

# Artificial Neural Networks (ANN)

Nikola Popović

[nipopovic@vision.ee.ethz.ch](mailto:nipopovic@vision.ee.ethz.ch)



**CVL** Computer  
Vision  
Lab

**ETH** zürich

# 1st wave of NN – The birth of the idea

- [1943] First artificial neuron model – W. McCulloch , W. Pitts
- [1949] Hebb's learning law – D. Hebb
- [1958] Perceptron – F. Rosenblatt
- [1962] Delta Learning Rule – B. Widrow, T. Hoff

## 2nd wave – Excitement again

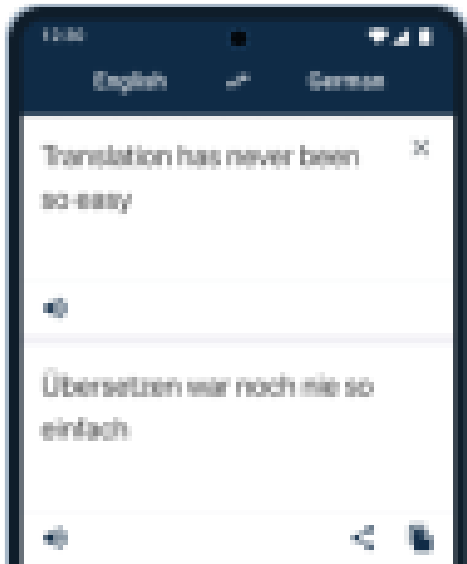
- [1986] Backpropagation popularized - D. Rumelhart, G. E. Hinton, R. Williams
- [1997] LSTM's - H. Sepp, J. Schmidhuber
- [1998] CNN's— Y. LeCun, L. Bottu, Y. Bengio, P. Haffner

# 3rd wave – ANN seem to work

- Much better hardware and software
- Much bigger data sets
- Some influential works
  - [2006] Deep Belief Networks – G. E. Hinton, S. Osindero, Y.-W. Teh
  - [2009] Speech processing – G. E. Hinton, L. Deng
  - [2012] Image classification – A. Krizhevsky, I. Sutskever, G. E. Hinton
- Many astonishing results followed

# Deep Learning applications

Machine translation



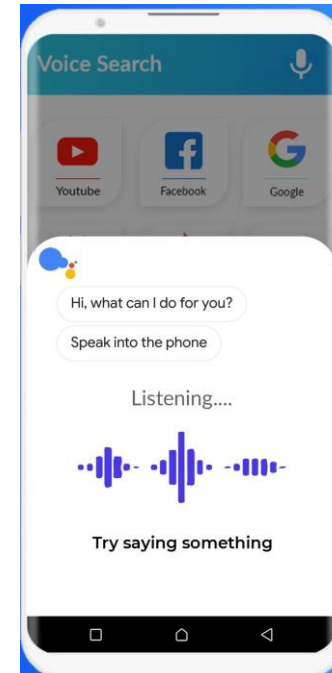
<https://www.deepl.com/translator>

Autonomous cars



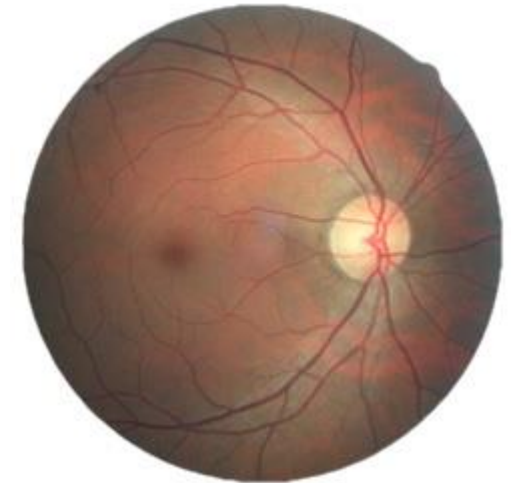
<https://www.ptolemus.com/topics/autonomous-vehicles/>

Speech recognition



[https://play.google.com/store/apps/details?id=com.prometheusinteractive.voice\\_launcher](https://play.google.com/store/apps/details?id=com.prometheusinteractive.voice_launcher)

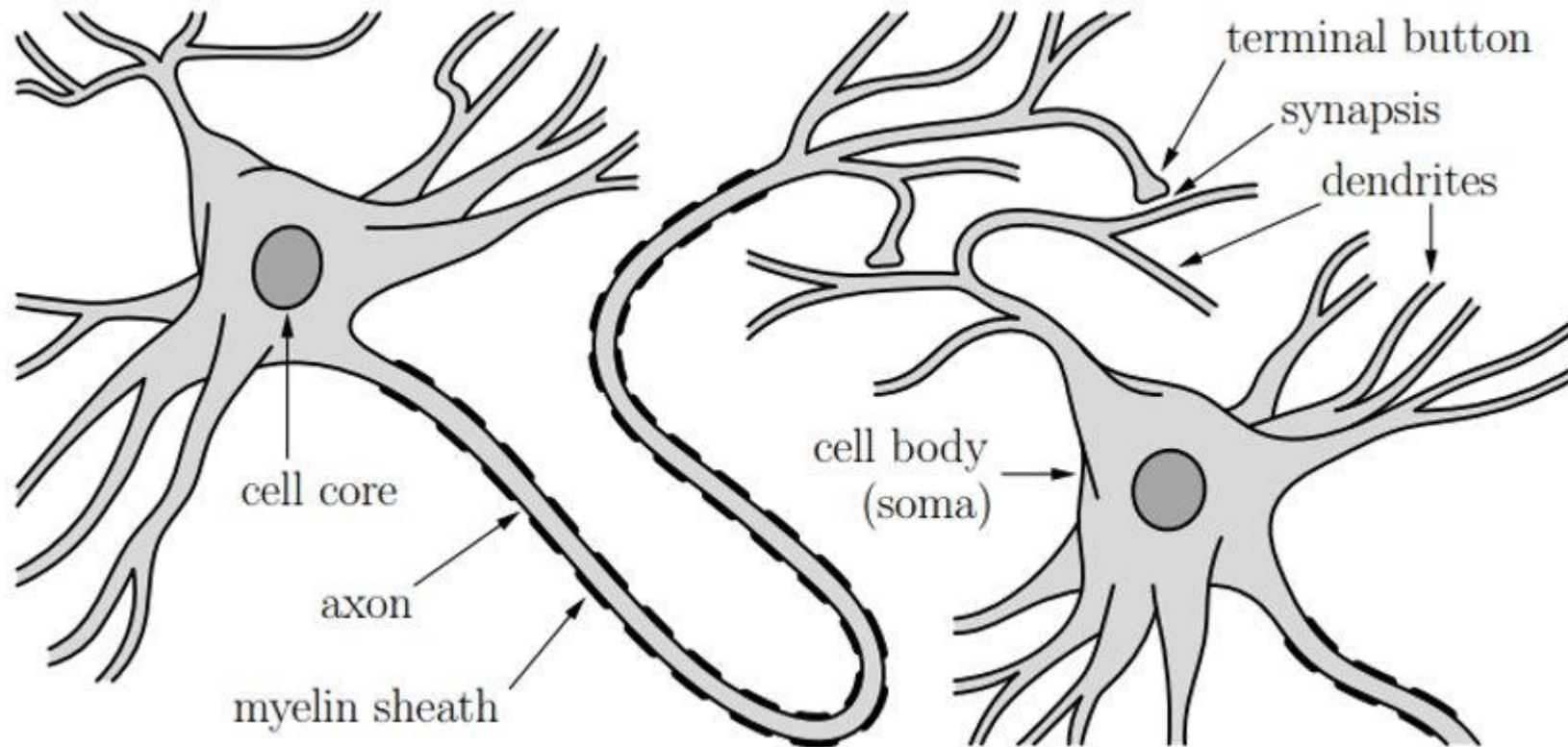
Diabetic retinopathy detection



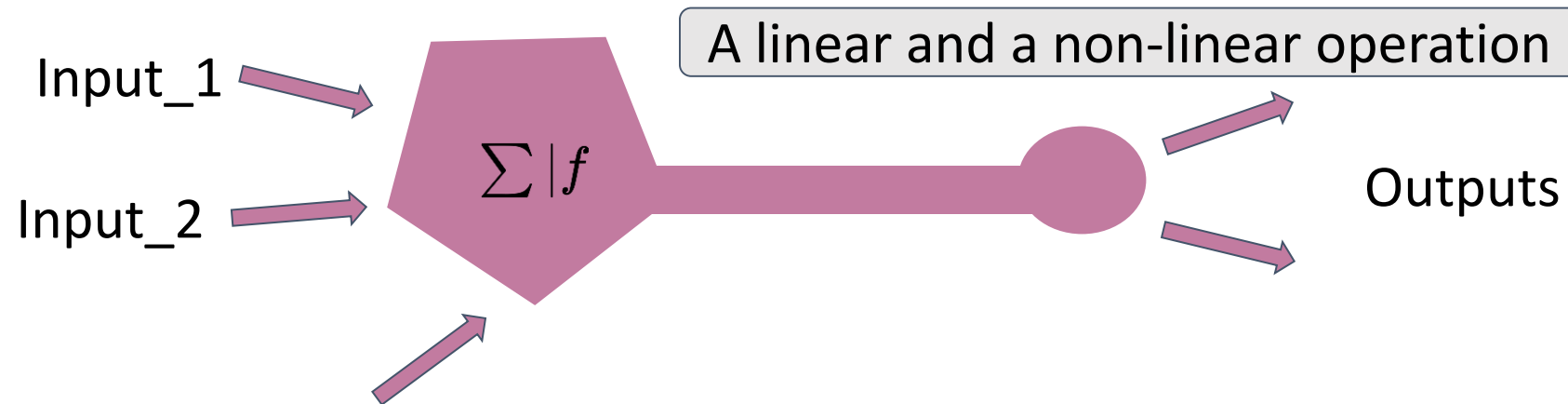
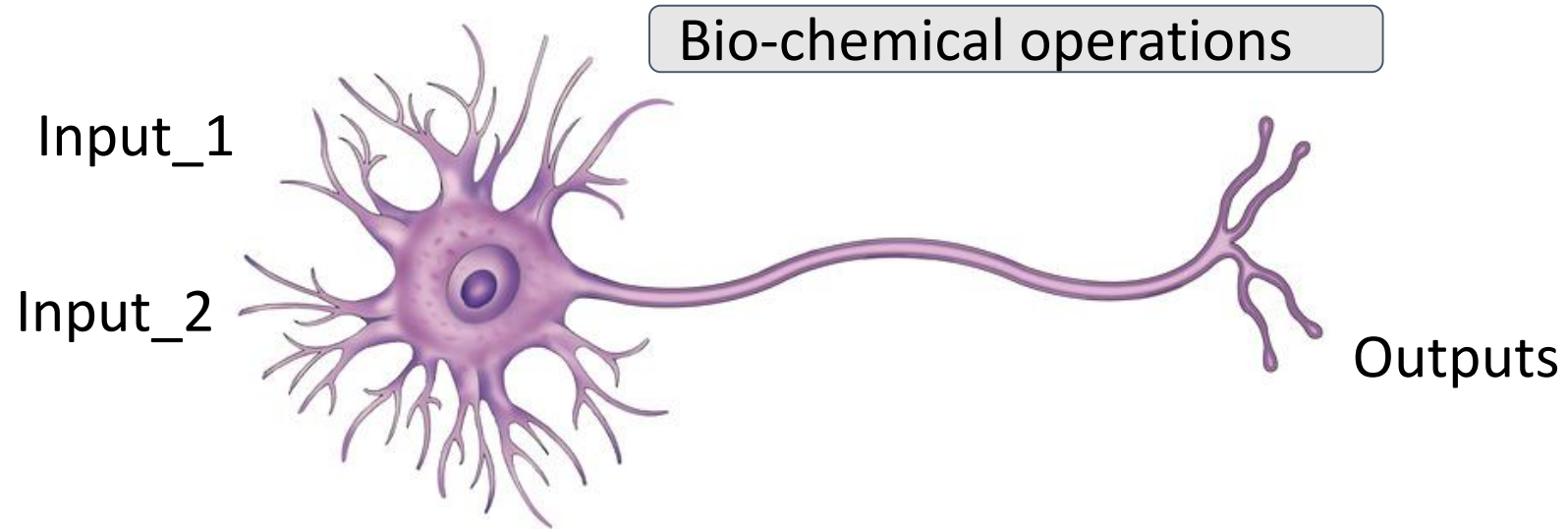
<https://www.kaggle.com/c/diabetic-retinopathy-detection>

# Biological and Artificial Neurons

- ANN were inspired by human nervous system



# Biological and Artificial Neurons



# One Artificial Neuron

- Input  $\mathbf{x} = [x_1, \dots, x_m]^T$
- Weights  $\mathbf{w} = [w_1, \dots, w_m]$
- Bias  $b$
- Pre-activation  $z = \sum_i x_i w_i + b = \mathbf{w}\mathbf{x} + b$
- Activation fn.  $f(\cdot)$
- Neuron output  $y = f(z) = f(\sum_i x_i w_i + b)$



# ANN

## Single layer

- Input – vector  $\mathbf{x} = [x_1, \dots, x_m]^T$
- Different weights for each of k neurons

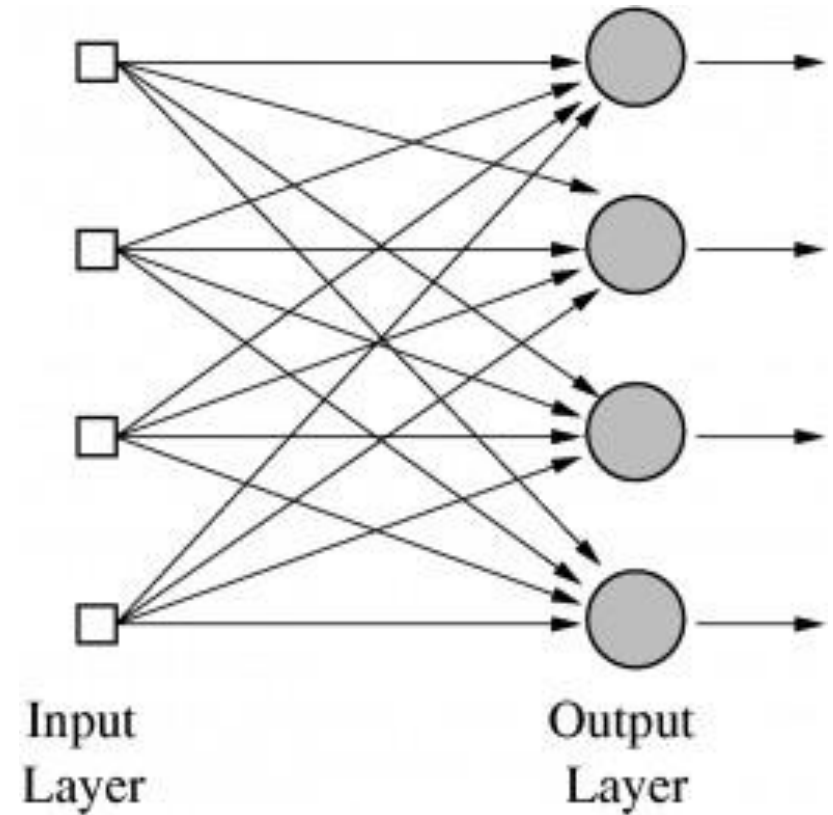
$$\mathbf{W} = \begin{bmatrix} w_{11} & \cdots & w_{1m} \\ \vdots & \ddots & \vdots \\ w_{k1} & \cdots & w_{km} \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1 \\ \dots \\ \mathbf{w}_k \end{bmatrix}$$

- Pre-activation vector

$$\mathbf{z} = [z_1, \dots, z_k]^T = \mathbf{W}\mathbf{x} + \mathbf{b}$$

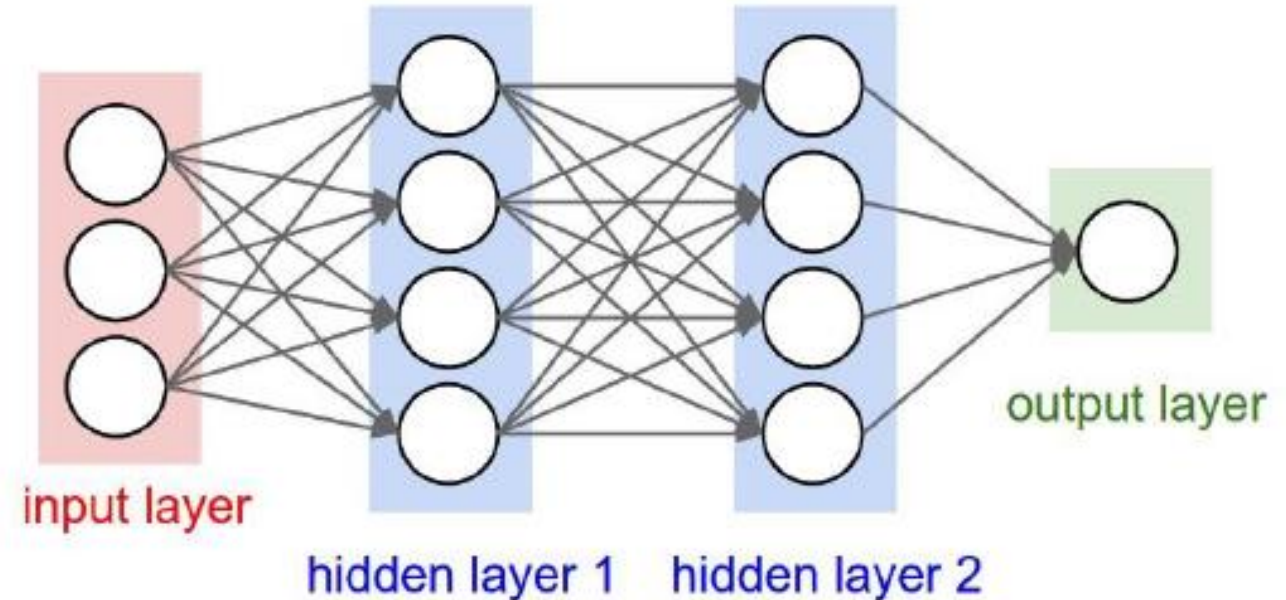
- Output vector

$$\mathbf{y} = [y_1, \dots, y_k]^T = f(\mathbf{W}\mathbf{x} + \mathbf{b})$$



# Multi Layer ANN

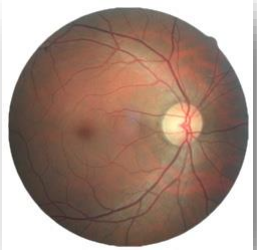
- Input layer
  - The input vector
- Hidden layer(s)
  - Intermediate computation
  - Extract features
- Output layer
  - Final computation and decision making
- Nowadays called multi layer perceptron (MLP)



