# DMML, 19 Mar 2019

## Neural networks

### Structure

### Backpropagation algorithm

Activation function : Sigmoid

Cost function : MSE

Training :

Mini-batches — stochastic gradient descent

Epoch — full training sample

# Hyperparameters

Structure — no. of layers, size of each layer

Batch size ⎤
Epochs  ⎦ Backpropagation

Activation function
Cost function
     └ By default, sigmoid + MSE
      Is this always a good choice?

Expectation of learning weights & biases

- Get closer to good value, rate slows down
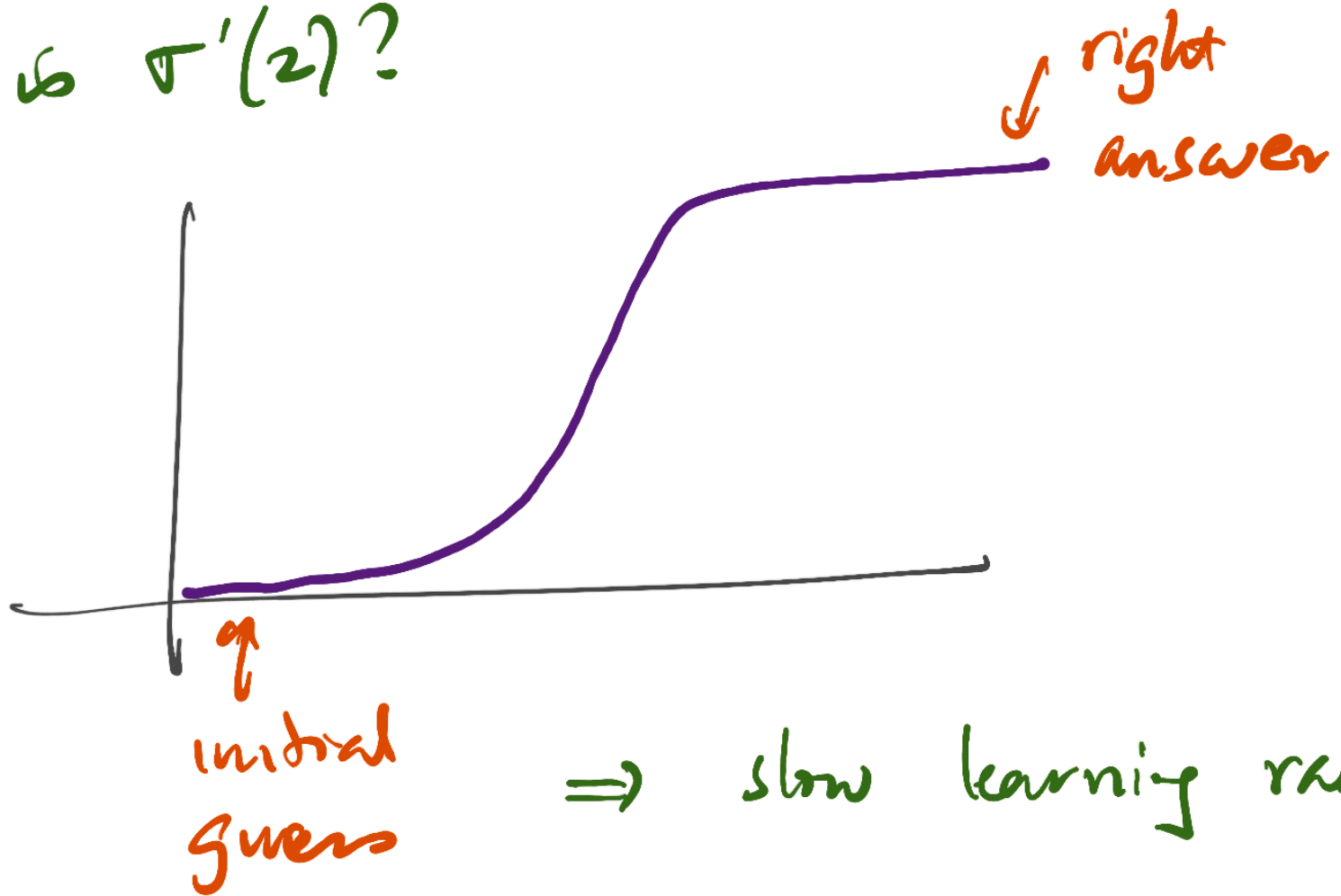- Very far away — progress is rapid

