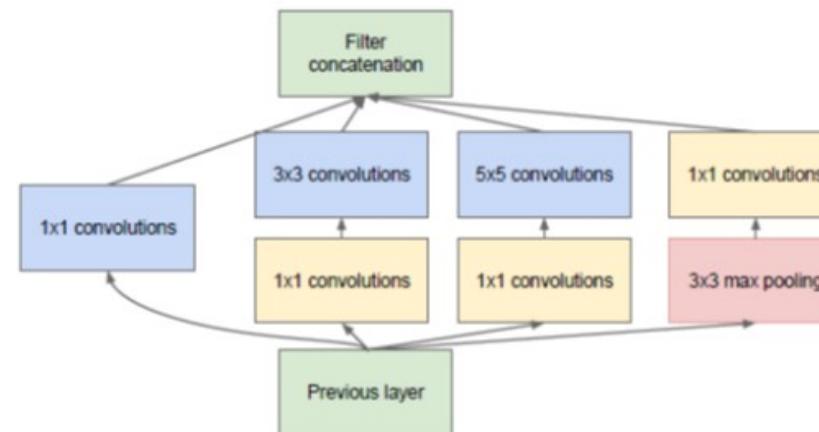
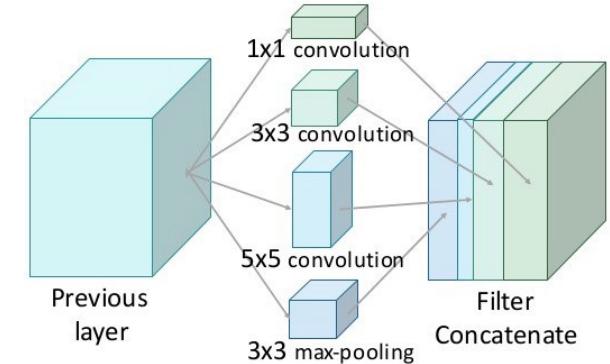


[Inception,15]
[ResNet,16]
[Huang,16]
[Xu,WWW]

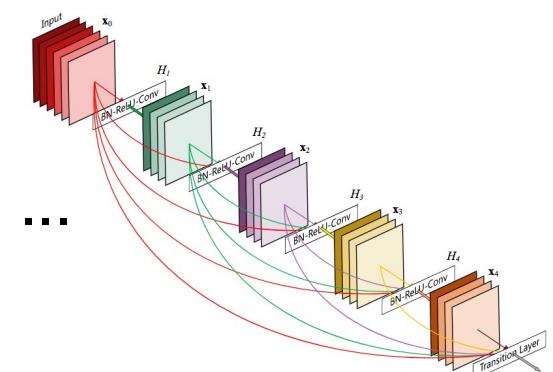
Inception:

How to choose filter size?

- Let the net decide which is best
- But bigger filters are expensive...
- Use 1x1 convs to reduce depth
- And fully connecteds require many parameters...
- Use global average pooling before to reduce width and height



DenseNet and others build on top, ...