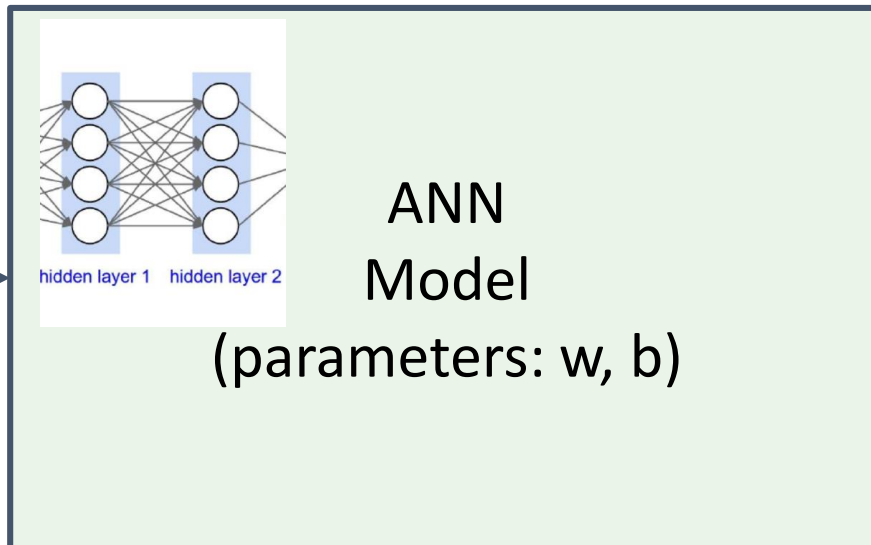
Stanford, CS231n course, Lecture 4 presentation

# High-level view

- Binary classification example



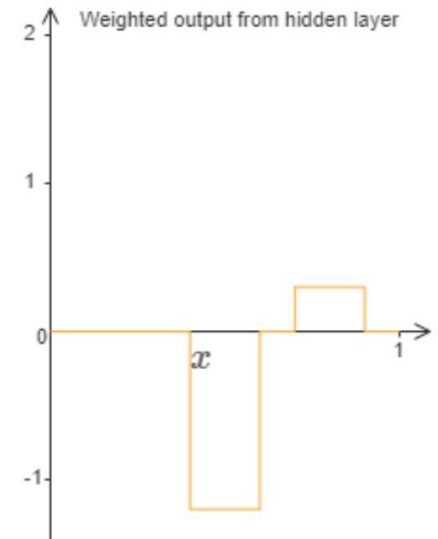
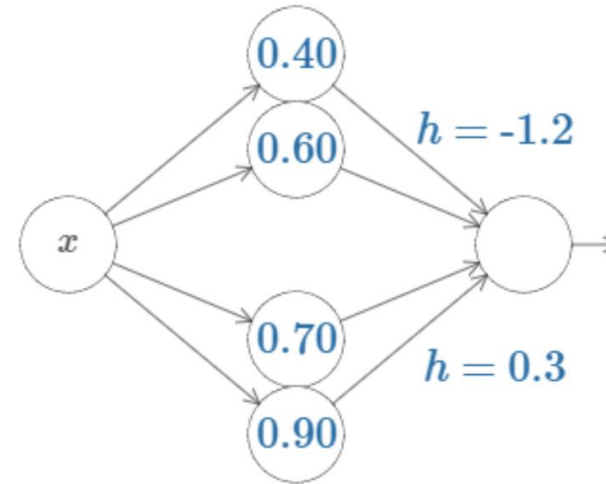
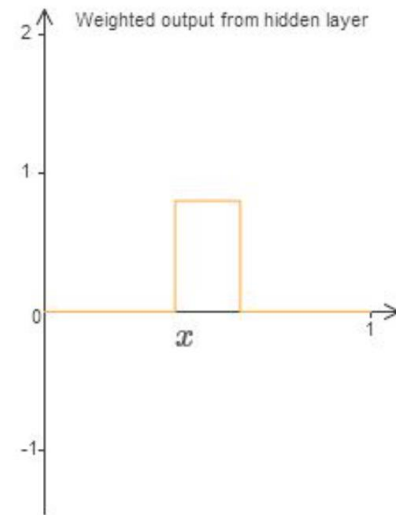
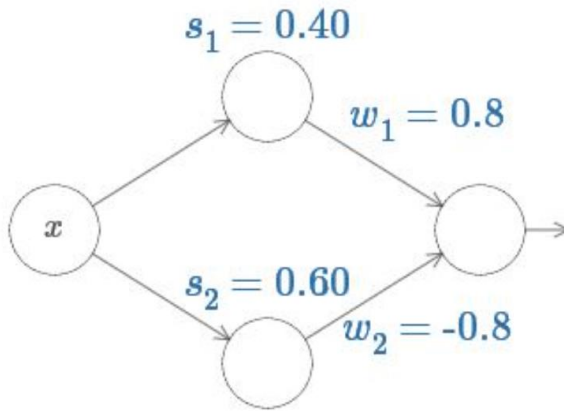
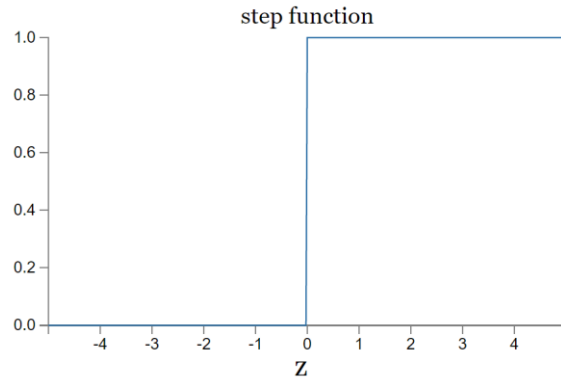
Input image  $x$   
(pixel rgb values)



$y$  = probability of an eye disease  
 $y = 0 \Rightarrow$  healthy  
 $y = 1 \Rightarrow$  disease

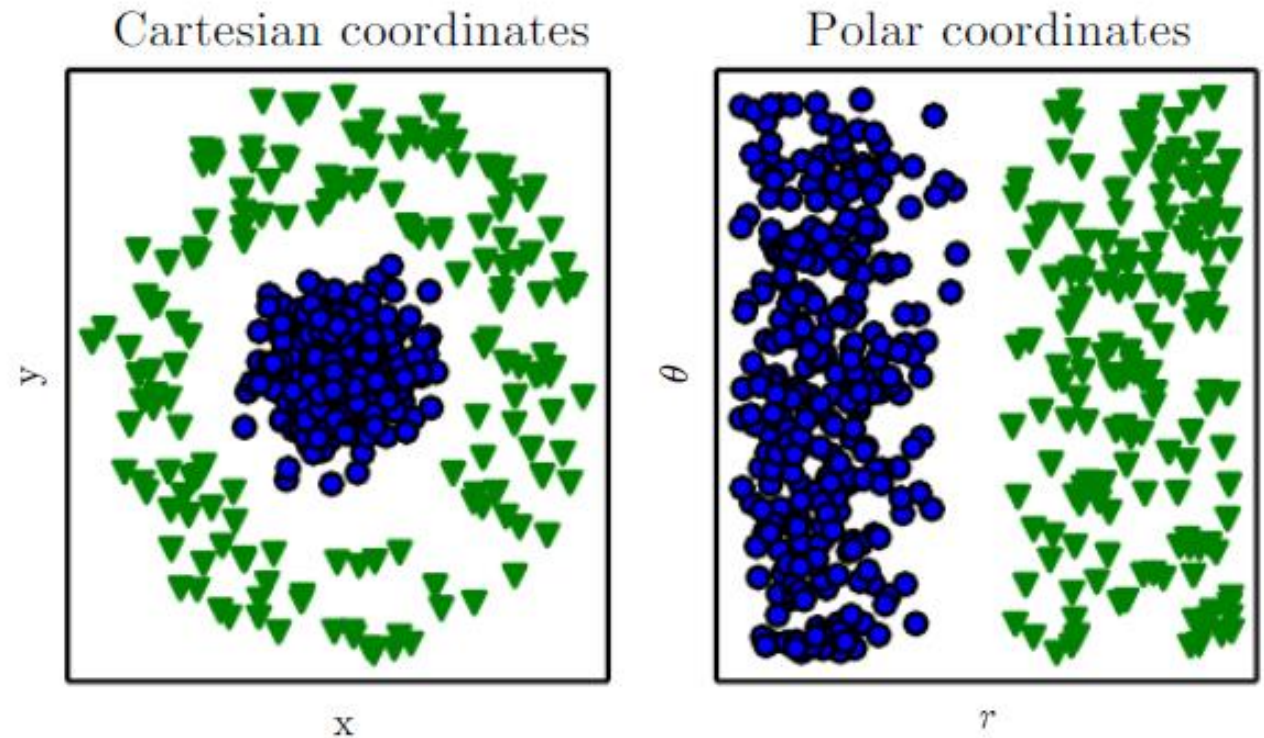
Output  $y$

# Intuition behind the power of ANN



# Motivation on need for good features

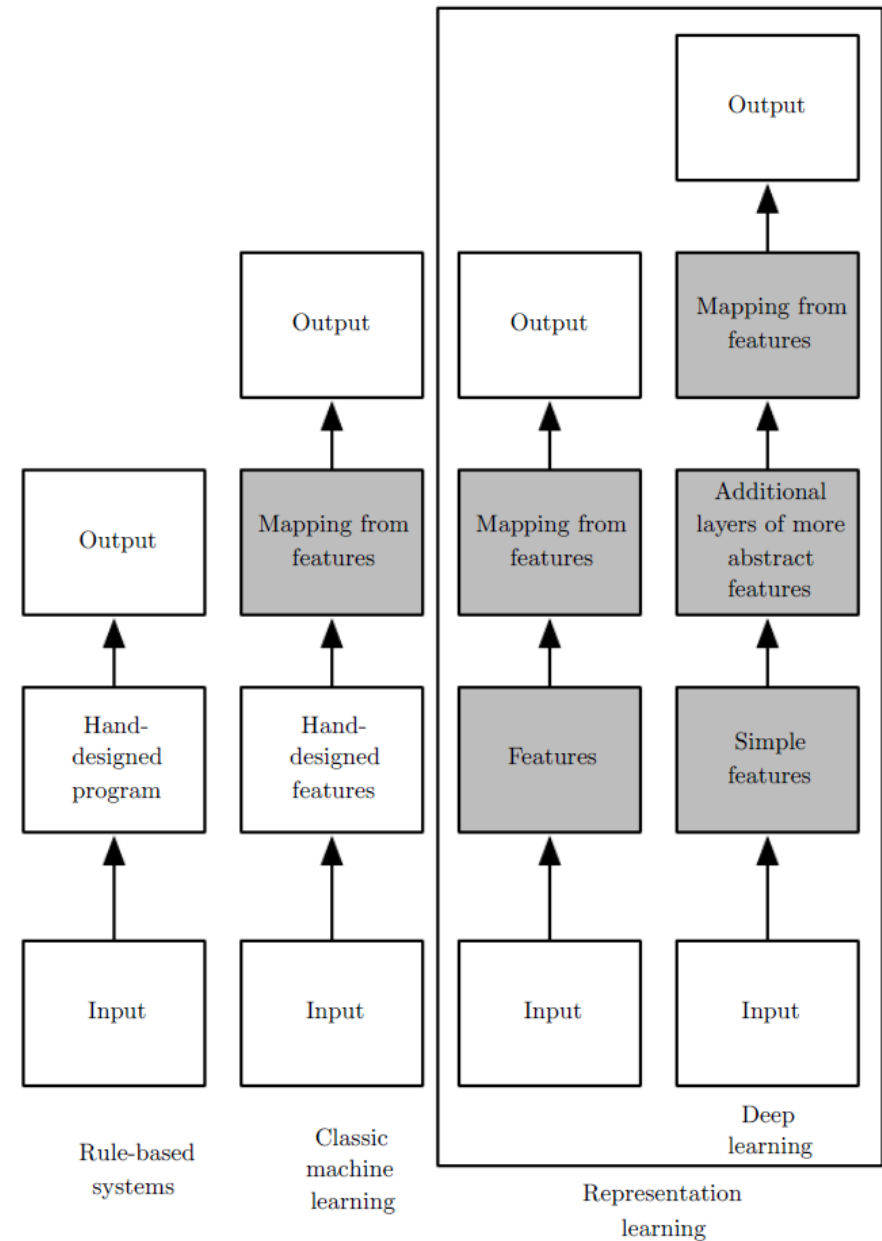
- The left plot cannot be separated with a linear classifier
- With good features (right plot) we can separate these two classes with a line



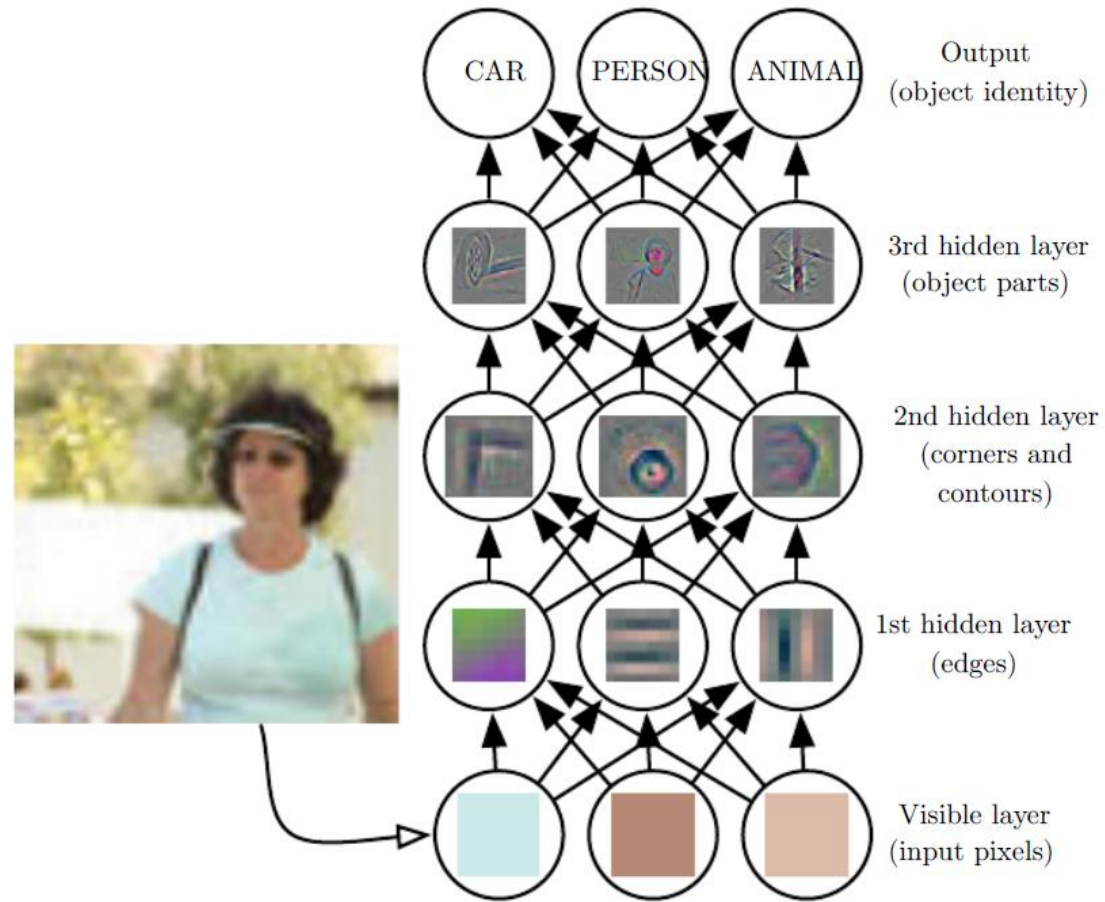
Goodfellow et al., Deep Learning, MIT Press 2016

# What is Deep Learning?

- Before – multiple stages
  - Designing features
  - Prediction algorithm (takes features as input)
- Now - Deep Learning
  - Raw data as input
  - Prediction as output

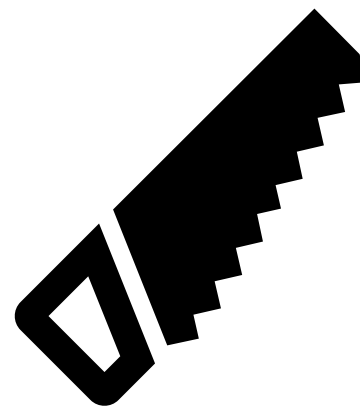
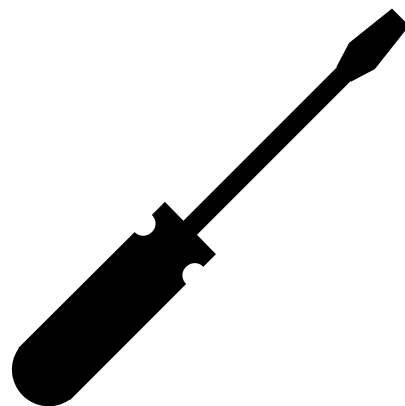
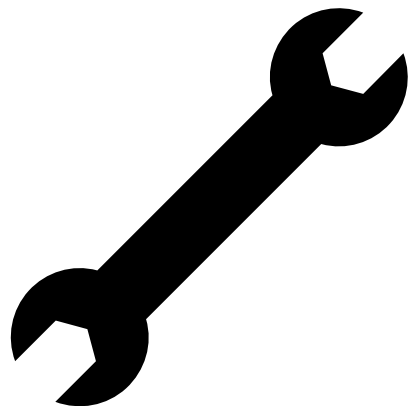
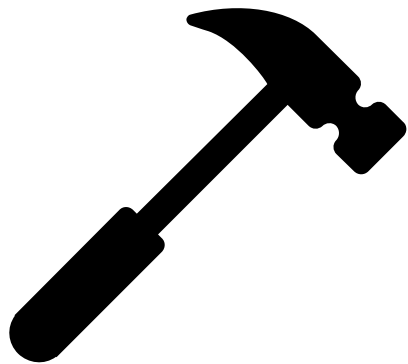


# Intuition behind Deep ANN



Goodfellow et al., Deep Learning, MIT Press 2016

ANN are tools





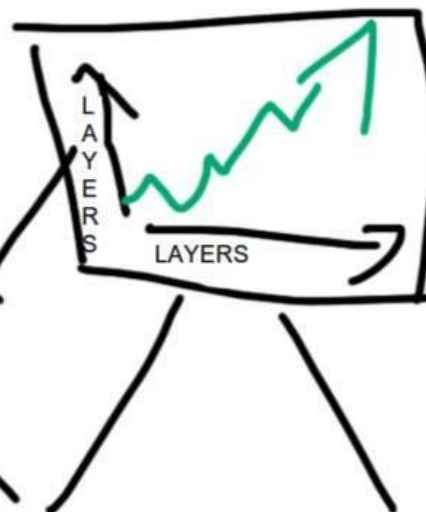
## STATISTICAL LEARNING

Gentlemen, our learner overgeneralizes because the VC-Dimension of our Kernel is too high, Get some experts and minimize the structural risk in a new one. Rework our loss function, make the next kernel stable, unbiased and consider using a soft margin

