# Backpropagation: A modular approach (Torch NN)
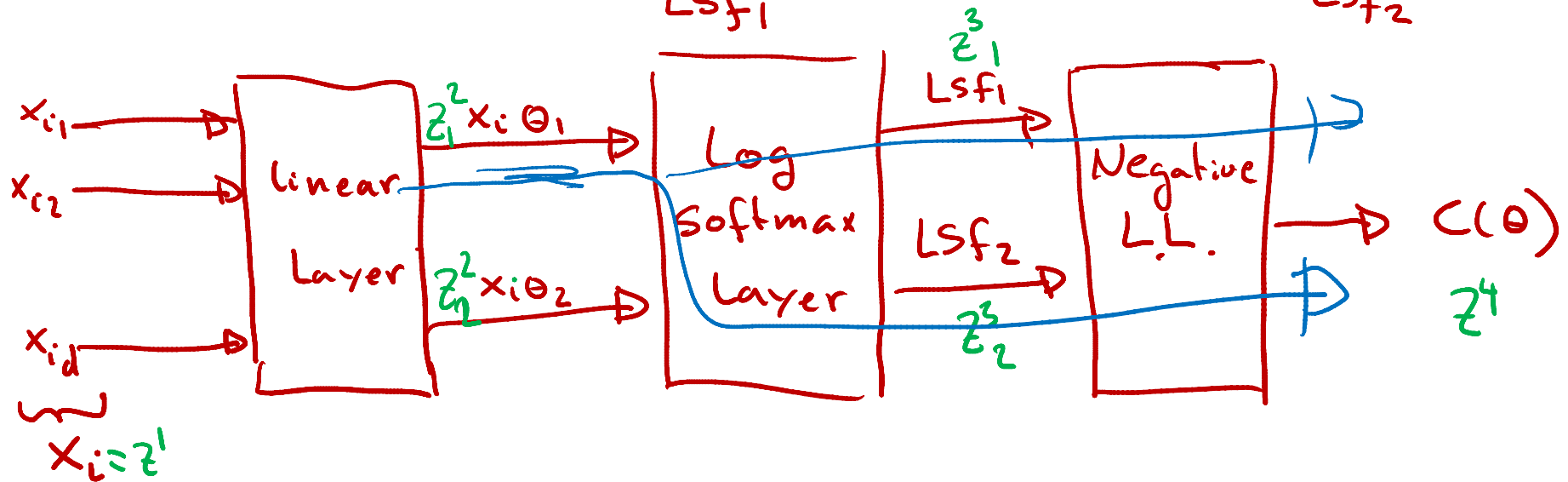
Nando de Freitas

UNIVERSITY OF OXFORD

# Outline of the lecture

This lecture describes modular ways of formulating and learning distributed representations of data. The objective is for you to learn:

- ❑ How to specify models such as logistic regression in layers.
- ❑ How to formulate layers and loss criterions.
- ❑ How well formulated <span style="color:red">local</span> rules results in correct <span style="color:red">global</span> rules.
- ❑ How back-propagation works.
- ❑ How this manifests itself in <span style="color:blue">Torch</span>.

$$C(\theta) = - \sum_{i=1}^{n} \mathbb{1}_0(Y_i) \log\left(\frac{e^{x_i\theta_1}}{e^{x_i\theta_1} + e^{x_i\theta_2}}\right) + \mathbb{1}_1(Y_i) \log\left(\frac{e^{x_i\theta_2}}{e^{x_i\theta_1} + e^{x_i\theta_2}}\right)$$