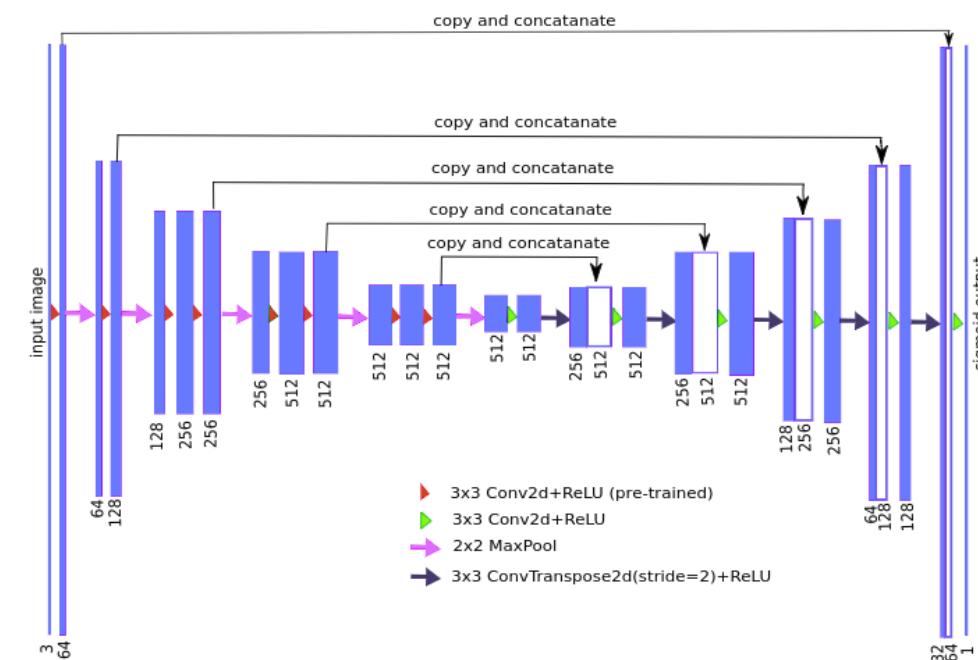
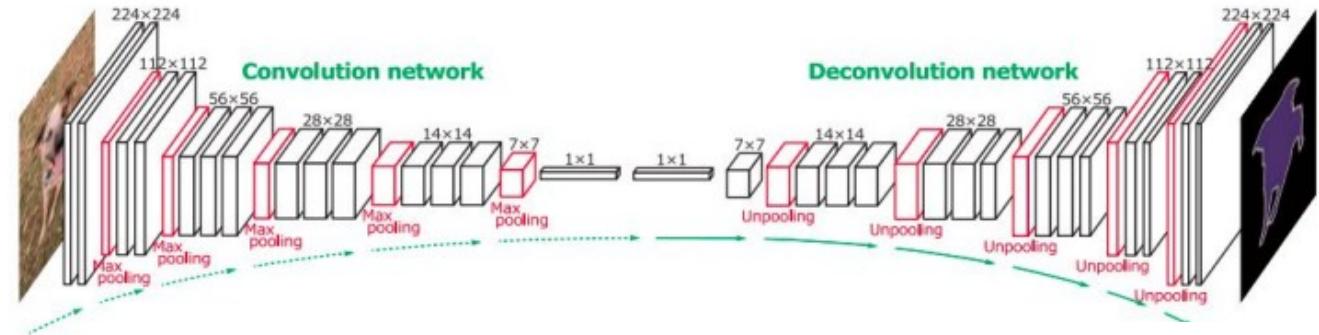
Convolutional Neural Networks

Limited connectivity
Convolution & weight sharing
Filters
Kernel size, stride and padding
Convolutional volumes
Pooling layers
Convolutional architectures
CNNs from the inside
CNN Applications



Different architectures can be done...

- Convolution – Transposed convolution (pixel-wise)