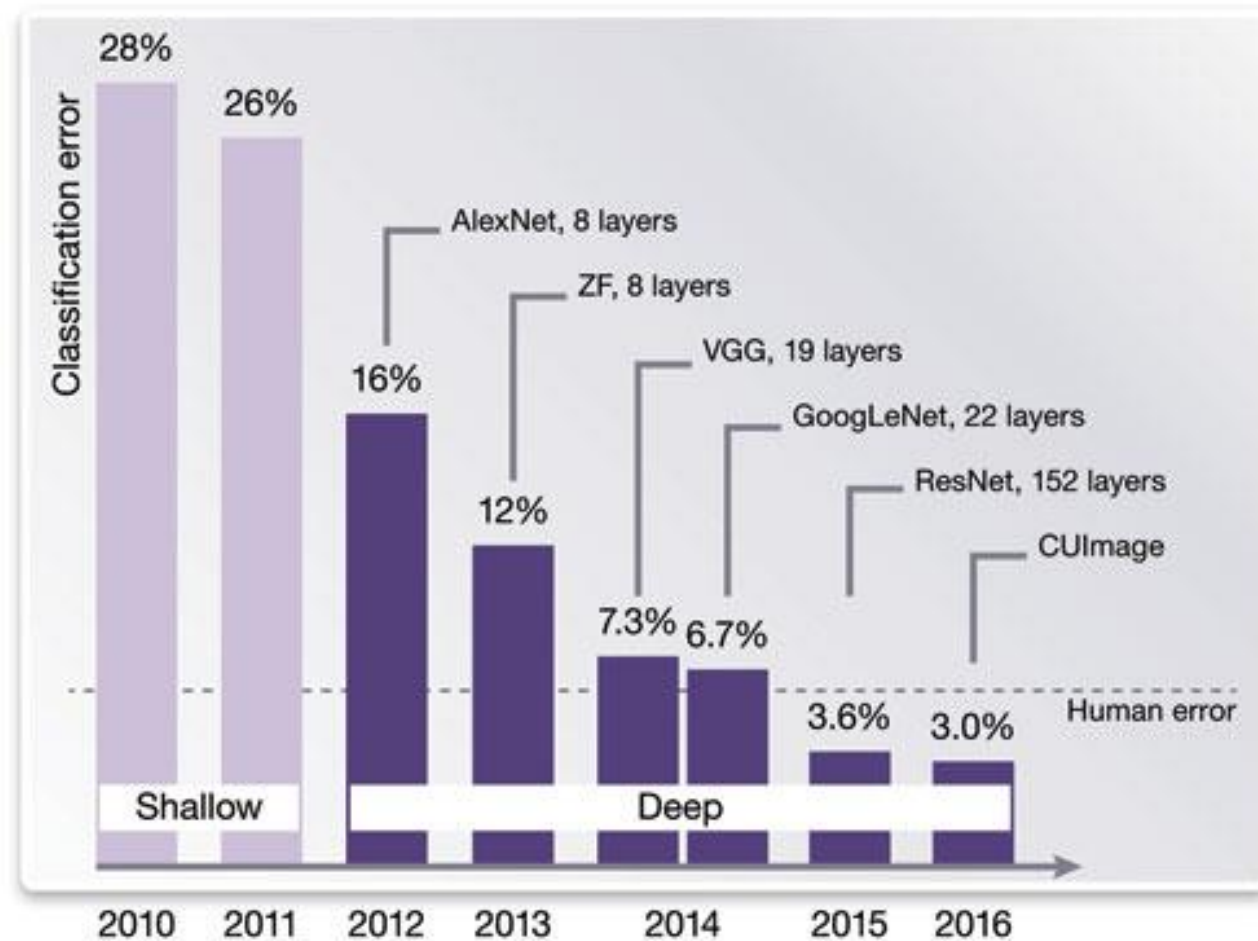


## NEURAL NETWORKS

STACK  
MORE  
LAYERS



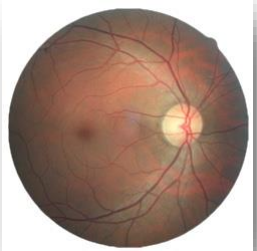
# ImageNet Challenge



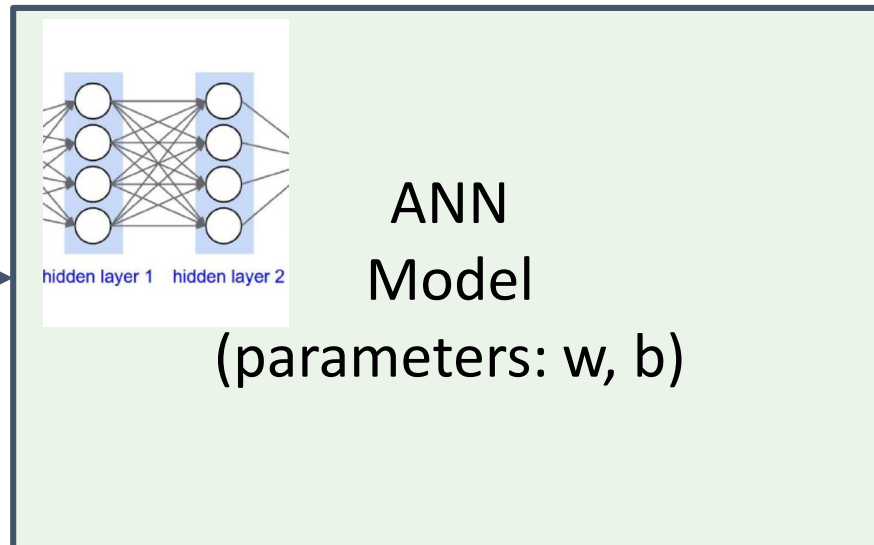
# Learning

How do we get the model parameters  $(w, b)$ ?

By designing loss functions and optimizing them.



Input image  $x$   
(pixel rgb values)



$y$  = probability of an eye disease  
 $y = 0 \Rightarrow$  healthy  
 $y = 1 \Rightarrow$  disease

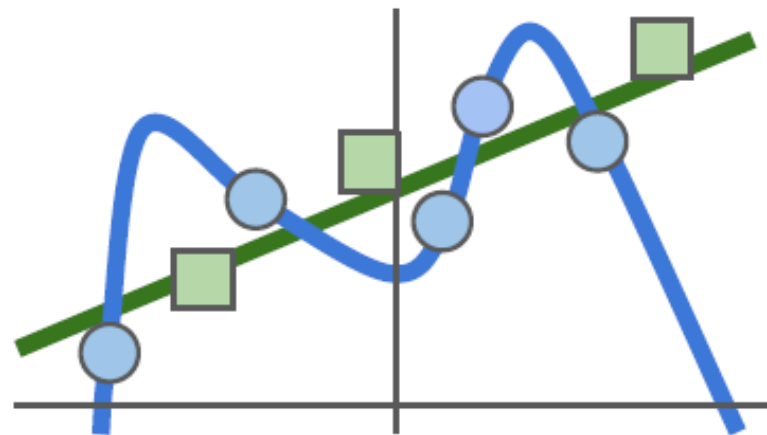
Output  $y$

# Supervised learning objective

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)$$

**Loss function**

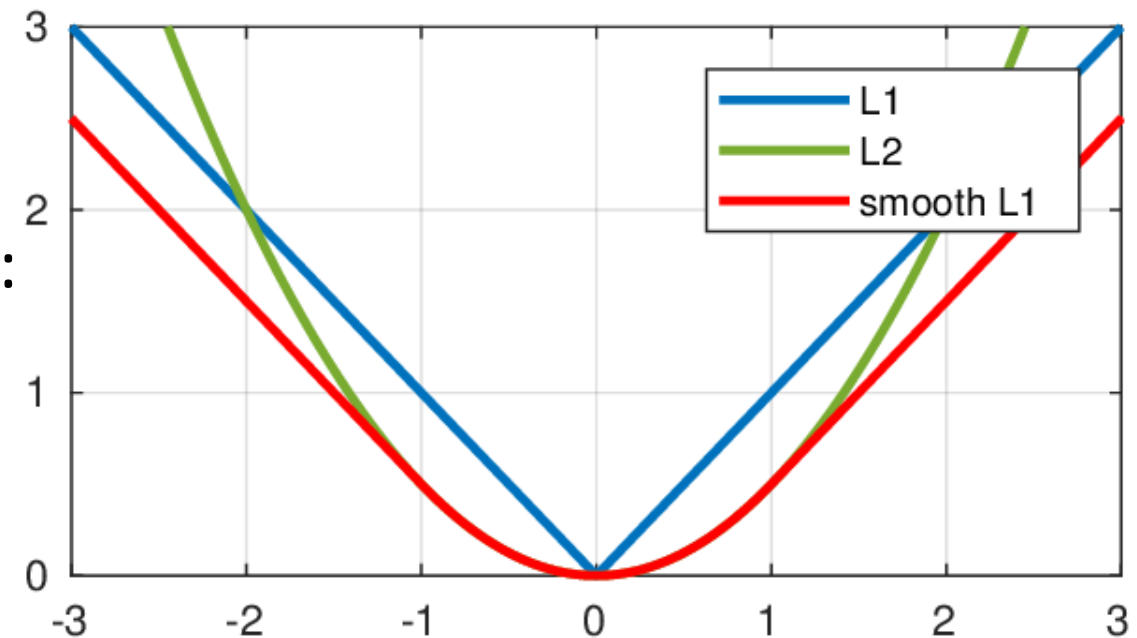
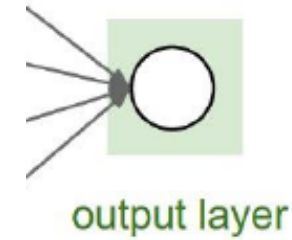
Model predictions should  
match training data



# Regression

- In case of one output:  $\hat{y}^{(i)} = f(x^{(i)}, W)$
- $L_2$  loss(MSE):  $L_i = (\hat{y}^{(i)} - y^{(i)})^2$
- $L_1$  loss:  $L_i = |\hat{y}^{(i)} - y^{(i)}|$
- Smooth  $L_1$  loss (less sensitive to outliers):

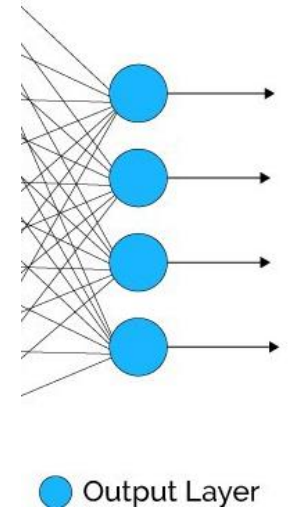
$$L_i = \begin{cases} L_2 \text{ loss} & |\hat{y}^{(i)} - y^{(i)}| \leq a \\ L_1 \text{ loss} & |\hat{y}^{(i)} - y^{(i)}| > a \end{cases}$$



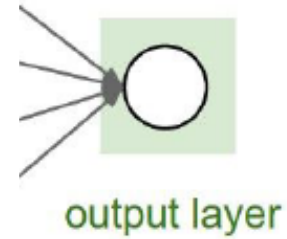
# Regression

- For multiple outputs, summ the losses for each output

$$L_i = \sum_j L_{ij}$$



# Sigmoid output layer (Classification)



- Binary classification problem with one output

- Network output:  $\hat{p} = P(Y = 1 | X = x^{(i)}) = \frac{1}{1 + e^{-net^{(i)}}}$ 
  - $net^{(i)} = f(x^{(i)}, W)$  (score for class 1)

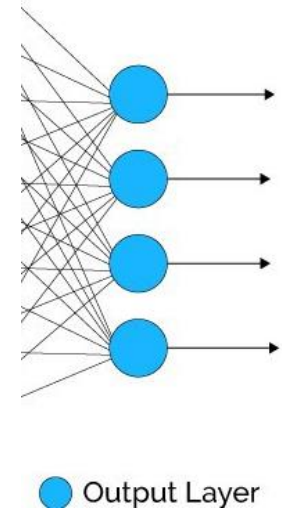
- Cross-entropy loss

$$\begin{aligned} L_i &= -\log P(Y = y^{(i)} | X = x^{(i)}) \\ &= -y^{(i)} \log(\hat{p}) - (1 - y^{(i)}) \log(1 - \hat{p}) \end{aligned}$$

# Softmax output layer (Classification)

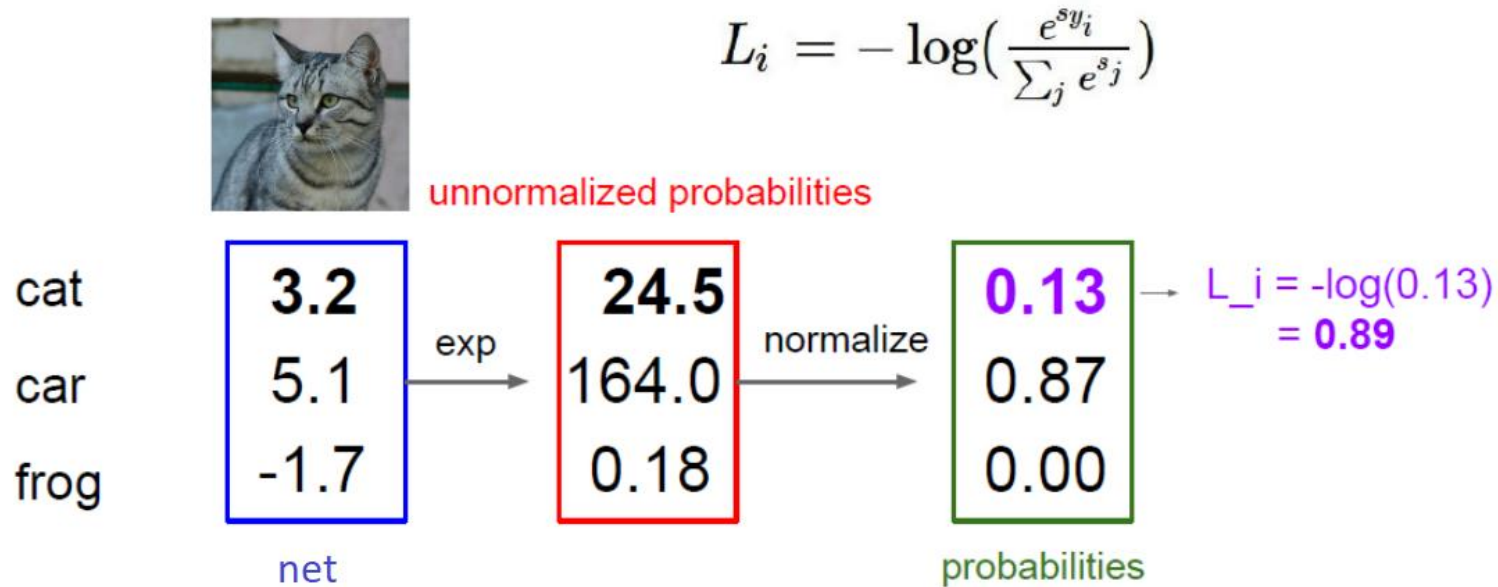
- Multi-class classes classification problem ( $m$  classes)
- Network outputs:  $\hat{p}_k = P(Y = k | X = x^{(i)}) = \frac{e^{net_k^{(i)}}}{\sum_{j=1}^m e^{net_j^{(i)}}}$ 
  - $net_k^{(i)} = f(x^{(i)}, W)$  – score for class  $k$
- Negative log likelihood

$$\begin{aligned} L_i &= -\log P(Y = y^{(i)} | X = x^{(i)}) \\ &= -\log \left( \frac{e^{net_{y^{(i)}}^{(i)}}}{\sum_{j=1}^m e^{net_j^{(i)}}} \right) \end{aligned}$$





# Softmax output layer (Classification)

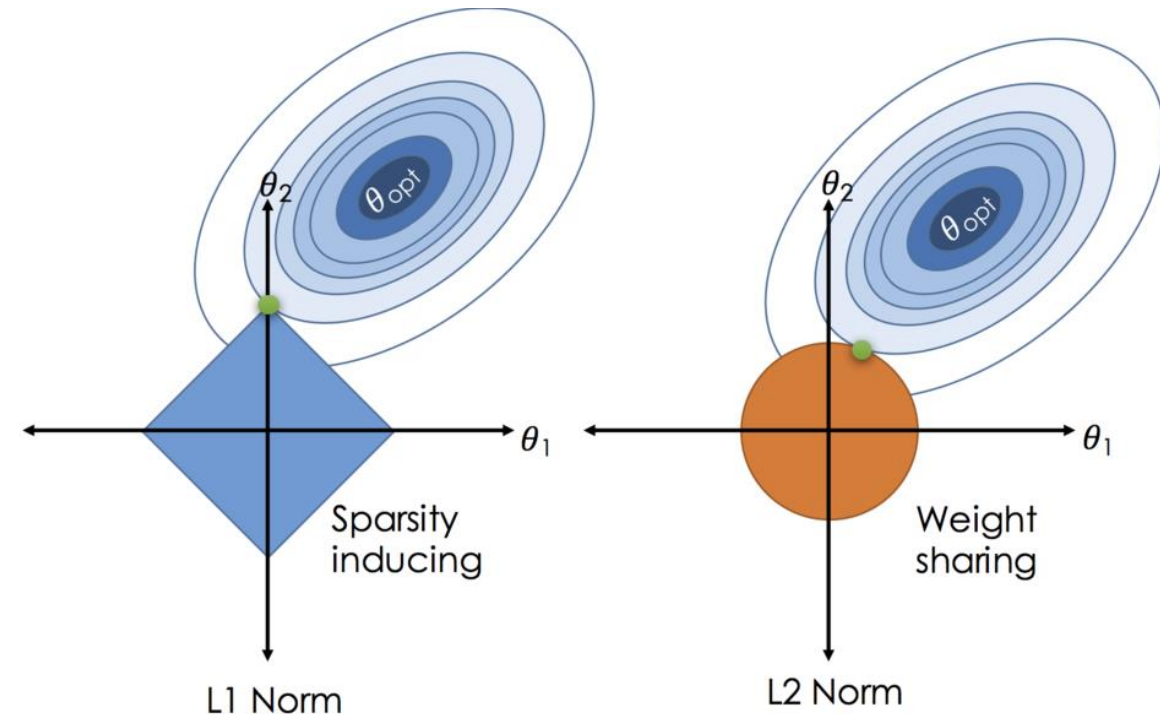
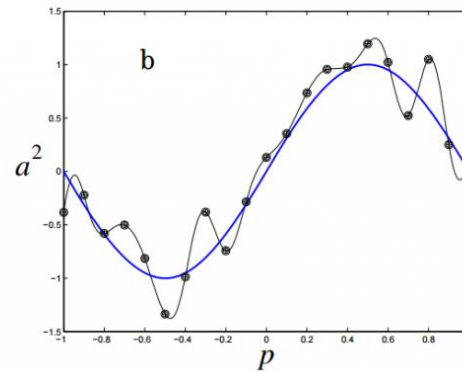
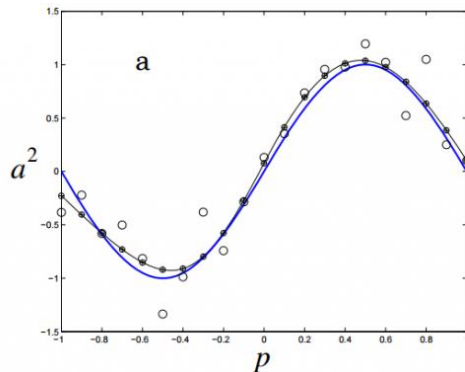


- Worst prediction  
 $L_i = -\log(0) \rightarrow +\infty$
- Best prediction  
 $L_i = -\log(1) = 0$

Slide taken from - Stanford, CS231n course, Lecture 3 presentation  
([http://cs231n.stanford.edu/slides/2017/cs231n\\_2017\\_lecture3.pdf](http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture3.pdf))

# Regularization

- $\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i) + \lambda R(W)$ 
  - $\lambda$  – regularization strength (hyperparameter)
- Commonly used
  - $L_2$  –  $R(W) = \sum W^2$  (Weight decay)
  - $L_1$  –  $R(W) = \sum |W|$



# Gradient Descent Optimization



Stanford, CS231n course, Lecture 3 presentation

# Gradient descent (GD)

- Technique for solving numerical optimization problems
- Algorithm:
  - 1) Initialize parameters ( $\mathbf{W}$  – weights,  $\mathbf{b}$  – biases)
  - 2) Repeat
    - Compute error gradient for each parameter respectively  $\frac{\partial L}{\partial w_{ij}}$   
 $i = 1..m, j = 1..k$
    - Update parameters using gradients  $w_{ij}(n + 1) = w_{ij}(n) - \alpha \frac{\partial L}{\partial w_{ij}}$
- Parameter  $\alpha$  is a **learning rate** and it should be chosen carefully!