Nielsen example
- Second assumption is not valid

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

$$\text{Cost} = \frac{1}{2n} \sum (y-a)^2$$

$\frac{\partial C}{\partial w}, \frac{\partial C}{\partial b}$ are proportional to $\sigma'(z)$

What is $\sigma'(z)$?

right answer

initial guess

$\Rightarrow$ slow learning rate

Want: $\dfrac{\partial C}{\partial w}, \dfrac{\partial C}{\partial b}$ proportional to $\|y-a\|$

By reverse engineering

Cross Entropy (not the "real" cross entropy)

$$\text{Cost} = \frac{1}{n} \sum_x y \ln a + (1-y) \ln(1-a)$$
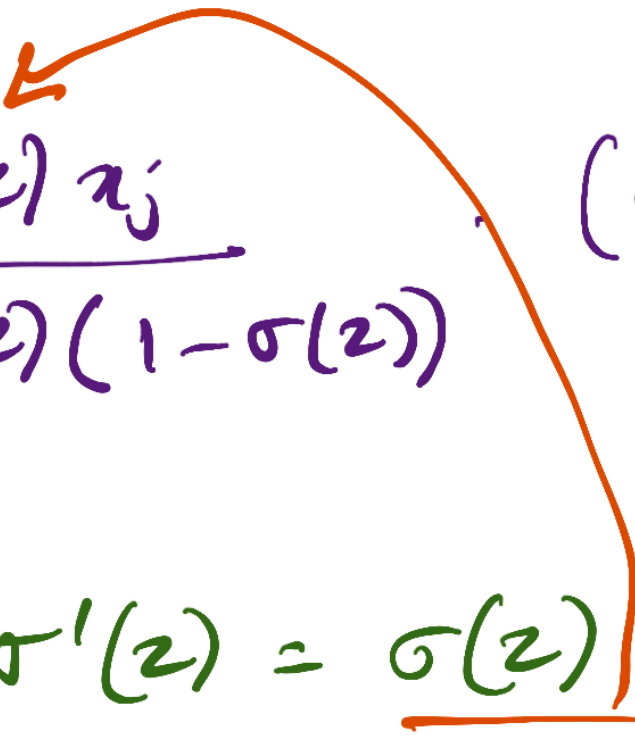
Why?   By reverse engg

Justify?   $y = 1 \Rightarrow C = 0$   $y = 0 \Rightarrow C \to 0$
          $a \approx 1$                $a \approx 0$

**Sigmoid**

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

Calculate $\frac{\partial C}{\partial w_j}$, $\frac{\partial C}{\partial b}$ for $C$ = cross entropy

$$\frac{\partial C}{\partial w_j} = \frac{1}{n} \sum_x \frac{\sigma'(z) x_j}{\sigma(z)(1-\sigma(z))} \cdot \underbrace{(\sigma(z) - y)}_{a - y}$$

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

cancels out

# Sigmoid + Cross Entropy

$$\frac{\partial C}{\partial w_j} = \frac{1}{n} \sum_x x_j \left(\sigma(z) - y\right)$$

$$\frac{\partial C}{\partial b} = \frac{1}{n} \sum_x \left(\sigma(z) - y\right)$$

