

# Advanced Machine Learning, 22 Aug 2019

Feedforward NN — nonlinearity

Cross entropy vs log likelihood

Output activation functions

$z$

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

0-1 output

Multiclass classification

$z_1, z_2, \dots, z_k$

$$\text{Softmax}_{\text{argmax } j} = \frac{\exp(z_j)}{\sum \exp(z_i)}$$

Softargmax

Crossentropy

$$\log P(z_i|x)$$

$$\log \frac{e^{z_j}}{\sum e^{z_i}} = z_j - \log \sum e^{z_i}$$

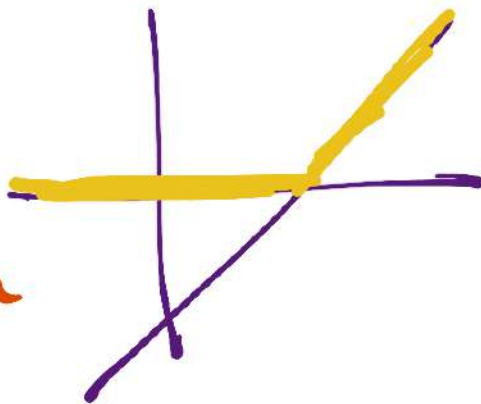
---

Interior nodes - nonlinearity

$$x \rightarrow \phi(x) \rightarrow z \rightarrow f(z)$$

RELU

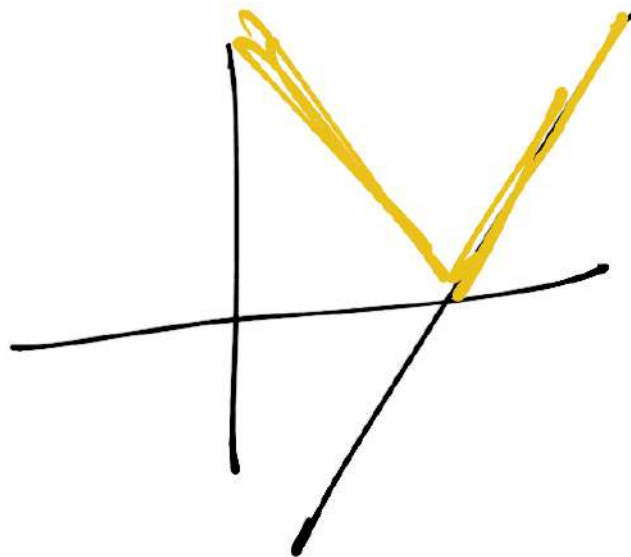
Rectified Linear  
Unit



$$\max(0, z)$$

Generalizations to handle zero gradient

Absolute value  
rectification



More generally

$$|z| = \max(0, z) - \min(0, z)$$
$$\max(0, z) + \alpha \min(0, z)$$

A purple arrow points from the  $\alpha$  in the second equation to the  $-$  sign in the first equation, with the text  $\alpha = -1$  written next to it.

# Architecture

How many layers?

Depth

How many nodes per layer?

Width

How are the layers connected?

↳ Assume completely connected

Recall:

Can construct a 1 hidden layer network to approximate any  $f$  using steps

Theoretically, 1 hidden layer suffices

- Width explodes

- Parameters explode

Tradeoff in using more layers is to  
reduce width per layer

Adding layers reduces parameters

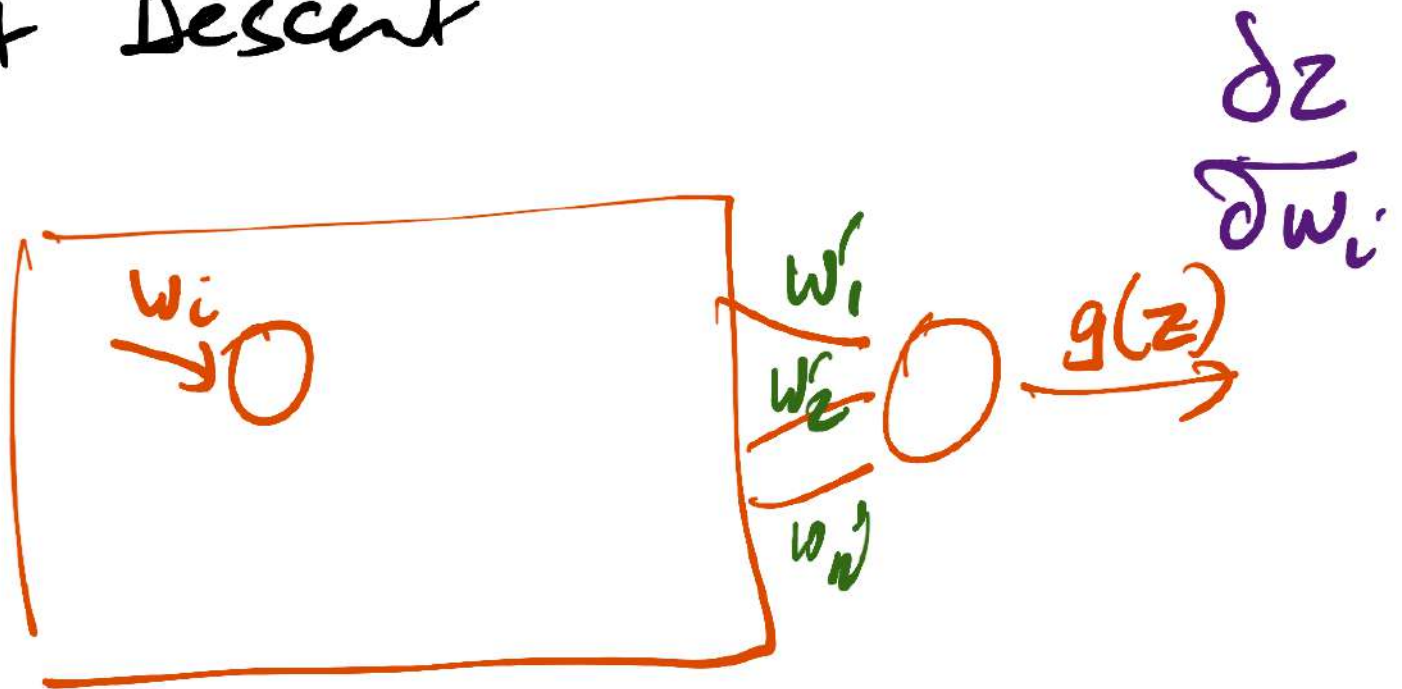
Each layer is an abstract "feature"

How many layers, how wide?

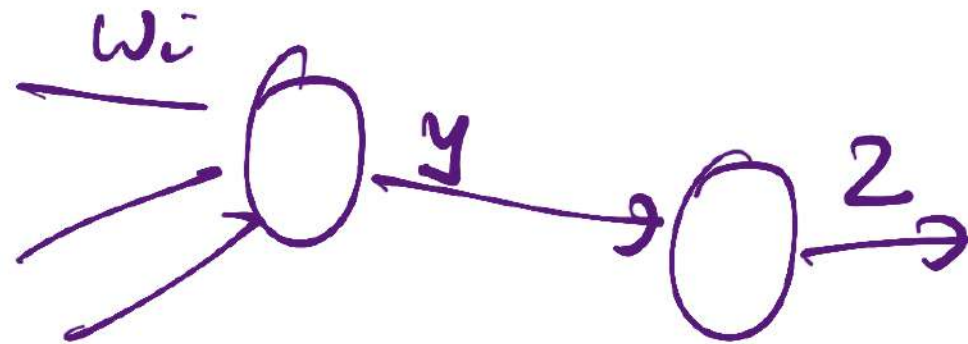
Empirical

Choose a value that you can compute!

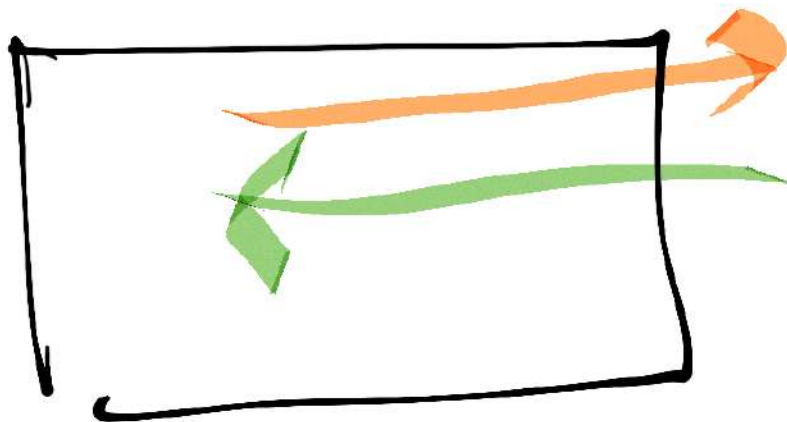
Gradient Descent



# Chain rule



$$\frac{\partial z}{\partial w_i} = \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial w_i}$$