# Gradient Descent

- Gradient Descent

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(f(x_i, W), y_i) + \lambda \nabla_W R(W)$$

- (Mini-)Batch GD

$$\nabla_W L(W) = \frac{1}{N_{mb}} \sum_{i=1}^{N_{mb}} \nabla_W L_i(f(x_i, W), y_i) + \lambda \nabla_W R(W)$$

- Stochastic (incremental) GD (SGD)

$$\nabla_W L(W) = \nabla_W L_i(f(x_i, W), y_i) + \lambda \nabla_W R(W)$$

# Gradient Descent Implementations

- Gradient Descent
  - Too slow for large datasets
  - We need to keep intermediate results for each datapoint forward pass
    - Memory problems
  - May be quickly stuck in local minimum
- Mini-batch GD
  - This is used in Deep Learning practice
  - Gradient approximation
  - May skip some local minima
  - $N_{mb}$  is usually 32/64/128/256 (depends on data/model/working memory size)

# Chain rule

- Backpropagation is based on the chain rule:

$$\begin{aligned}h &= f_1(z) \\g &= f_2(h) \\y &= f_3(g) \\y &= f_3\left(f_2\left(f_1(z)\right)\right)\end{aligned}$$

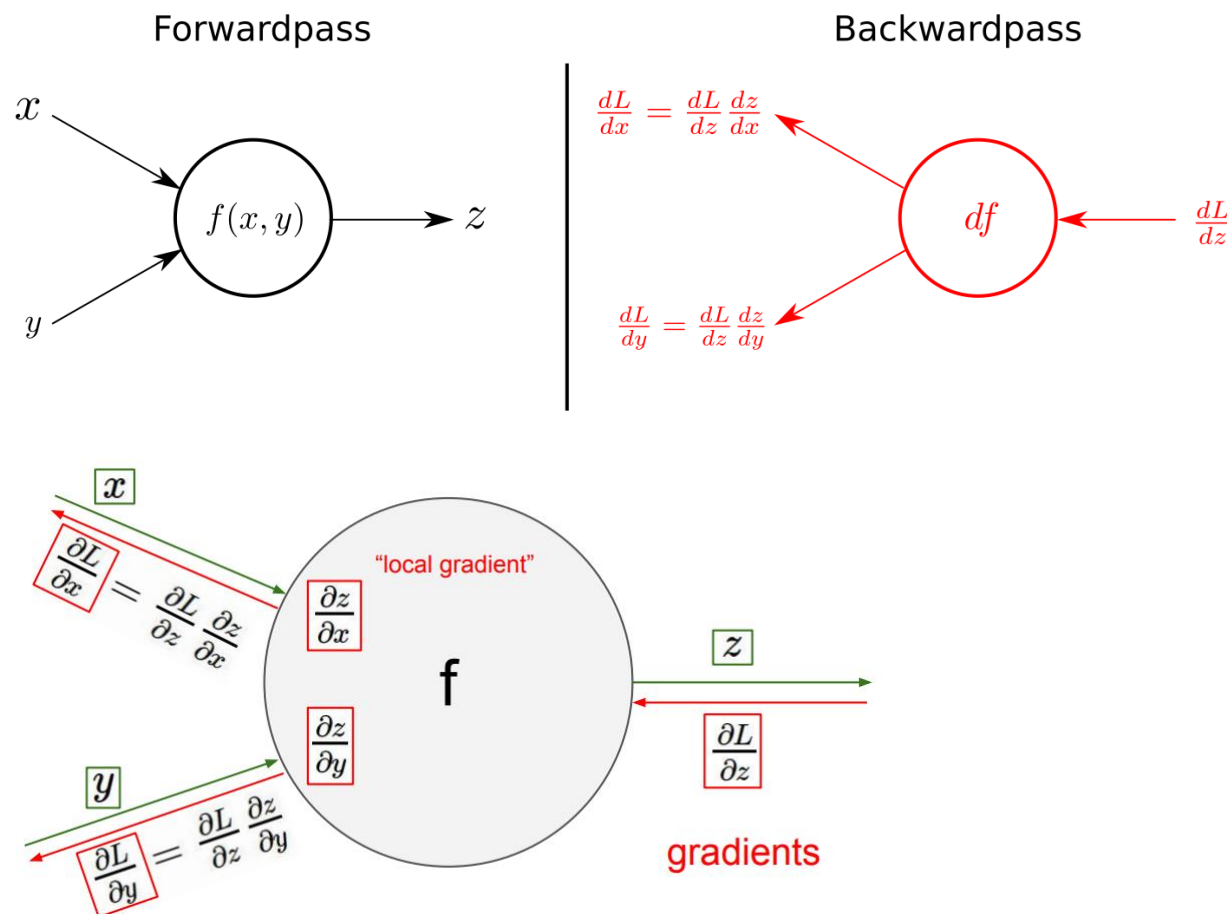
$$\frac{\partial y}{\partial z} = \frac{\partial y}{\partial g} \frac{\partial g}{\partial h} \frac{\partial h}{\partial z}$$



# Backpropagation

## Learning representations by back-propagating errors

David E. Rumelhart\*, Geoffrey E. Hinton† & Ronald J. Williams\* ©1986 Nature Publishing Group



$$E = \frac{1}{2} \sum_c \sum_j (y_{j,c} - d_{j,c})^2 \quad (3)$$

The backward pass starts by computing  $\partial E / \partial y$  for each of the output units. Differentiating equation (3) for a particular case,  $c$ , and suppressing the index  $c$  gives

$$\partial E / \partial y_j = y_j - d_j \quad (4)$$

We can then apply the chain rule to compute  $\partial E / \partial x_j$

$$\partial E / \partial x_j = \partial E / \partial y_j \cdot dy_j / dx_j$$

Differentiating equation (2) to get the value of  $dy_j / dx_j$  and substituting gives

$$\partial E / \partial x_j = \partial E / \partial y_j \cdot y_j (1 - y_j) \quad (5)$$

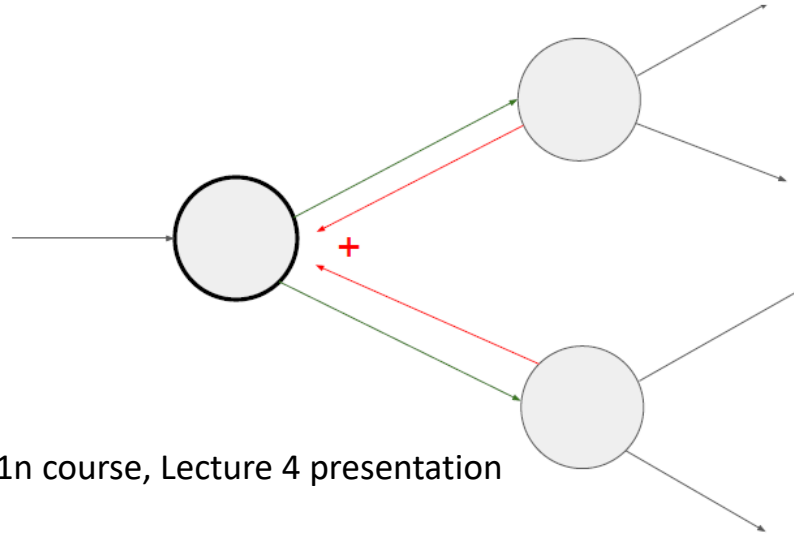
This means that we know how a change in the total input  $x$  to an output unit will affect the error. But this total input is just a linear function of the states of the lower level units and it is also a linear function of the weights on the connections, so it is easy to compute how the error will be affected by changing these states and weights. For a weight  $w_{ji}$ , from  $i$  to  $j$  the derivative is

$$\begin{aligned} \partial E / \partial w_{ji} &= \partial E / \partial x_j \cdot \partial x_j / \partial w_{ji} \\ &= \partial E / \partial x_j \cdot y_i \end{aligned} \quad (6)$$



# Computational graphs

- Neuron with multiple outputs – gradients add



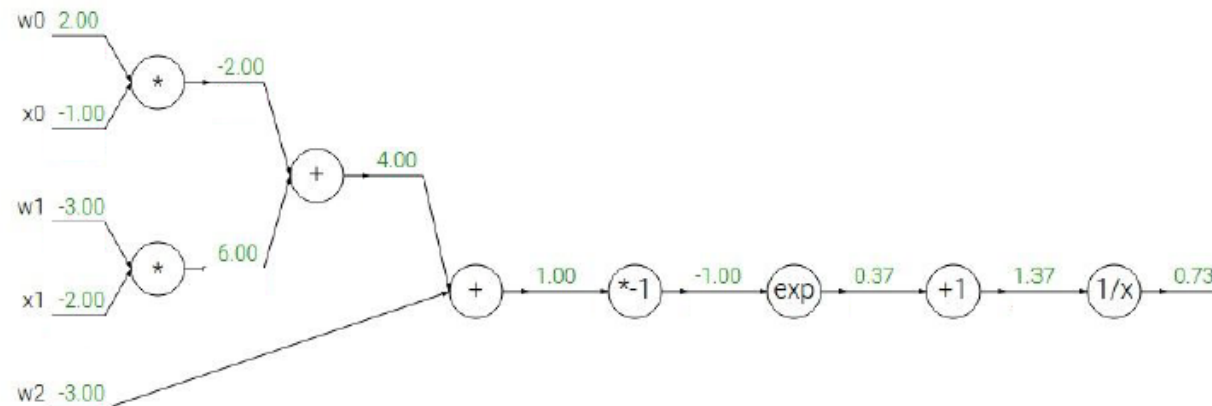
Stanford, CS231n course, Lecture 4 presentation

- Forward and backward passes in these graphs can be vectorized (Jacobian matrix instead of gradient)

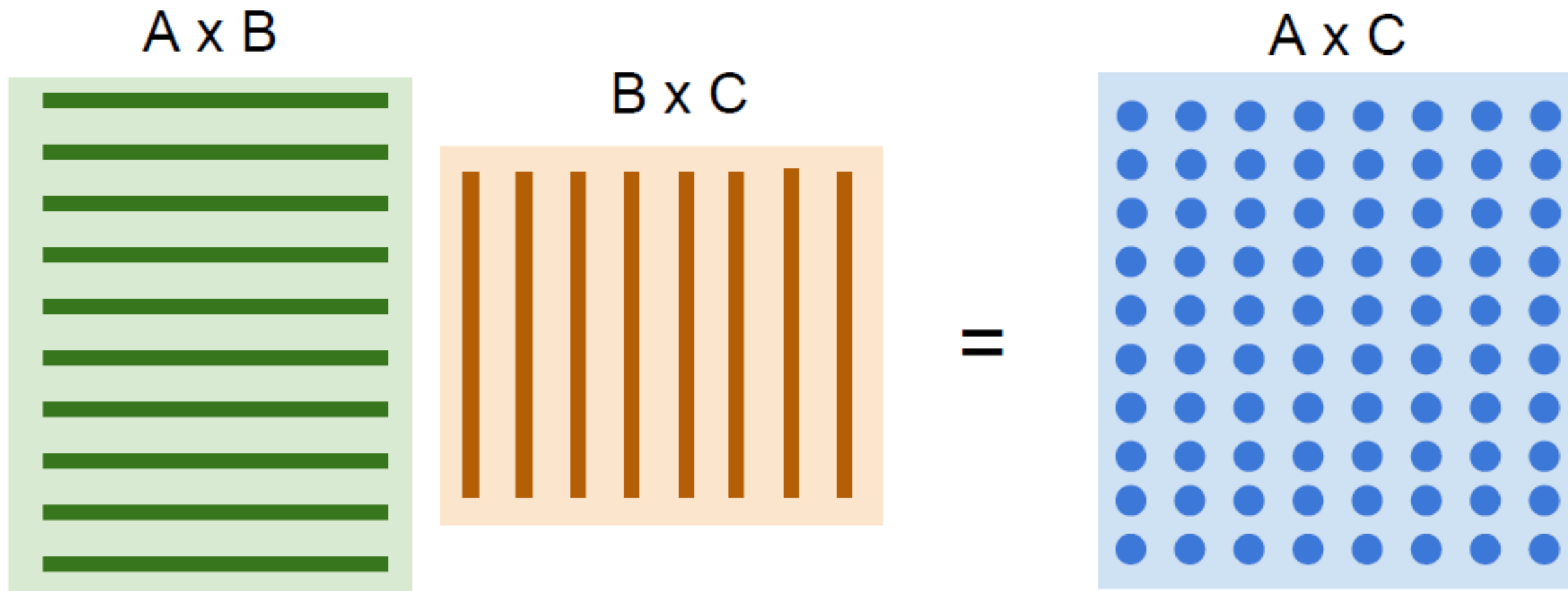
# Computational graphs

- ANN will be very large: impractical to calculate all gradient formulas by hand
- Implementation of DL algorithms comes down to implementing computational graphs and calculating their partial derivatives
  - This is called backpropagation

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



Matrix multiplications can be parallelized  
 $\Rightarrow$  GPUs



# Common practices

- Usually we start with certain amount of labeled samples
  - $(x_i, y_i); i = 1, \dots, n$
- The dataset is divided into three partitions
  - Training set – used in training
  - Validation set – used to tune hyperparameters
  - Test set – used to measure quality of the network
- Neural network are “data” hungry
- Randomization
  - Present samples in random order

# Some terminology

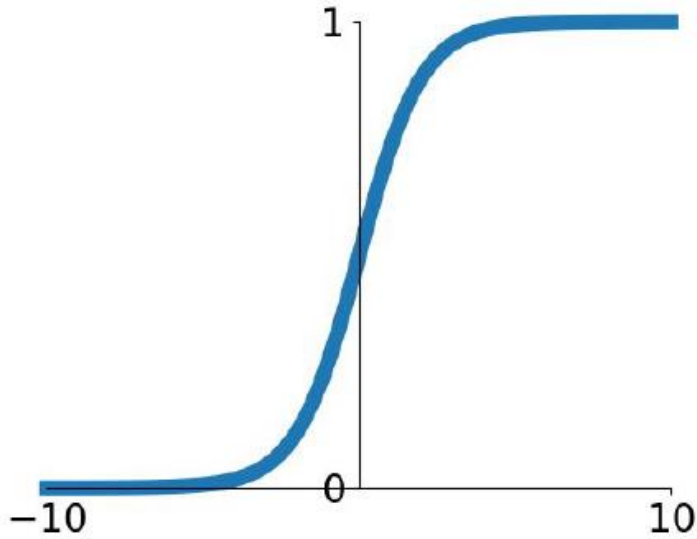
- **Forward pass:** Calculate the network outputs  $y_i$  based on the inputs  $x_i$
- **Backward pass:** Calculate gradients  $\frac{\partial L}{\partial w_{ij}}$ , by using the chain rule and intermediate calculations from the forward pass
- **Backpropagation:** Update network weights  $w_{ij}$  by using the calculated gradients  $\frac{\partial L}{\partial w_{ij}}$
- **Mini-batch:** examples whose gradients are averaged before backpropagation
- **Epoch:** The ENTIRE dataset is passed ONCE forwards and backwards
  - A complete training goes through multiple epoch

You can train an ANN now





# Activation Functions - Sigmoid



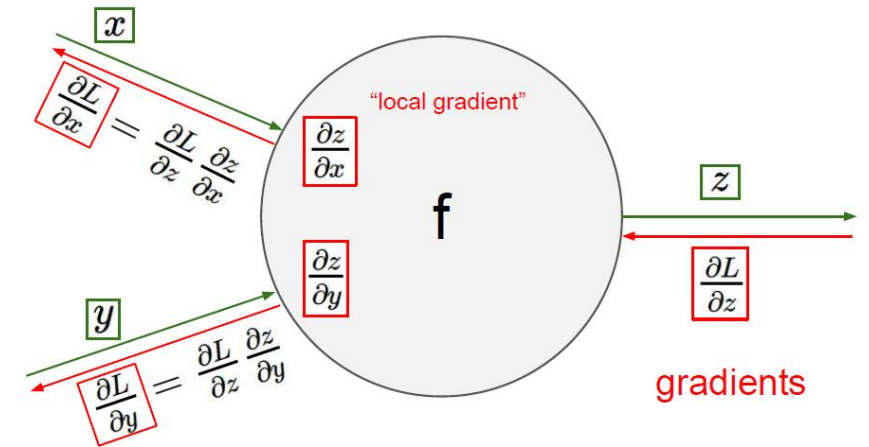
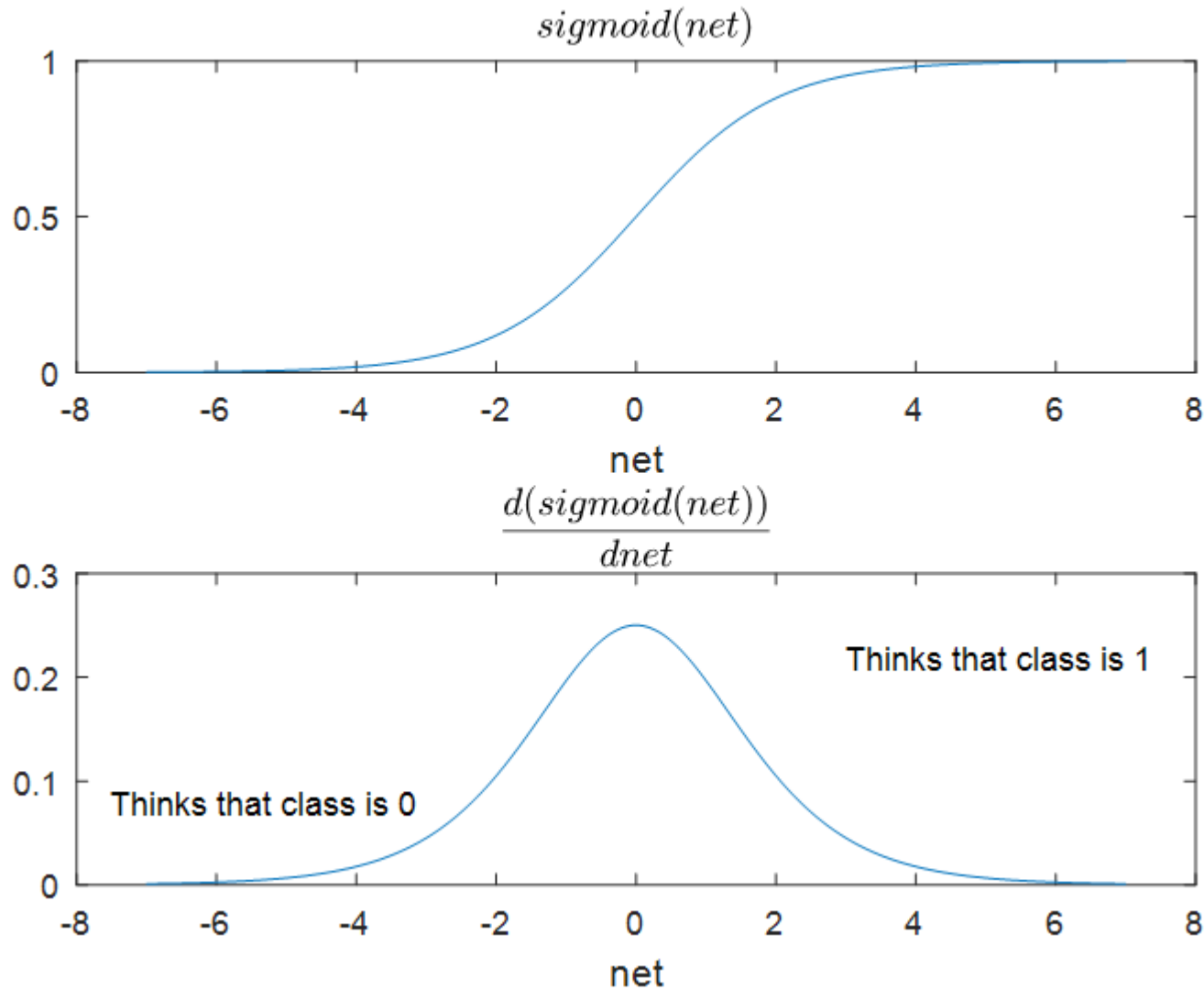
- Saturation “kills” the gradients **X**