Convolutional Neural Networks

Limited connectivity

Convolution & weight sharing

Filters

Kernel size, stride and padding

Convolutional volumes

Pooling layers

Convolutional architectures

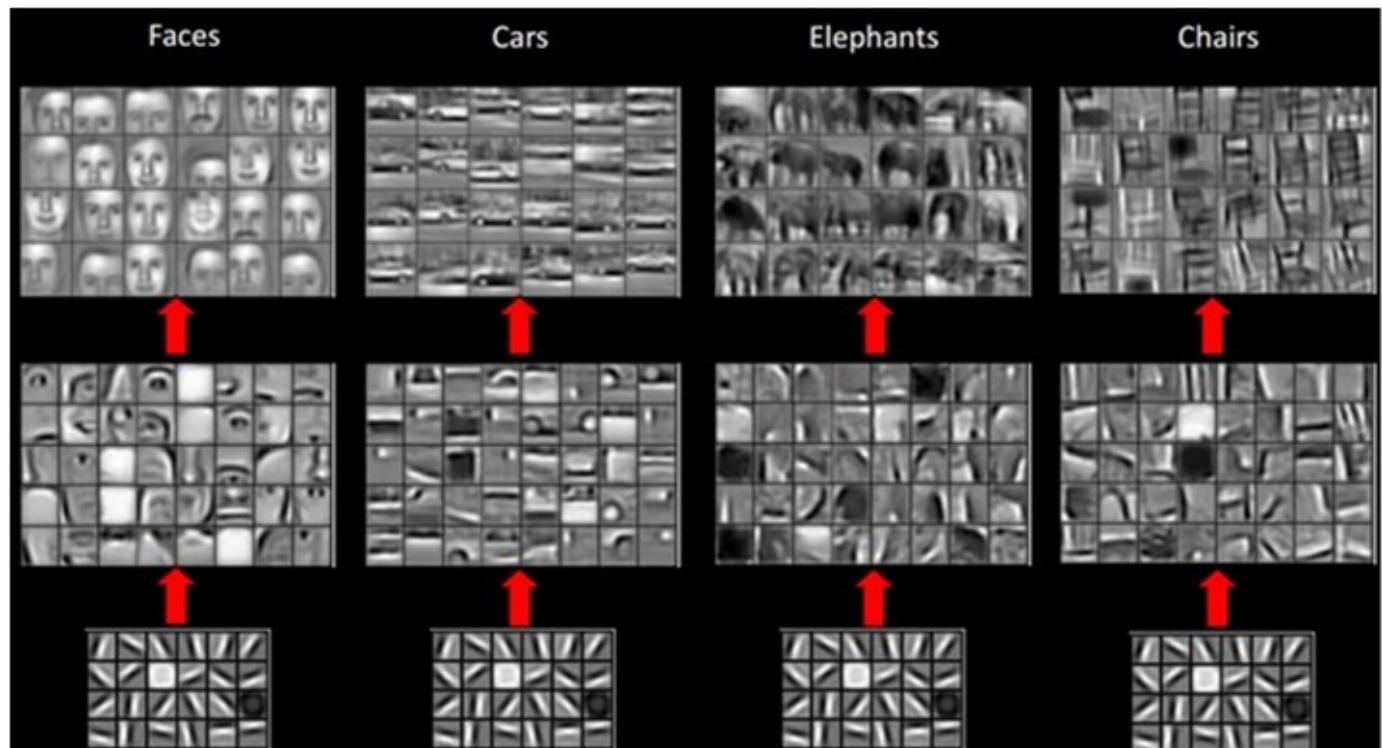
CNNs from the inside

CNN Applications



What do filters learn?

- More depth more complexity
-
- Close to input: More relevant to data, less to task
- Close to output: More relevant to task, less to data



Convolutional Neural Networks

Limited connectivity

Convolution & weight sharing

Filters

Kernel size, stride and padding

Convolutional volumes

Pooling layers

Convolutional architectures

CNNs from the inside

CNN Applications



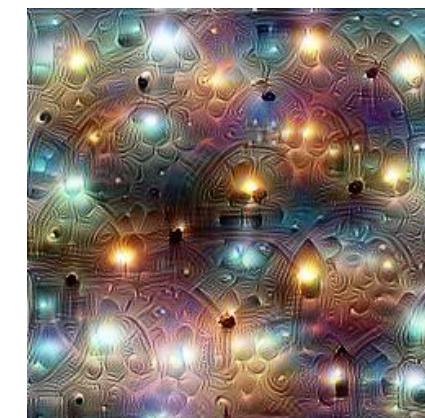
- Low-level, channel wise



- Channel-wise vs neuron-wise



- High-level



Convolutional Neural Networks

Limited connectivity

Convolution & weight

sharing

Filters

Kernel size, stride and
padding [Gatys,15]
[Kiros,14]

Convolutional volumes

Pooling layers

Convolutional
architectures

CNNs from the inside

CNN Applications



Style transfer



there is a cat
sitting on a shelf .

a plate with a fork
and a piece of cake .

a black and white
photo of a window .

a young boy standing
on a parking lot
next to cars .

a wooden table
and chairs arranged
in a room .



a kitchen with
stainless steel
appliances .

this is a herd
of cattle out
in the field .

a car is parked
in the middle
of nowhere .

a ferry boat on
a marina with a
group of people .

a little boy with
a bunch of friends
on the street .



a giraffe is standing
next to a fence
in a field .
(hallucination)

the two birds are
trying to be seen
in the water .
(counting)

a parked car while
driving down the road .
(contradiction)

the handlebars
are trying to ride
a bike rack .
(nonsensical)

a woman and
a bottle of wine
in a garden .
(gender)

Multimodal pipelines

Convolutional Neural Networks

Limited connectivity

Convolution & weight

sharing

Filters

Kernel size, stride and
padding

Convolutional volumes

Pooling layers

Convolutional
architectures

CNNs from the inside

CNN Applications



- Image colorization

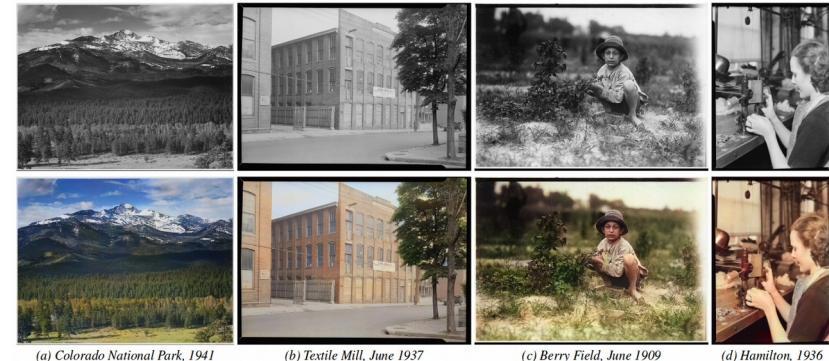


Image segmentation