Another option for activation

0 - 9 digit recognition — 10 outputs, take max

## Softmax

$a_0$
$a_1$
$\vdots$
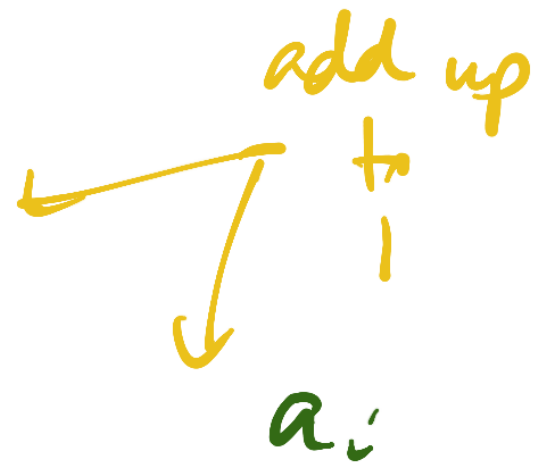$a_9$

$$\frac{a_i}{\sum a_j}$$

Instead Softmax:  $$\frac{e^{a_i}}{\sum e^{a_j}}$$

add up to 1

Use a different cost function

$$-\ln(a) \qquad \text{"log likelihood"}$$

In practice

Sigmoid + Cross Entropy

Softmax + log likelihood

# Regularization

Smoothen out the learning process

Penalize large weights

$$\text{Add} \quad \text{---} \quad + \quad \underline{\underline{\lambda \Sigma \ w^2}}$$

regularization term

to cost function

# Deep learning

Theoretically, one hidden layer can approximate any continuous function

- High fan in / fan out of signals

Impractical

Build layers that make incremental progress

Image recognition

- layer detects horizontal edges ⎫
- layer " vertical edges ⎬ features
      ⋮                          ⎭

Combine layers to classify

- Multiple layers → deep networks

# Deep learning

- Neural network with "many" layers
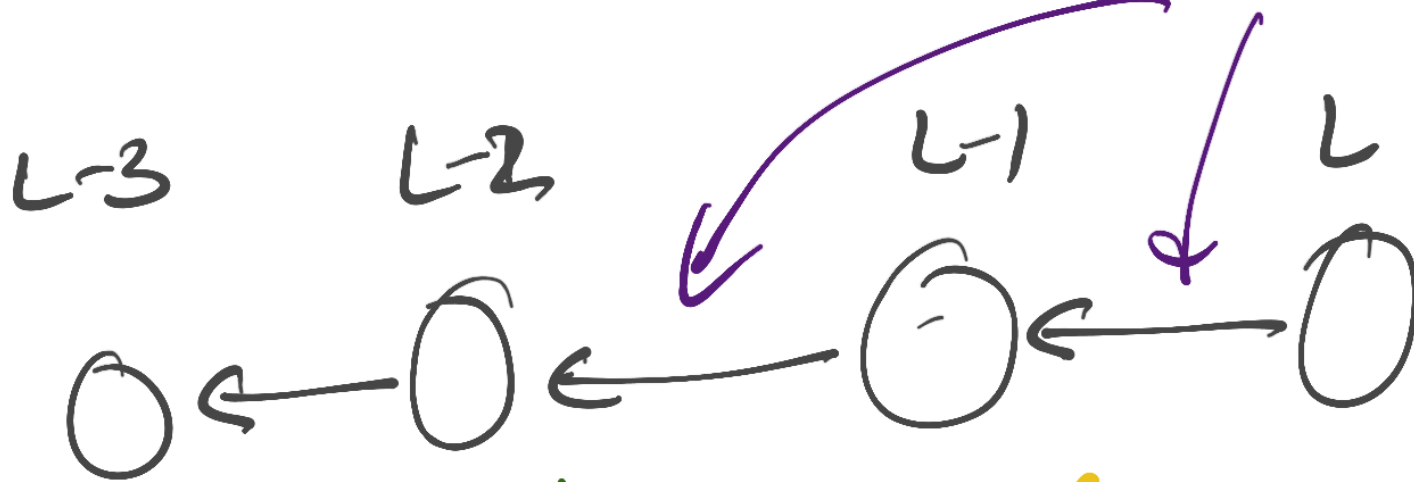
By hand, "many" $\leq 10$