Stanford, CS231n course, Lecture 6 presentation

$$\sigma(net) = \frac{1}{1 + e^{-net}}$$

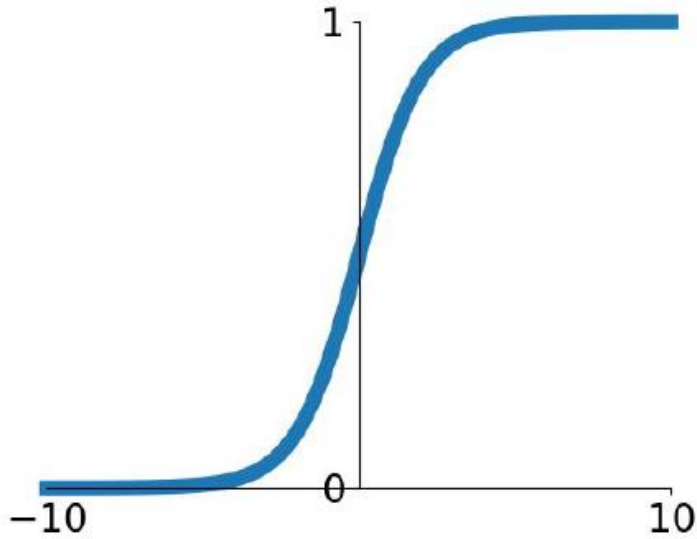


# Activation Functions - Sigmoid



Stanford, CS231n course, Lecture 4 presentation

# Activation Functions - Sigmoid



Stanford, CS231n course, Lecture 6 presentation

- Saturation “kills” the gradients **X**
- Sigmoid always outputs a positive number **X**

$$\sigma(net) = \frac{1}{1 + e^{-net}}$$

# Activation Functions - Sigmoid

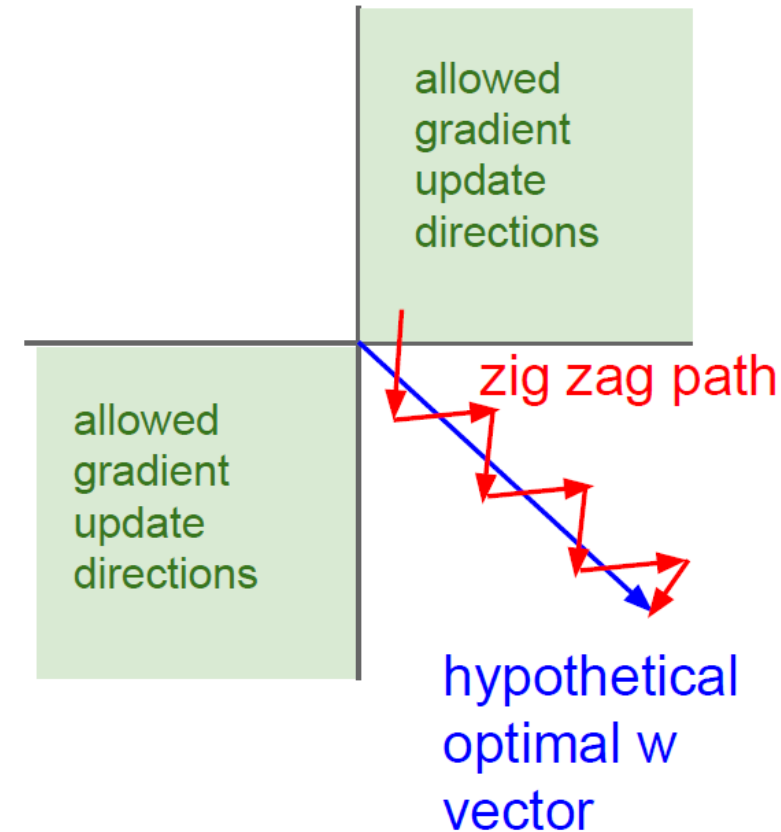
- $\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial w_i}$ 
  - $z = \sum_i x_i w_i + b$  (pre-activation)

1. When inputs  $x_i$  are always positive

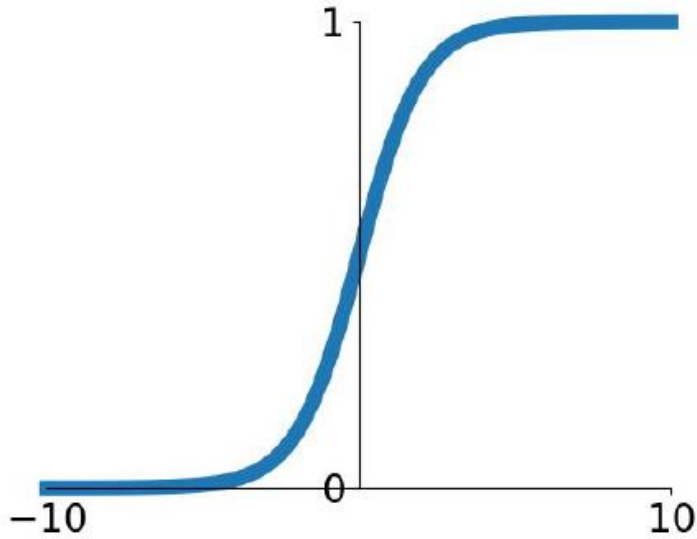
- $\Rightarrow \frac{\partial z}{\partial w_i} = x_i$  is always positive

2. Each  $w_i$  uses the same gradient  $\frac{\partial L}{\partial z}$

$\Rightarrow$  Gradient update rule is positive/negative for all  $w_i$



# Activation Functions - Sigmoid

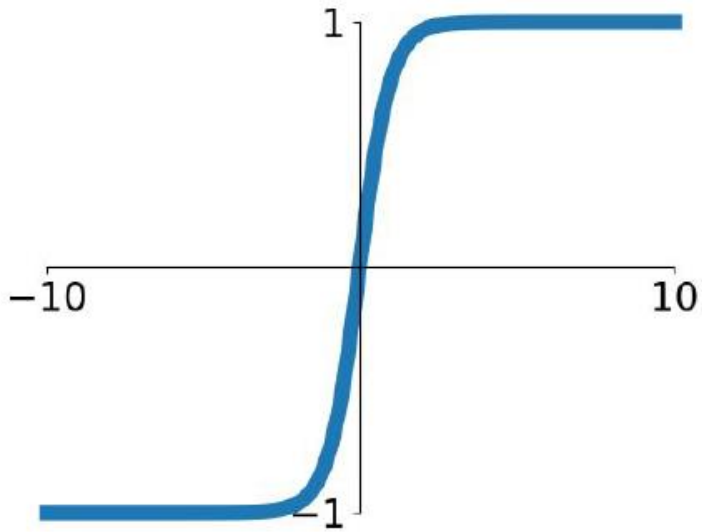


Stanford, CS231n course, Lecture 6 presentation

- Saturation “kills” the gradients **X**
- Sigmoid always outputs a positive number **X**
- $\exp( )$  is a bit computationally expensive **X**
  - Not a big problem

$$\sigma(net) = \frac{1}{1 + e^{-net}}$$

# Activation Functions – Hyperbolic Tangent

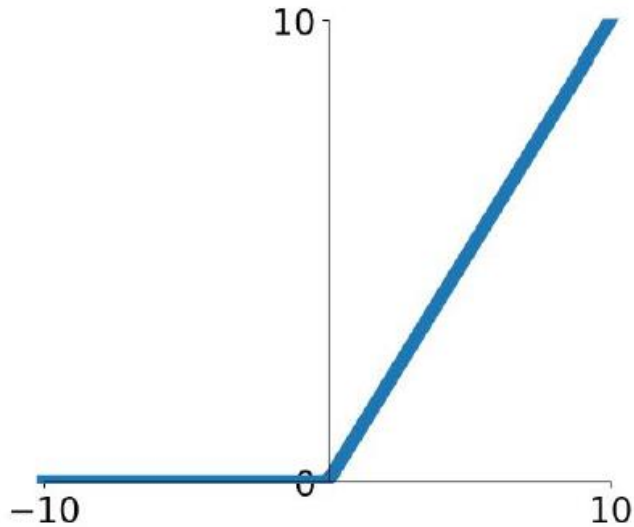


Stanford, CS231n course, Lecture 6 presentation

$$\sigma(net) = \tanh(net)$$

- Zero centered ✓
- Saturation “kills” the gradients ✗
- A bit computationally expensive ✗  
(not a big problem)

# Activation Functions – ReLU

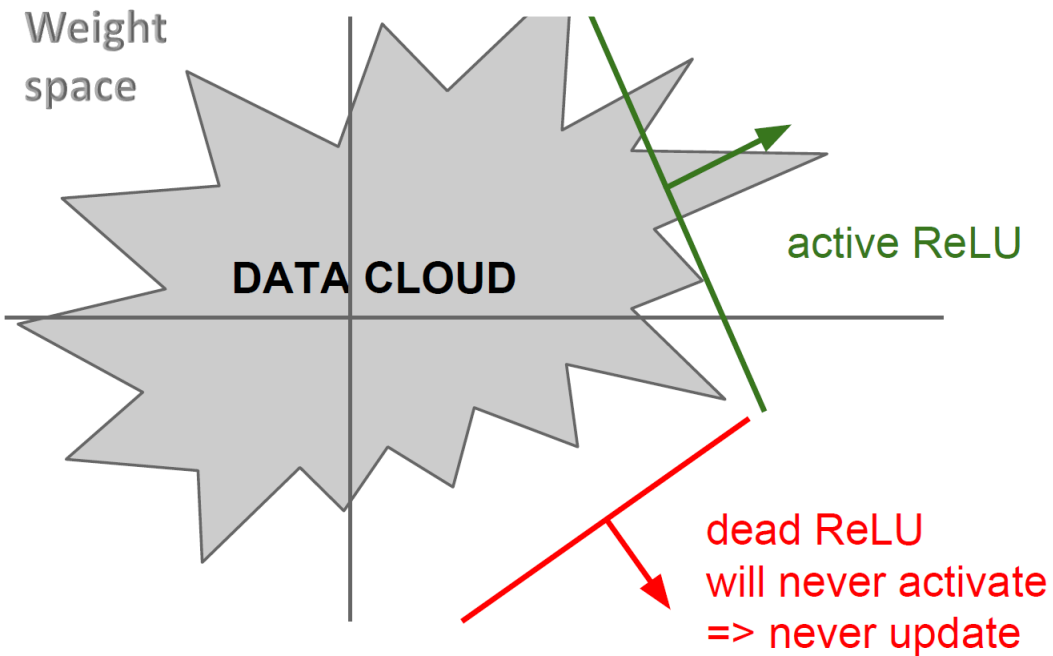


Stanford, CS231n course, Lecture 6 presentation

- Doesn't saturate in + region ✓
- Computationally efficient ✓
- Converges much faster than sigmoid/tanh in practice ✓
- Output not zero-centered ✗
- Gradient is 0 in – region ✗ (saturation)

$$\sigma(net) = \max(0, net)$$

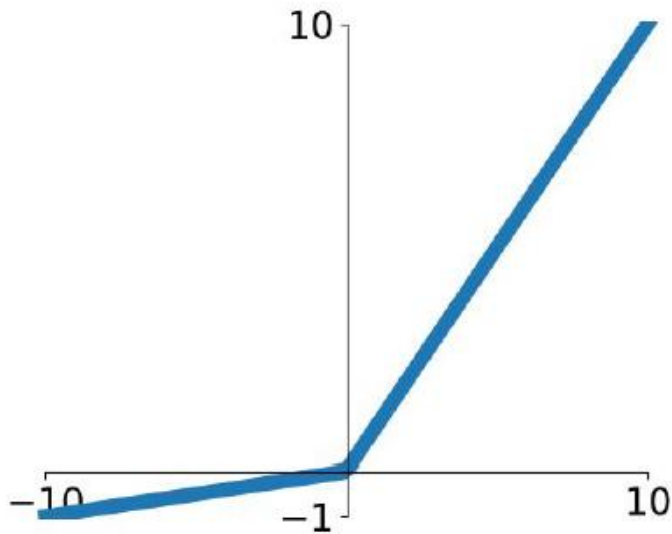
# Activation Functions – ReLU



Stanford, CS231n course, Lecture 6 presentation

- If datapoints don't activate the ReLU, it won't update weights in a BP step
  - ReLUs with bad initialized weights will stay dead
- Learning rate too high
  - ReLU could fall into the dead ReLU region
- Initializing ReLU biases with small positive numbers often helps

# Activation Functions – Leaky ReLU



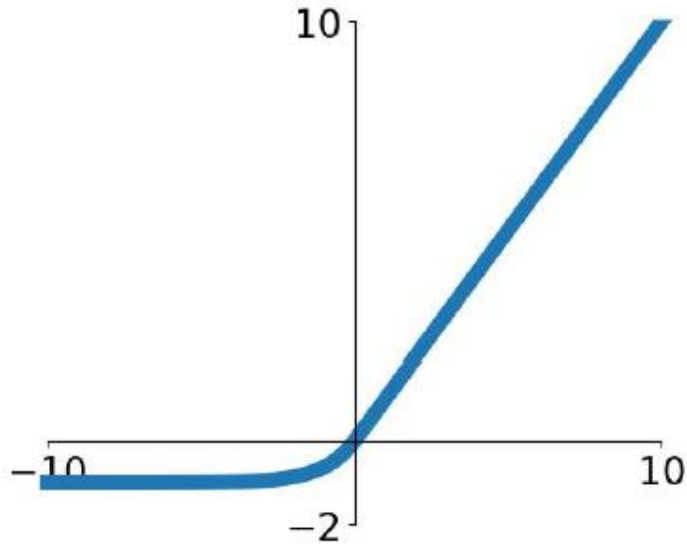
Stanford, CS231n course, Lecture 6 presentation

$$\sigma(net) = \max(0.01net, net)$$

- Doesn't saturate in + region ✓
- Computationally efficient ✓
- Converges much faster than sigmoid/tanh in practice ✓
- Doesn't saturate in - region ✓  
(will not die)
- Slope in the – region could be parametrized  
(learn it with backprop)



# Activation Functions – ELU



Stanford, CS231n course, Lecture 6 presentation

- Similar to leaky ReLU ✓
- Saturation in the – region adds some robustness to noise compared to leaky ReLU ✓
- $\exp( )$  is a bit computationally expensive ✗

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

# Activation Functions – Maxout

$$\sigma(net) = (w_1^T x + b_1, w_2^T x + b_2)$$

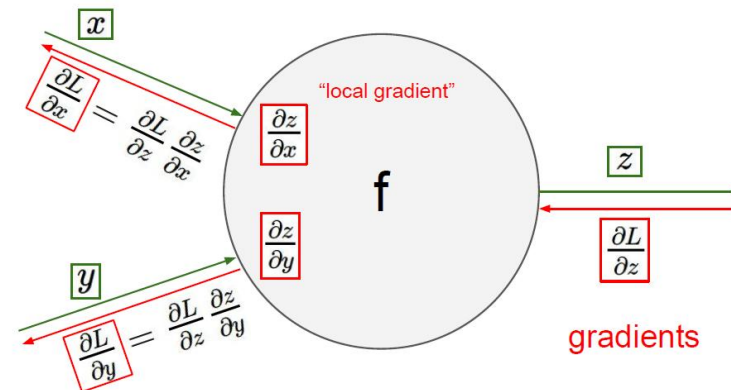
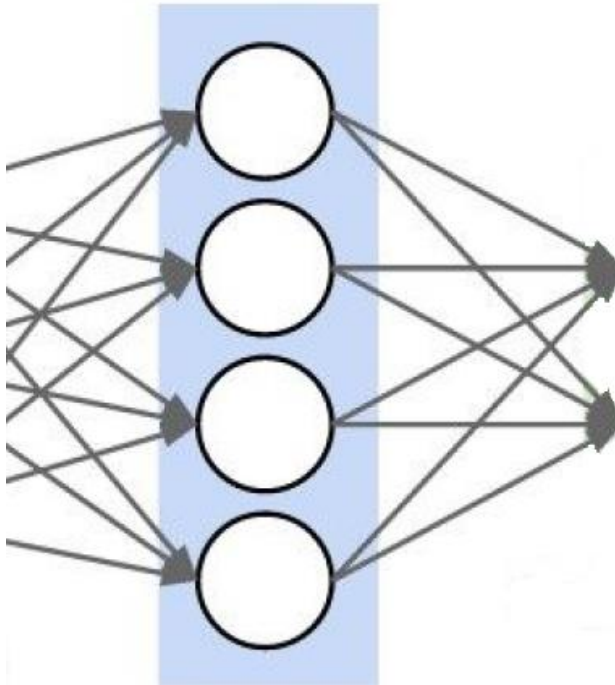
- Generalizes ReLU and Leaky ReLU ✓
- Doesn't saturate ✓
- Doesn't die ✓
- Twice as many parameters as ReLU ✗

# Activation Functions – Practical advice

- Use **ReLU**. Be careful with your learning rates
- Try out **Leaky ReLU / Maxout / ELU**
- Try out **tanh** but don't expect much
- **Don't use sigmoid**