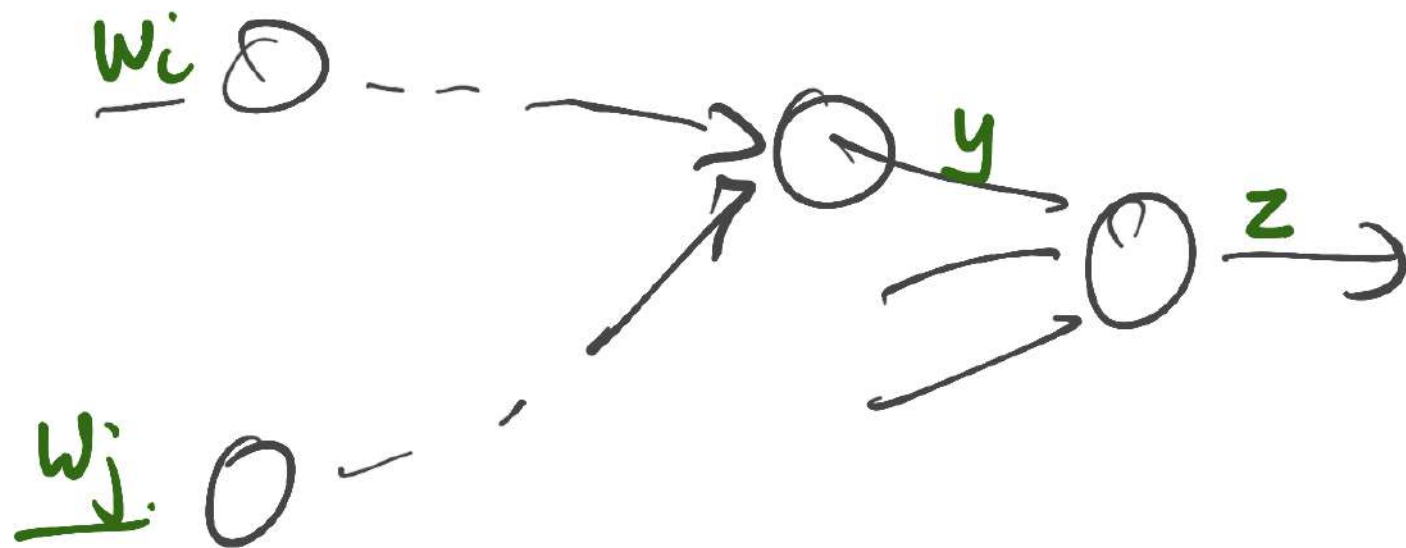


Compute outputs forward  
Compute gradients  
backwards

Backwards calculation exploits the fact that the network is a directed acyclic graph

Keep calculated values stored -

Dynamic Programming





Backpropagation can be done

- Numerically
  - Symbolically
  - Torch
  - Tensorflow
- 

- Loss function
- Activation function
- Structure of network
- Gradient descent

# Optimizing the process of learning

1. Tune the learning to generalize better

2. Speed it up

1. Regularization

Add a cost based on model complexity

Postpone / Skip

## 2. Speed up the computation

Minimizing loss  
w.r.t training  
data

$$E_{(x,y) \sim P_{\text{data}}} L(f(x; \theta), y)$$

Diagram illustrating the loss function  $L(f(x; \theta), y)$ . The input  $x$  is labeled "input" and the real output  $y$  is labeled "real output". The function  $f$  takes  $x$  and parameters  $\theta$  (labeled "Loss function parameters") as input and produces the predicted output  $\hat{y}$ . The loss  $L$  is calculated between the predicted output  $\hat{y}$  and the real output  $y$ . The expectation  $E$  is taken over the data distribution  $P_{\text{data}}$ .

What we really  
want

Minimizing risk

$$E_{(x,y) \sim P_{\text{data}}} L(f(x; \theta), y)$$

The expectation  $E$  and the data distribution  $P_{\text{data}}$  are underlined in yellow.