Brute force, $> 100$

## Slowdown in learning rate

Intrinsic to backpropagation

$$\delta_j^l = \frac{\partial C}{\partial z_j^l}$$

Express $\delta_j^{\ell}$ in terms of $\delta_j^{\ell+1}$

$$\delta_j^{\ell} = \sum_k \delta_k^{\ell+1} \cdot W_{kj}^{\ell+1} \sigma'(z_j^{\ell})$$

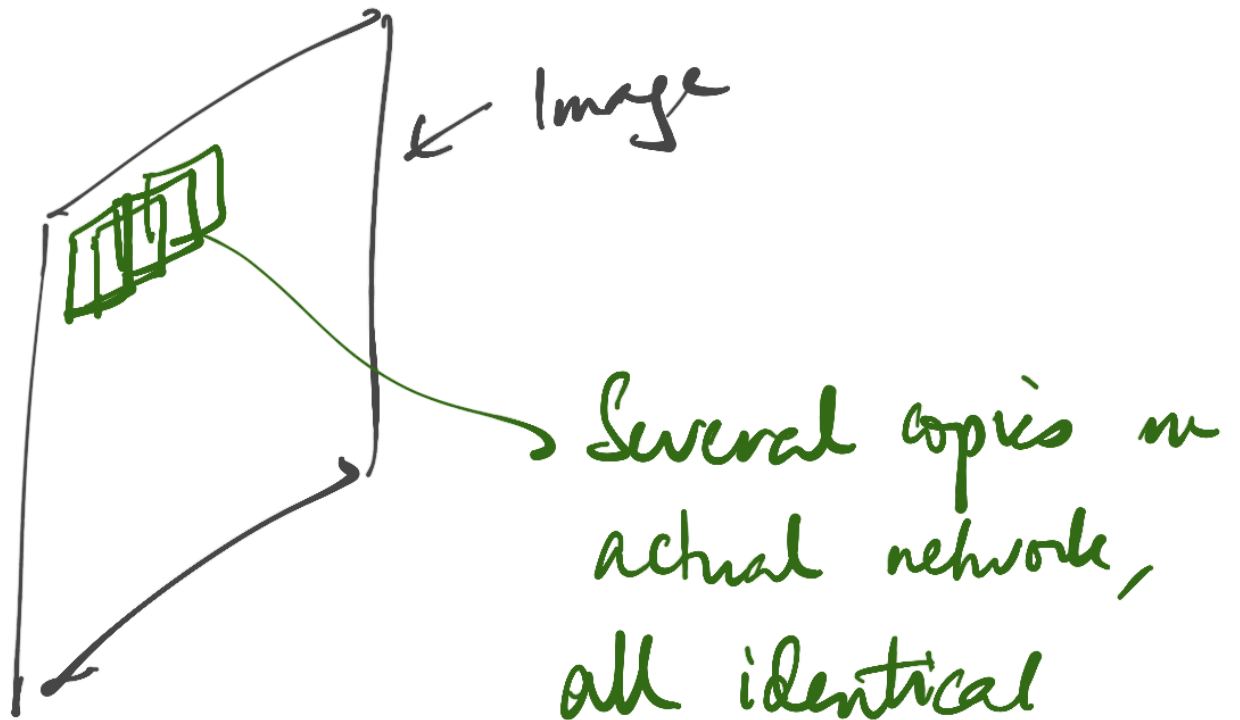L-3          L-2          L-1          L

As we move left, $\delta_j^{\ell}$ "vanishes"

"Vanishing gradient" problem

# Solution?

Add structure to the network

Design a small "filter" that checks for a feature in a small part of input



← Image

Several copies in actual network, all identical

All copies of our filter have exactly same weights & biases

Convolution - Passing a filter over an input

Convolutional Neural Networks

L Layer is a collection of identical small filters