# Weight Initialization – $W=0$

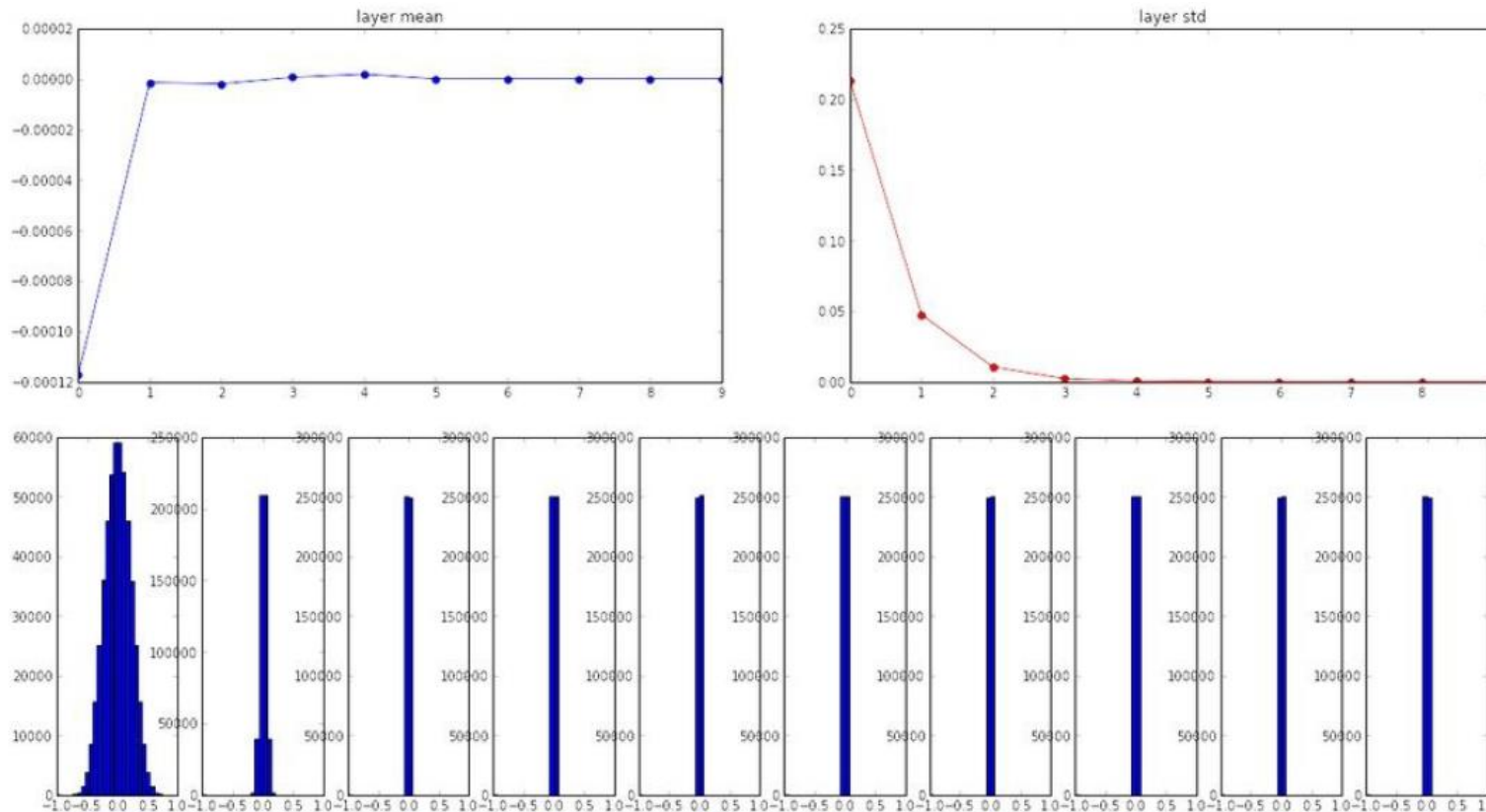
- What will happen if we initialize all weights to 0?



Stanford, CS231n course, Lecture 4 presentation

# Weight Initialization – Small Random Values

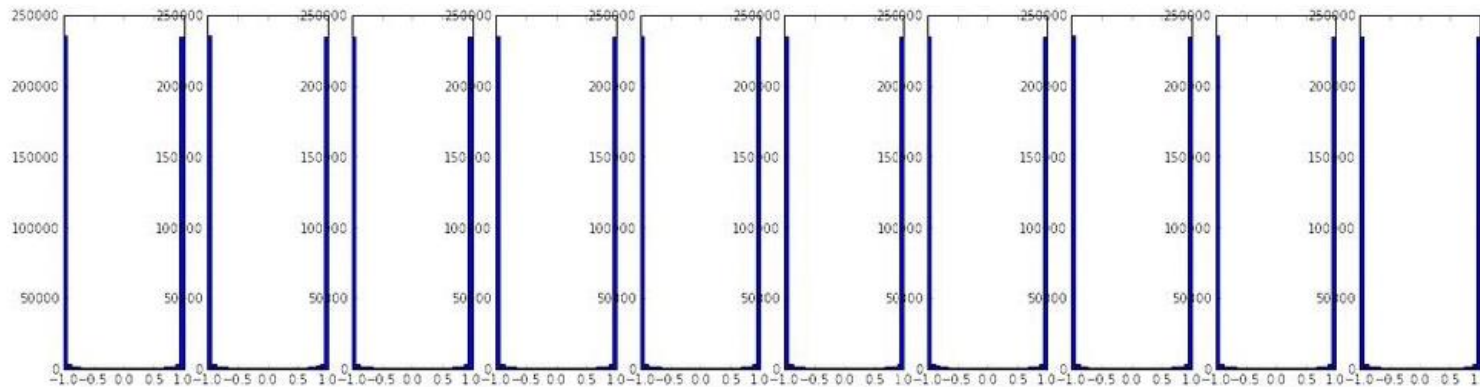
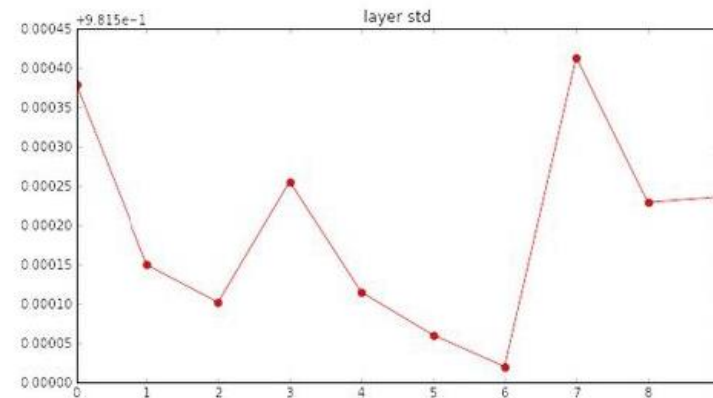
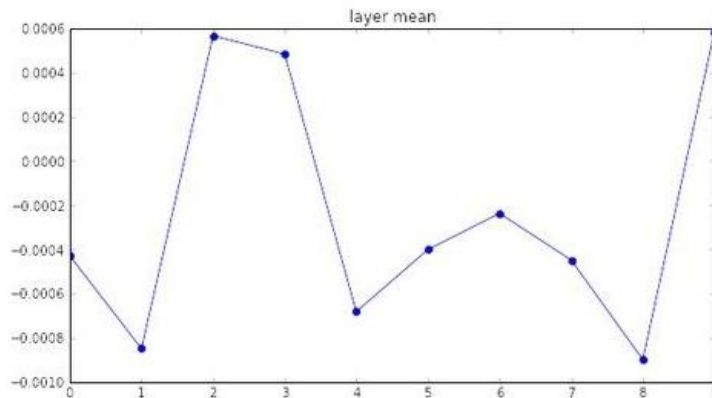
Tanh activations



- Deeper into the network activations become zero
- During backpropagation each layer is multiplied by  $w$ 
  - Slow learning

# Weight Initialization – Big Random Values

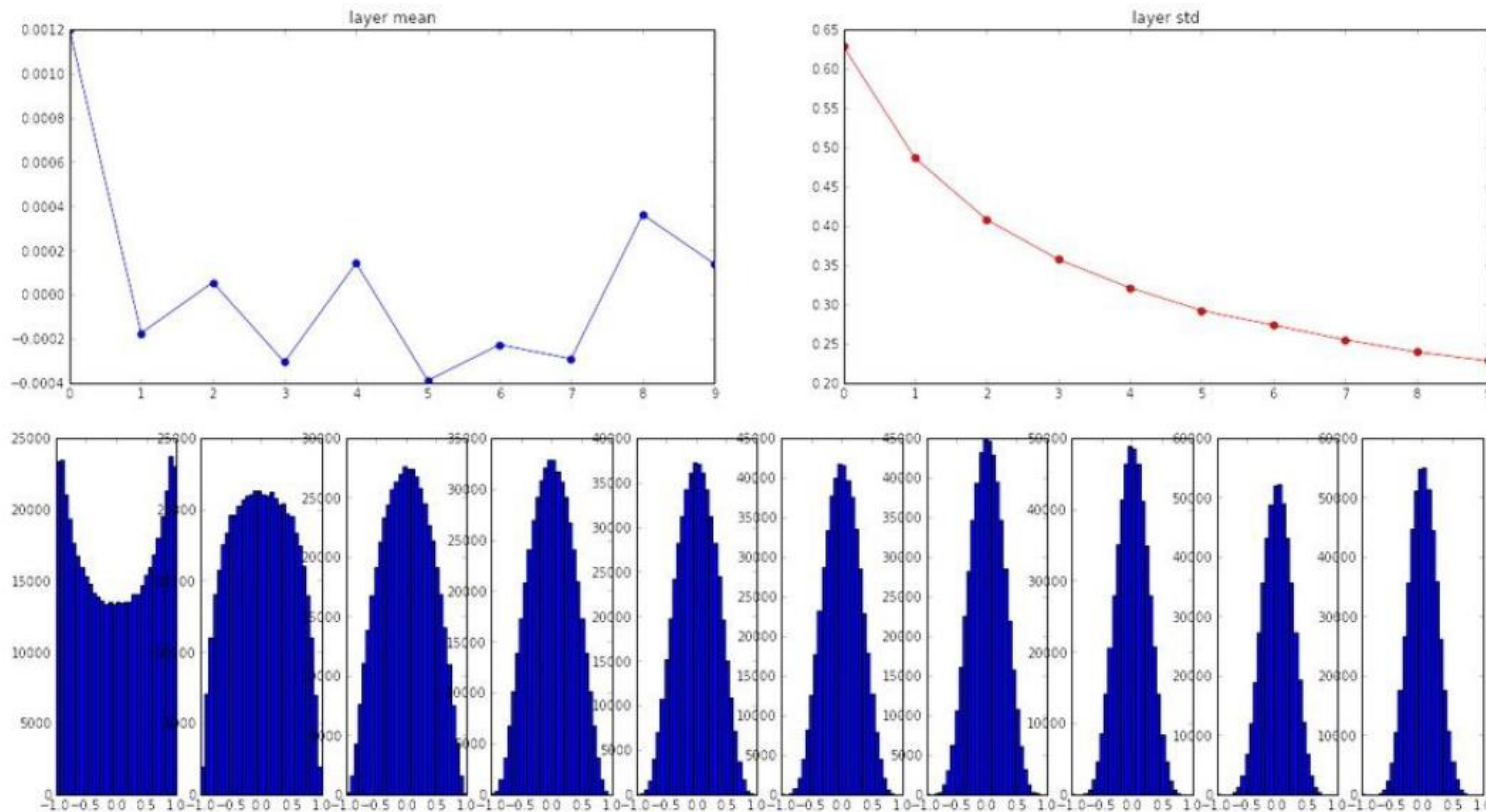
Tanh activations



- Most of the activations go into saturation.
- Gradients will be small because of the saturations
  - Slow learning

# Weight Initialization – Xavier

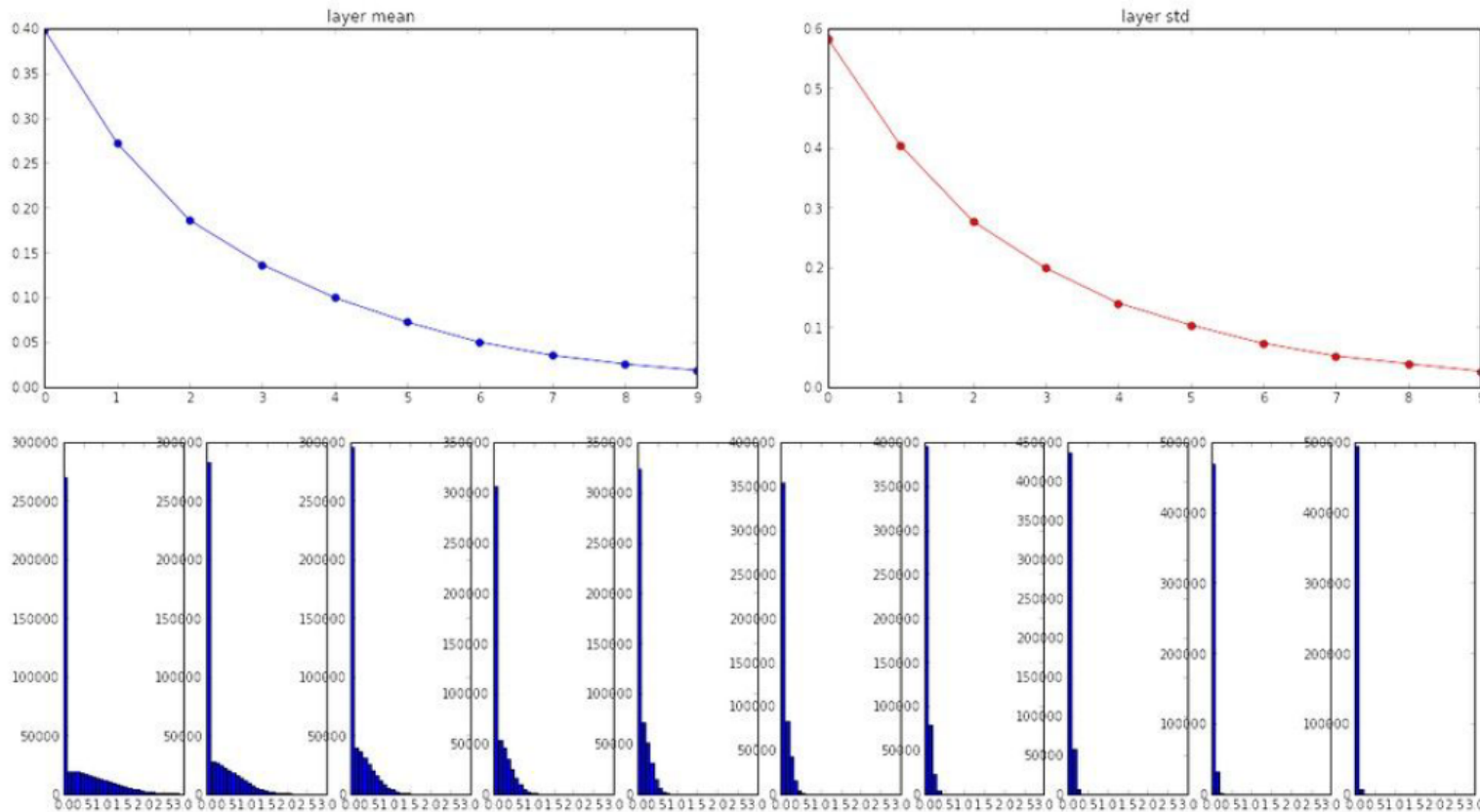
Tanh activations



- Glorot et al., 2010
  - Wants the same variance at input & output of a layer
- Layer initialization:
  - Typically Gaussian or Uniform
$$\mu(W)=0$$
$$\sigma(W) = \frac{1}{\sqrt{n_{in} + n_{out}}}$$
  - $n_{in}$  – number of units in the previous layer
  - $n_{out}$  – number of units in the next layer

# Weight Initialization – Xavier

ReLU activations

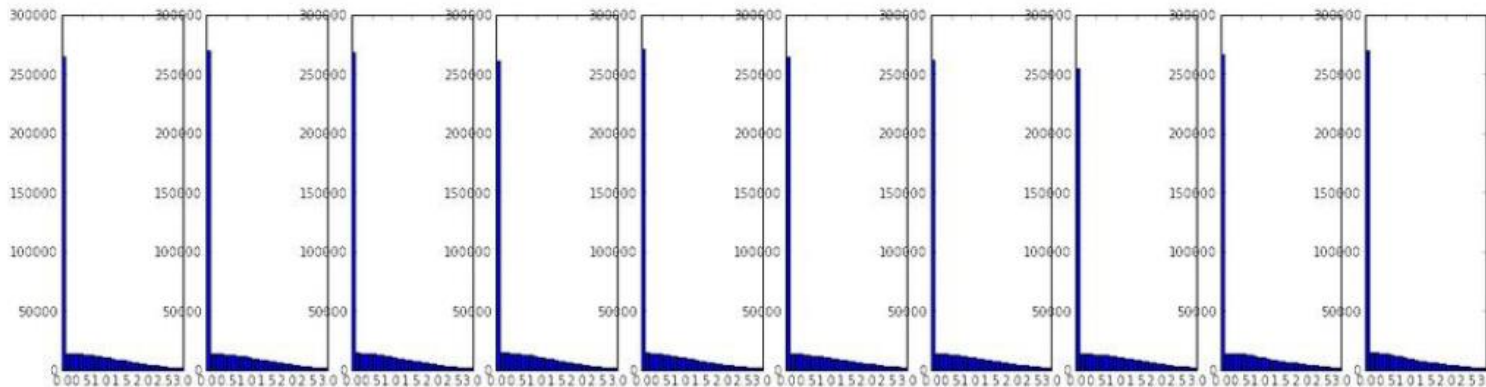
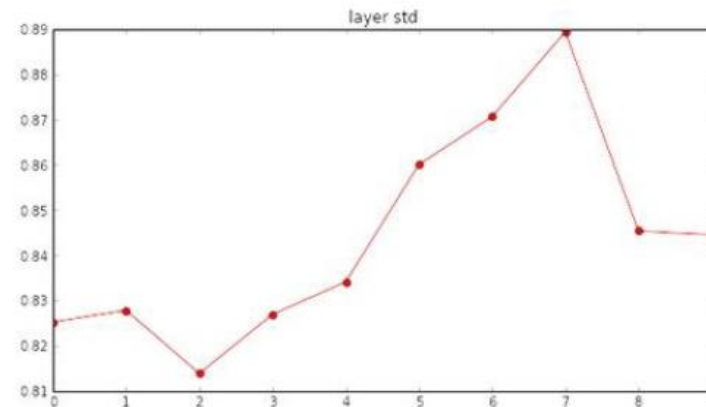
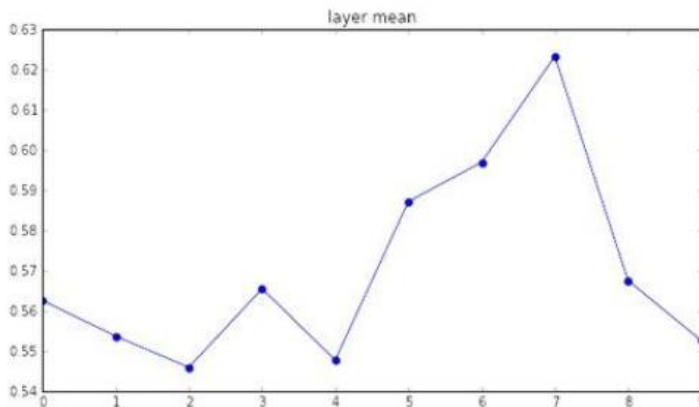