Natural loss function is 0-1 loss

0 if  $f(x) = y$

1 if  $f(x) \neq y$

↓ "surrogate" loss functions

log likelihood

# Gradient descent

Most cases, loss function is not convex

Batch - take mean over entire  
input

Single input - adjust per input

Mini batches - Stochastic Gradient  
Descent

# Statistically

Sample of size  $n$

True mean  $\mu$

Variance

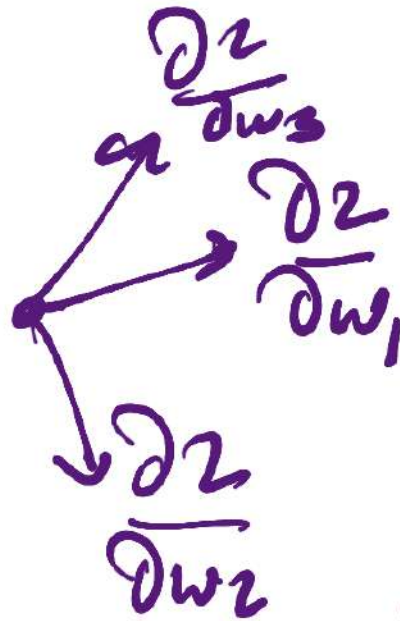
$$\frac{\mu}{\sqrt{n}}$$

Important to choose minibatches  
as random samples

# Issues

1. Multiple minima
2. Plateaus
3. Ill-conditioning

Learning rate  $\epsilon$



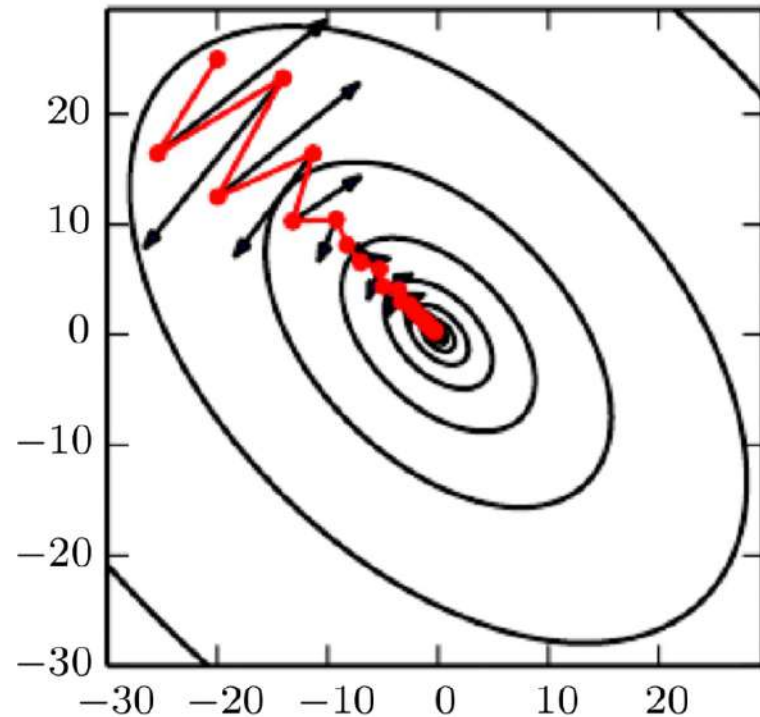
Values are very different

Add a parameter called "momentum"

Adjusting  $\epsilon$

Rate increases if direction does not change.

Rate decreases if direction keeps changing.





Variations Nesterov momentum

Hyperparameter - weight of the  
momentum term

---

Directly adapt  $\epsilon$

Adaptive learning rate

Delta-bar-delta

Decrease  $\epsilon$  if sign  
changes  
Increase  $\epsilon$  otherwise

Delta-lan-delta

- sign wrt entire history

Narrow the history to a recent window

AdaGrad

$$- \frac{1}{\sqrt{\text{sum of all gradient}}}$$

Bigger past values give bigger changes w.  $\epsilon$

RMS Prop - exponential decay of  
Delta-bar-Delta

Adam - Adaptive moments

RMS Prop + momentum

---

Initialization

Black magic

Identical hidden nodes - different values

Assign randomly

— Normal

— Uniform

Use small weight