- Doesn't work well on ReLU



# Weight Initialization – Recent Recommendation

ReLU activations



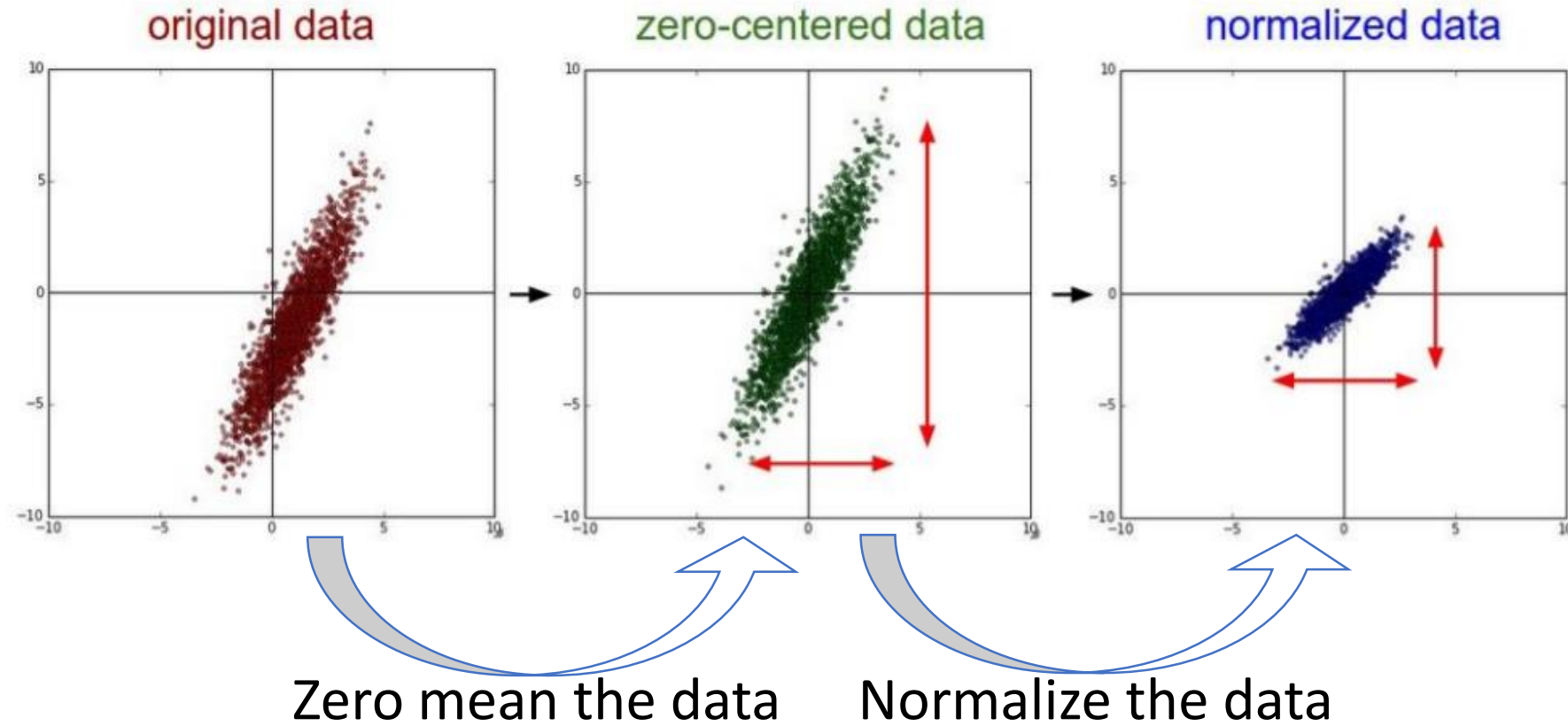
- He et al., 2015
- Layer initialization:
  - Typically Gaussian or Uniform
$$\mu(W)=0$$
$$\sigma(W) = \frac{2}{n_{in}}$$
  - $n_{in}$  – number of units in the previous layer

# Weight Initialization

- Having a good initialization scheme is important in practice
- Some networks couldn't be trained at all without a good initialization
- Still an open area of research

# Data preprocessing

Stanford, CS231n course, Lecture 6 presentation



All features are in the same range  $\Rightarrow$  they all contribute equally

# In practice

Training phase:

- Subtract the mean for each image
- Normalize with standard deviation

Testing phase:

- Subtract empirical mean (obtained from training data)

Only helps for the first layer!

# Batch normalization

$$\hat{x}^{(k)} = \frac{x^{(k)} - E\{x^{(k)}\}}{\sqrt{\text{var}\{x^{(k)}\}}}$$

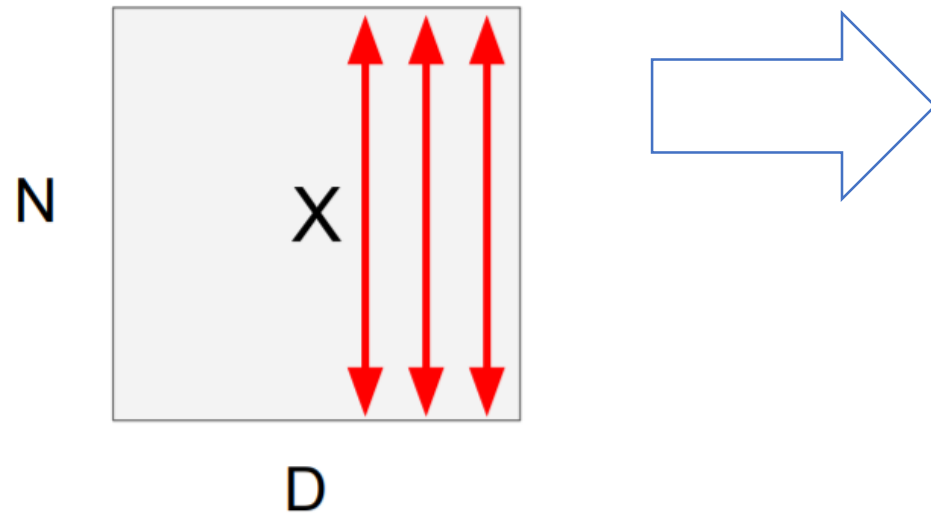
$$y^{(k)} = \gamma^{(k)} x^{(k)} + \beta^{(k)}$$

- $\beta$  &  $\gamma$  are hyperparameters that can change during training
- Network can learn to undo the normalization of not suitable

$$\gamma^{(k)} = \text{var}\{x^{(k)}\}$$

$$\beta^{(k)} = E\{x^{(k)}\}$$

# Batch normalization

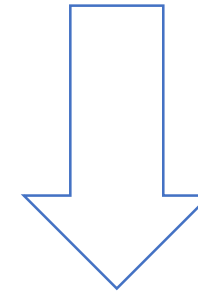


$x^{(k)}$  – current batch

$N$  – training examples in the batch

$D$  – dimension of each batch

Compute empirical mean and variance independently for each dimension



$\mathbf{x}: N \times D$

Normalize



$\boldsymbol{\mu}, \boldsymbol{\sigma}: 1 \times D$

$\boldsymbol{\gamma}, \boldsymbol{\beta}: 1 \times D$

$\mathbf{y} = \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta}$

# Batch normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

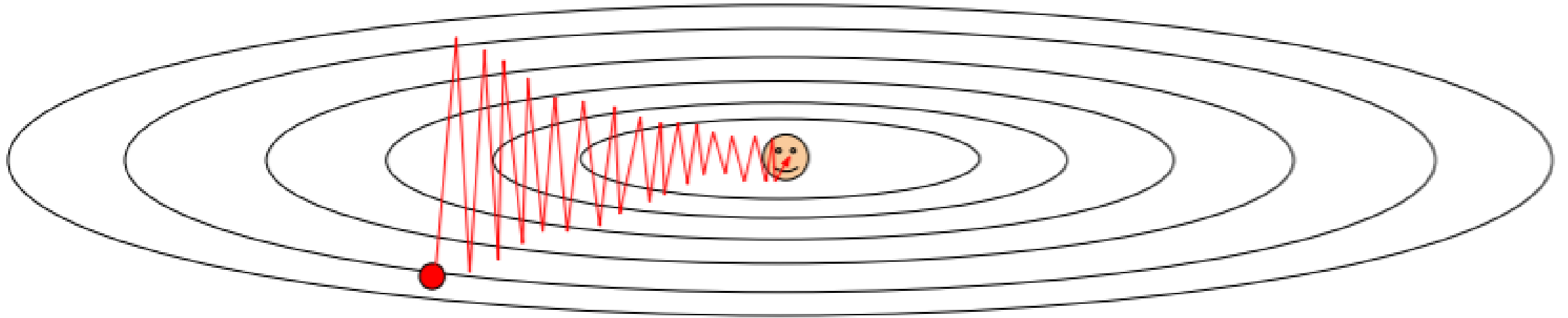
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

- More robust to:
  - Higher learning rates
  - Different kinds of weight init.
- Improves gradient flow
  - Easier to train
- Some kind of regularization
- It is not the same during training and testing

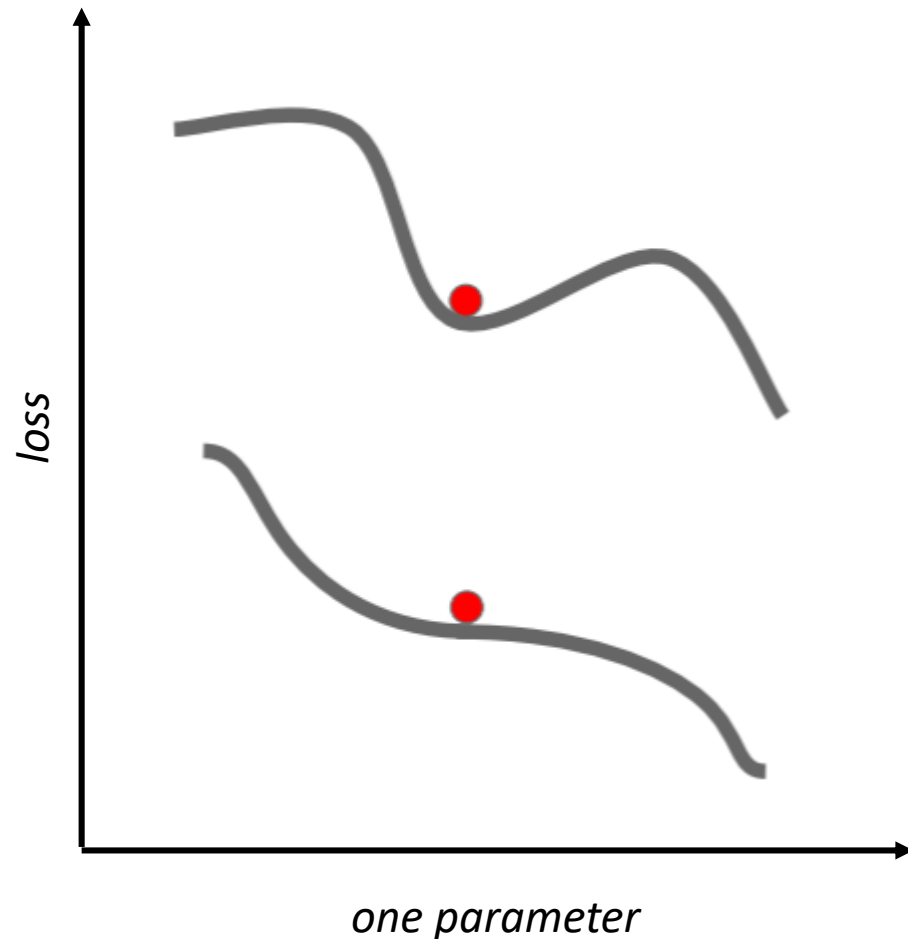
# 1<sup>o</sup> problem with SGD



- Different sensitivity to different dimensions
  - Slow progress along shallow dimension
  - Zigzag behavior along steep dimension

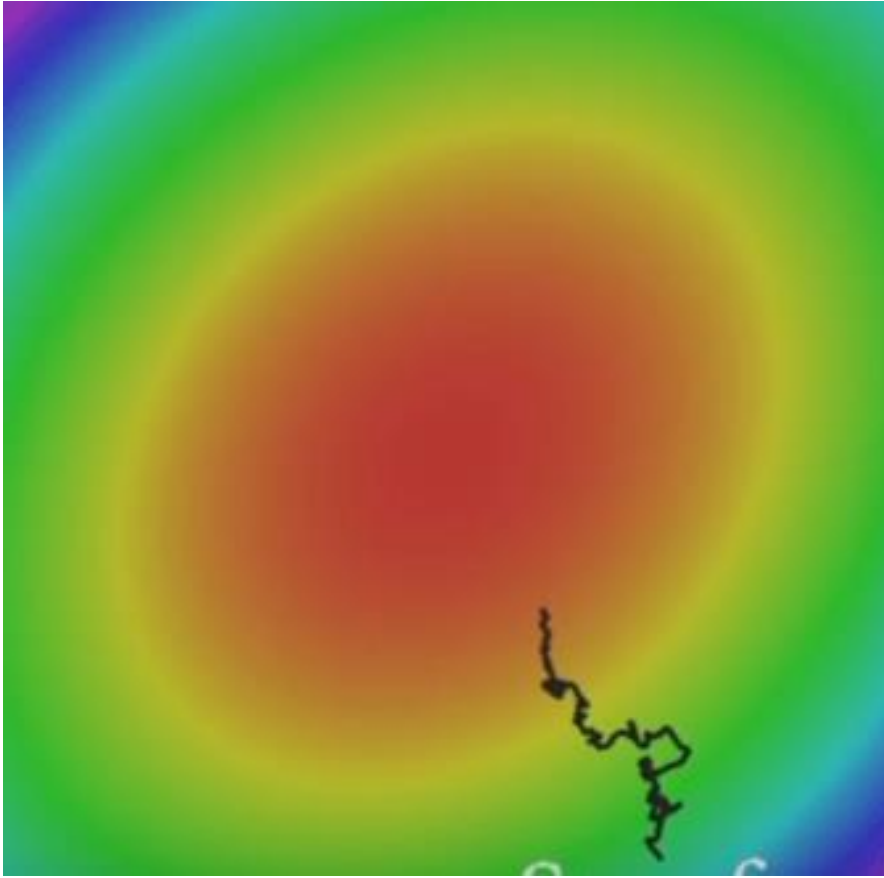


## 2° problem with SGD



- Can get stuck in local minima
- Local minima is not a big problem for high dimensions
- Saddle point is more common problem
  - Gradient is very small

# 3° Problem with SGD



Stanford, CS231n course, Lecture 7 presentation

- We usually use minibatches
  - We are only getting the noisy estimation of the gradient in the current minibatch

# Momentum

SGD

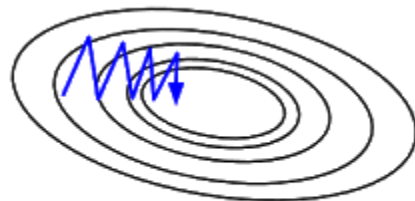
$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

SGD + Momentum

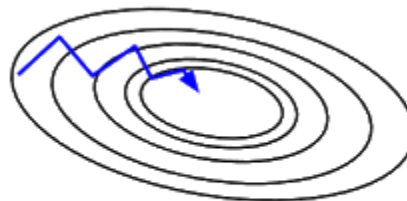
$$v_{t+1} = \rho v_t + \alpha \nabla f(x_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1}$$

- SGD – always stepping in the direction of the gradient
- SGD + Momentum – stepping in the direction of velocity

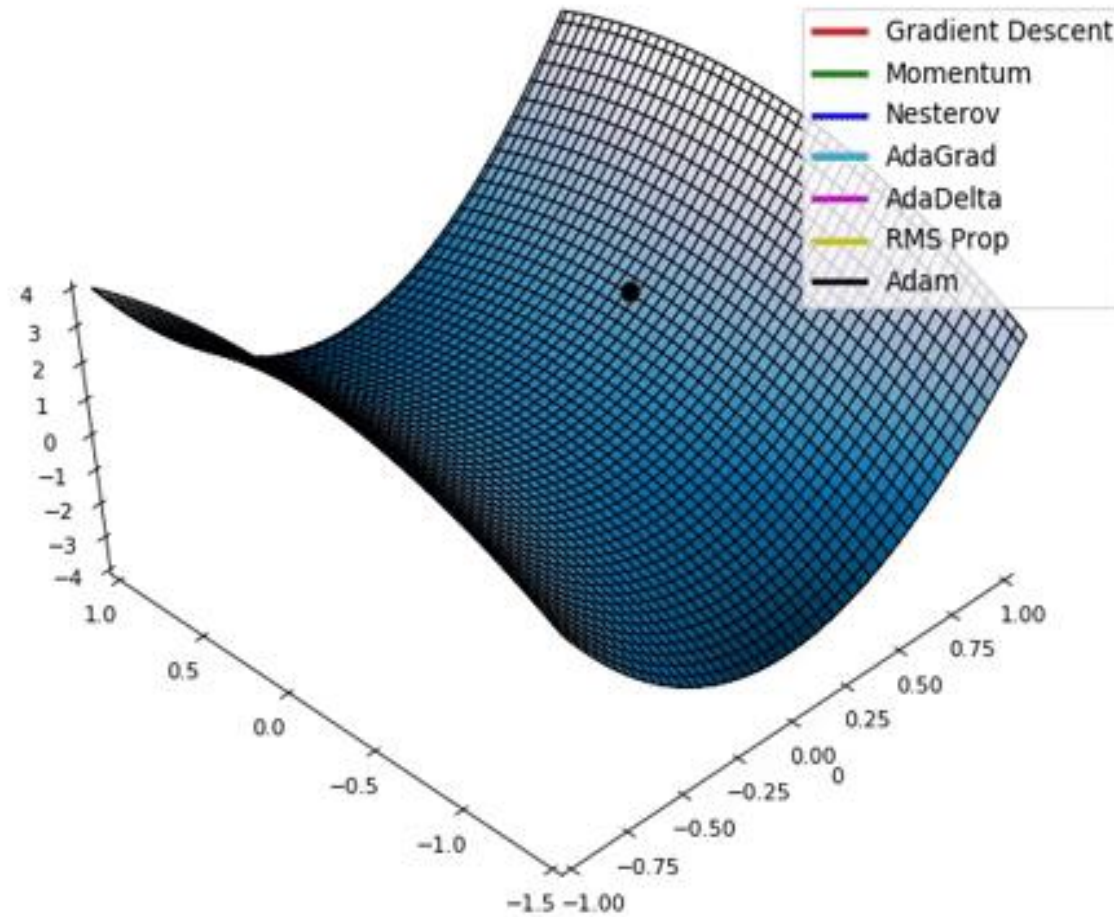


Stochastic Gradient Descent  
(SGD) without Momentum



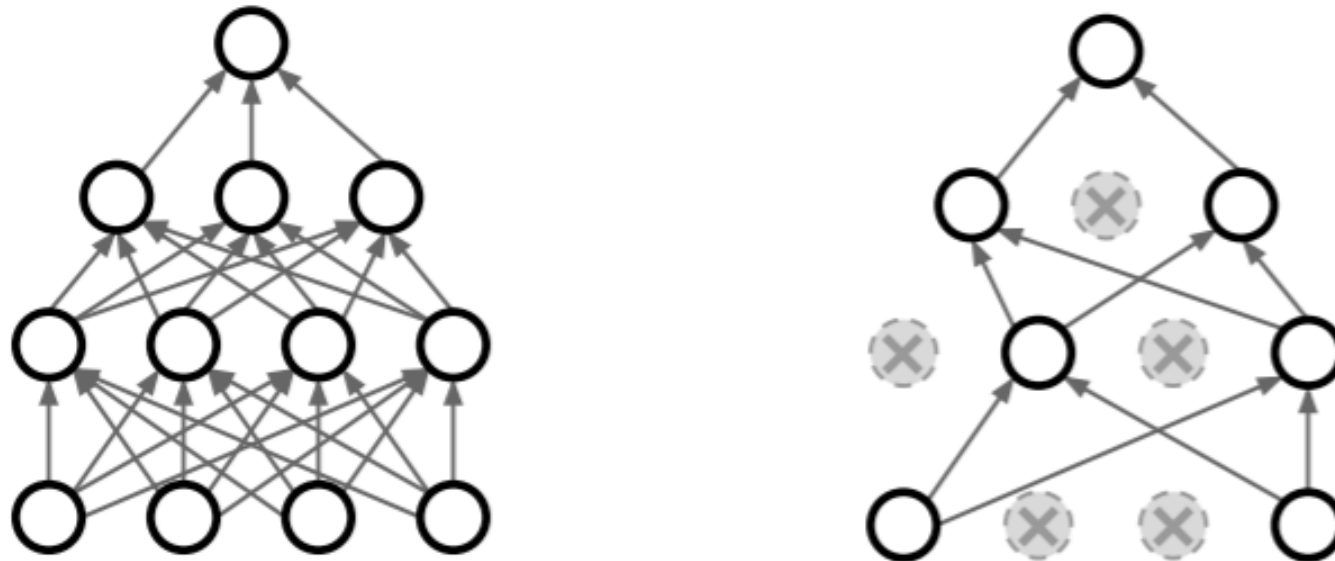
Stochastic Gradient Descent  
(SGD) with Momentum

# More sophisticated options



# Regularization - dropout

- In each pass through the NN we randomly set some of the neurons to zero, with probability  $p$  one layer at the time
- Dropping probability  $p$  is a hyperparameter
- Each dropout gives us a different subset of the NN
  - Dropout is like learning whole ensemble of networks at the same time



# Dropout – why is it good?

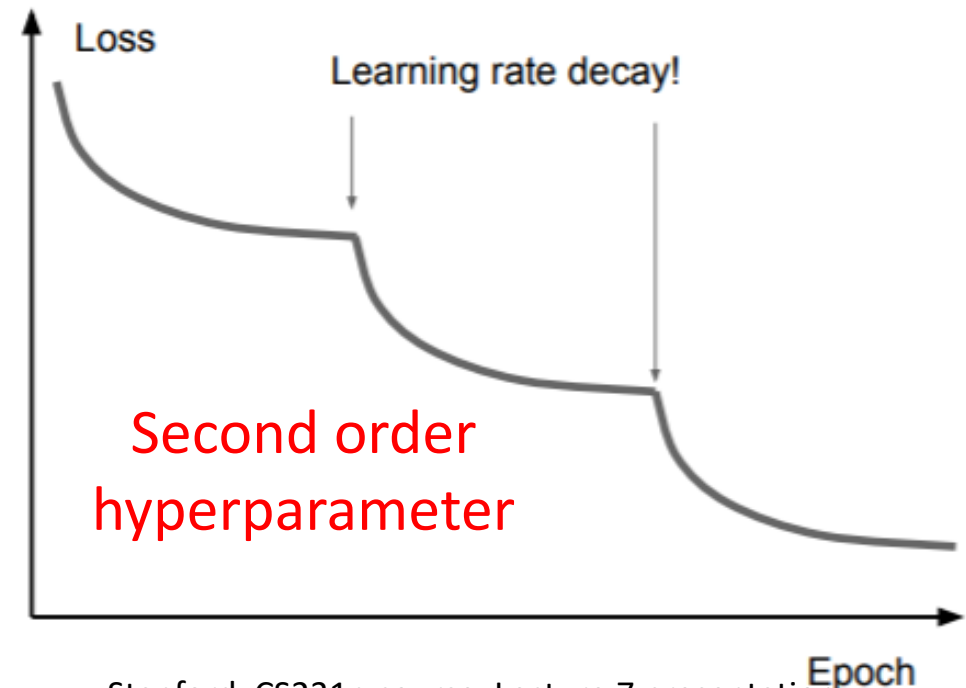
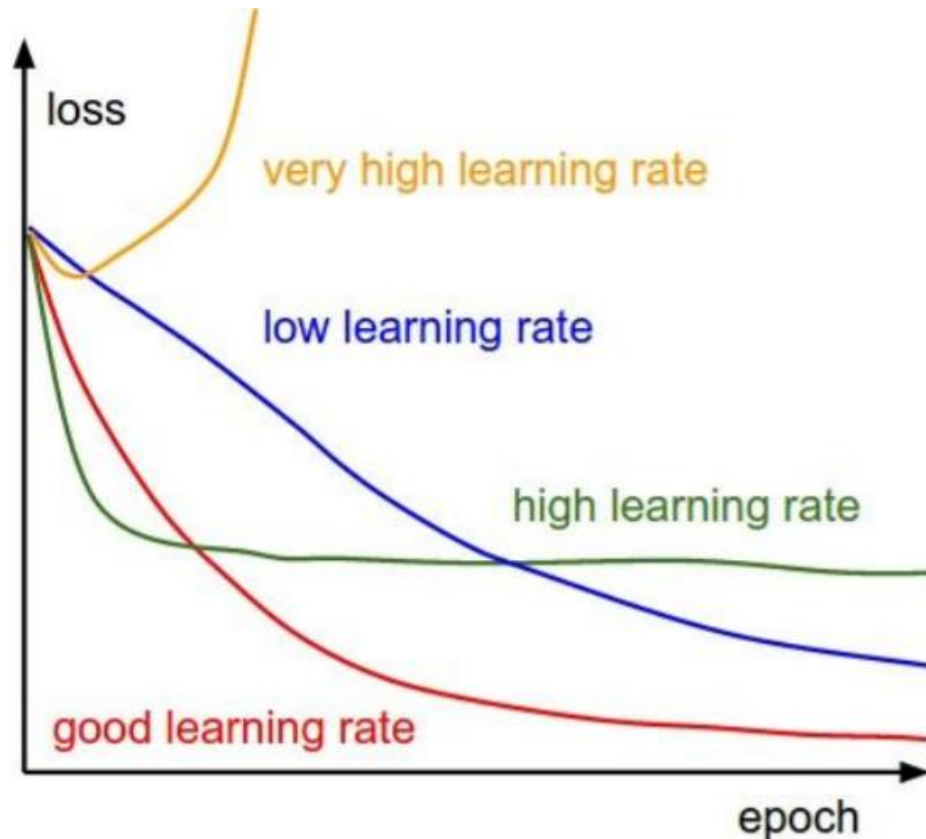
- Helps to prevent co-adaptation of features
- The algorithm doesn't depend too much on one feature



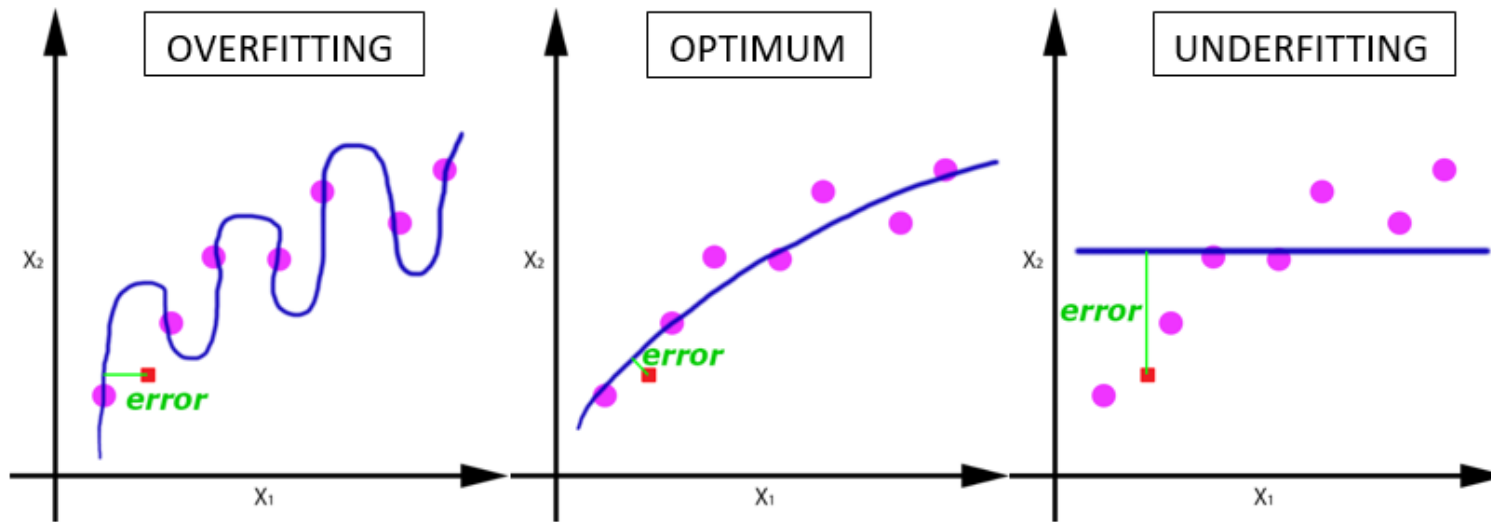
# Learning rate adjusting

- SGD, SGD + Momentum, AdaGrad, RMSProp and Adam have learning rate as hyperparameter

⇒ Learning rate decay over time: start with bigger value, and decrease it over time

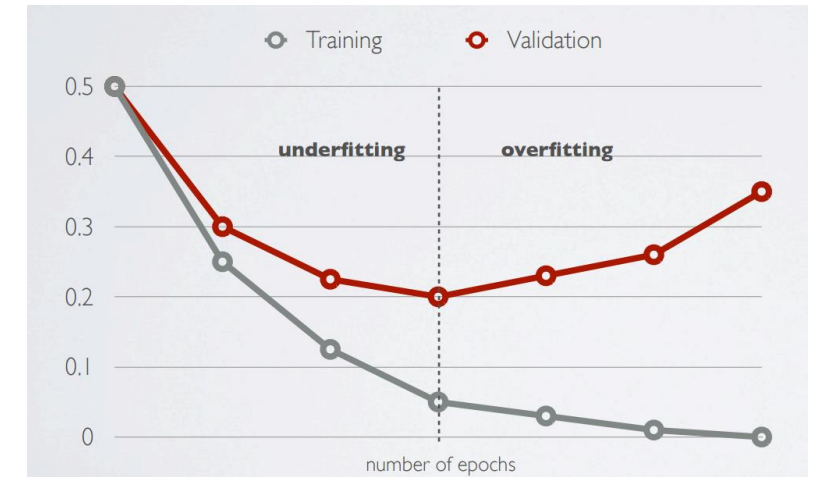


# What is Overfitting?



## Fighting overfitting

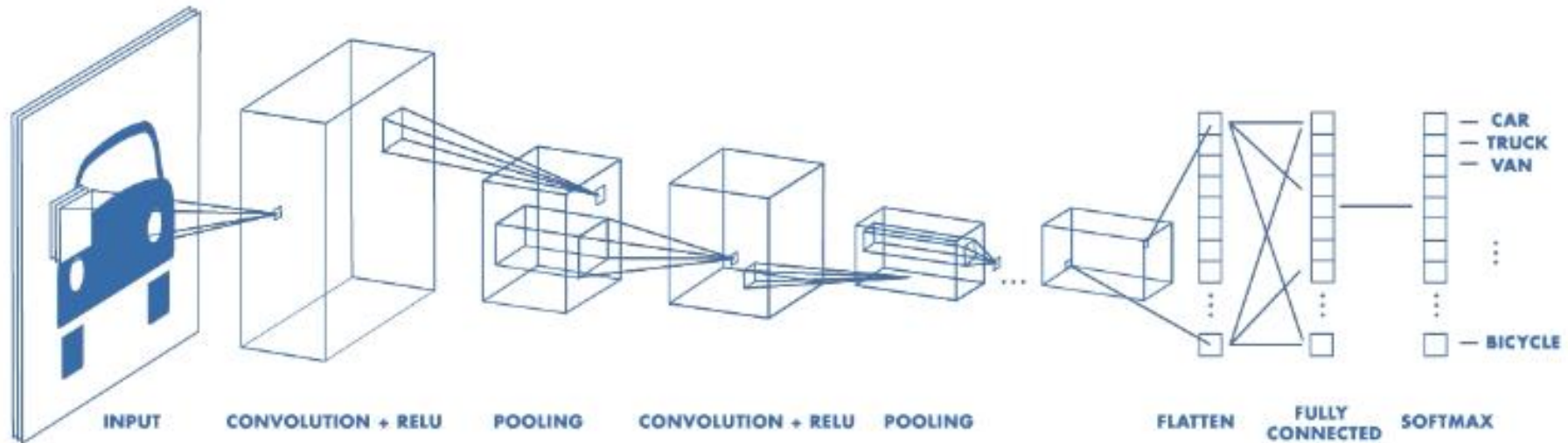
- Smaller network
- Weight decay
- Dropout
- Early stopping
- More data
  - Data augmentation





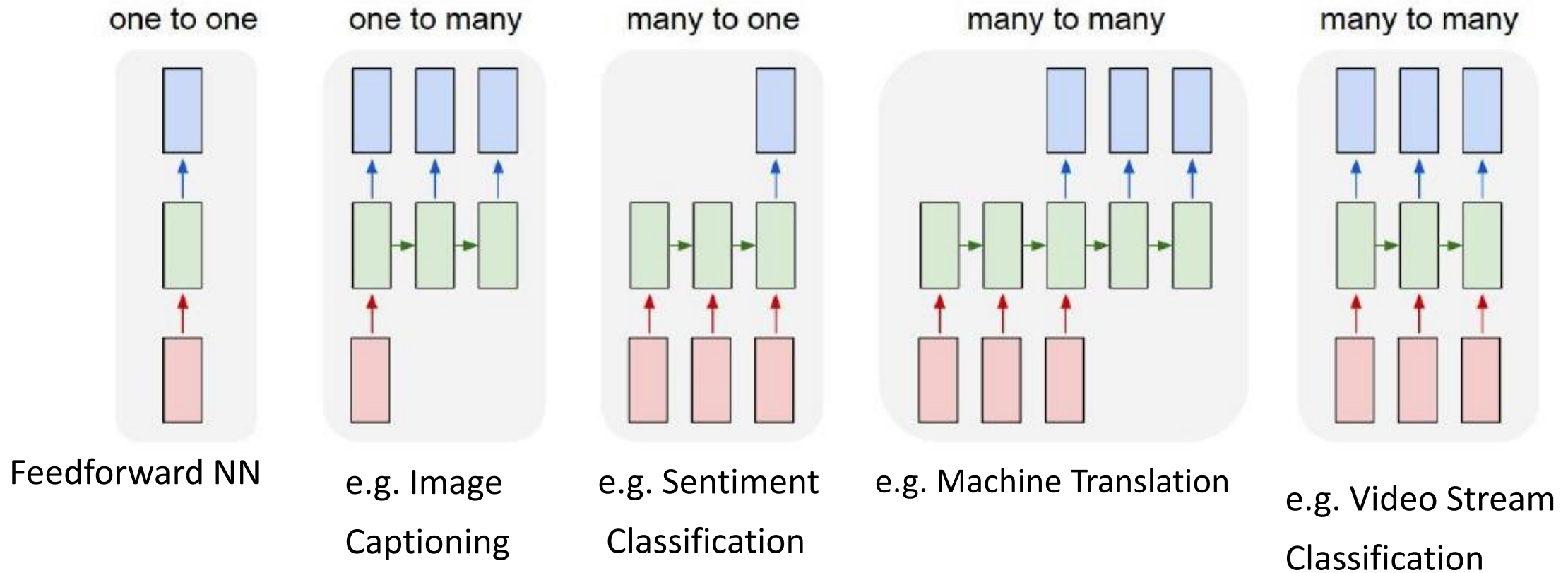
# Popular Deep Learning models

- **Convolutional Neural Networks (CNN)** – process Image-like data



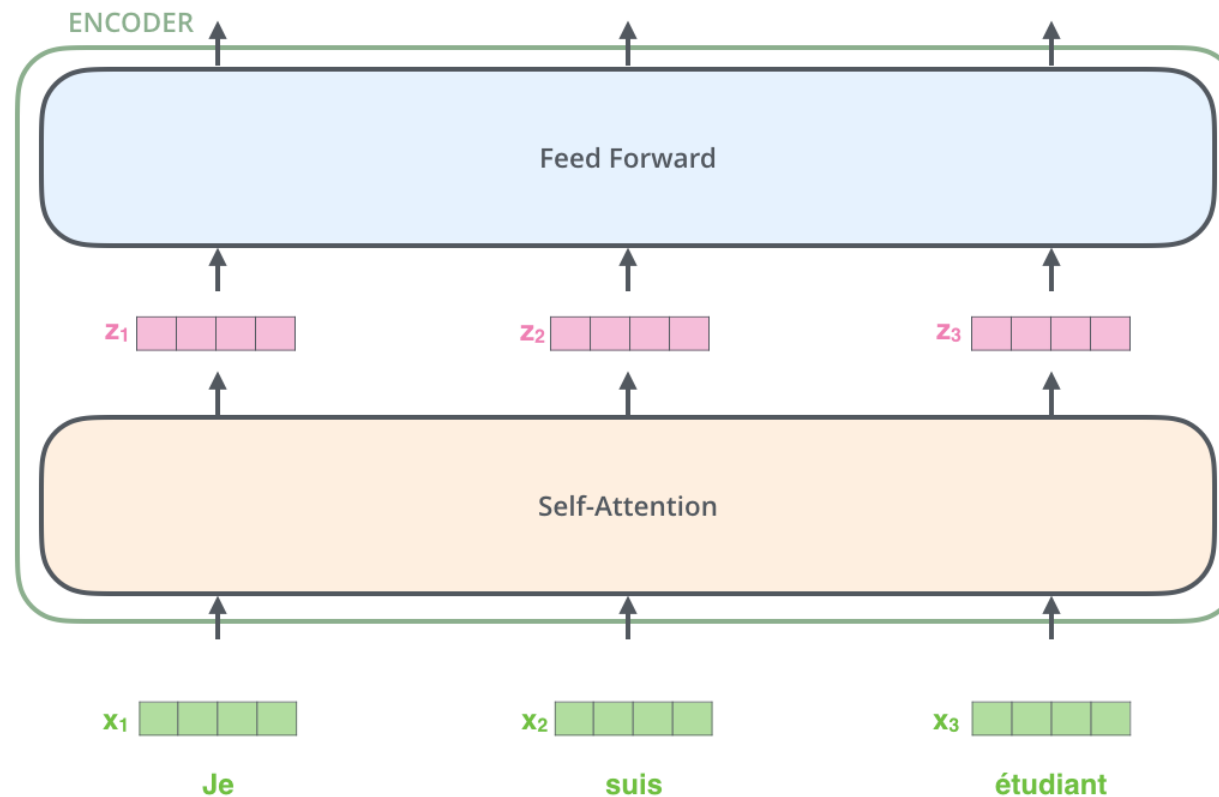
# Popular Deep Learning models

- **Recurrent Neural Networks (RNN)** –process sequences



# Popular Deep Learning models

- **Transformers** –process sets and sequences



# Some terminology

- Optimizer
- Loss
- Weights
- Backprop
- Learning rate
- Batch size
- Weight decay
- Dropout
- Hyperparameters
- Batch Norm
- Layer
- Activations
- Augmentation
- Softmax
- .....

Thank you!