$\underbrace{\phantom{xxxxxxxxxxxxxxxx}}_{LSf_1}$  $\underbrace{\phantom{xxxxxxxxxxxxxxxx}}_{LSf_2}$



$$z_1^2 = x_i\theta_1 \qquad z_2^2 = x_i\theta_2$$

$$z_1^3 = \log\left(\frac{e^{z_1^2}}{e^{z_1^2} + e^{z_2^2}}\right) \qquad z_2^3 = \log\left(\frac{e^{z_2^2}}{e^{z_1^2} + e^{z_2^2}}\right)$$

$$z^4 = \sum_i \mathbb{1}_0(Y_i) z_1^3 + \mathbb{1}_1(Y_i) z_2^3$$
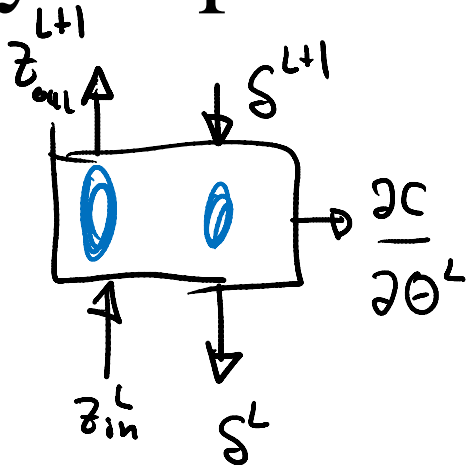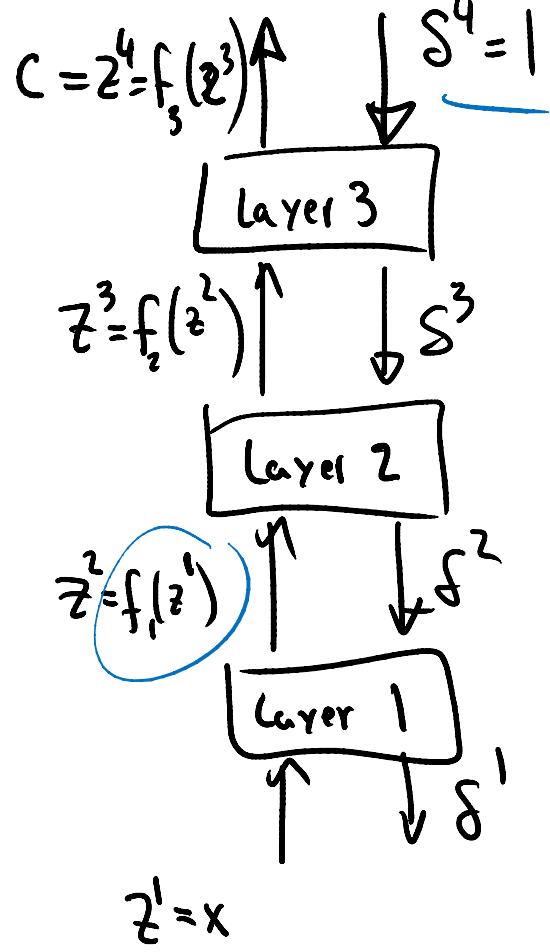
# Derivative using the chain rule

$$C(\theta) = -\sum_{i=1}^{n} \mathbb{I}_0(Y_i) \log\left(\frac{e^{x_i\theta_1}}{e^{x_i\theta_1} + e^{x_i\theta_2}}\right) + \mathbb{I}_1(Y_i) \log\left(\frac{e^{x_i\theta_2}}{e^{x_i\theta_1} + e^{x_i\theta_2}}\right)$$

$$C(\theta) = z^4\left\{ z_1^3\left[z_1^2(\theta_1, z^1), z_2^2(\theta_2, z^1)\right], z_2^3\left[z_1^2(\theta_1, z^1) z_2^2(\theta_2, z^1)\right]\right\}$$

$$\frac{\partial C(\theta)}{\partial \theta_1} = \frac{\partial z^4}{\partial z_1^3} \frac{\partial z_1^3}{\partial z_1^2} \frac{\partial z_1^2}{\partial \theta_1} + \frac{\partial z^4}{\partial z_1^3} \frac{\partial z_1^3}{\partial z_2^2} \frac{\partial z_2^2}{\partial \theta_1}$$

$$+ \frac{\partial z^4}{\partial z_2^3} \frac{\partial z_2^3}{\partial z_1^2} \frac{\partial z_1^2}{\partial \theta_1} + \frac{\partial z^4}{\partial z_2^3} \frac{\partial z_2^3}{\partial z_2^2} \frac{\partial z_2^2}{\partial \theta_1}$$
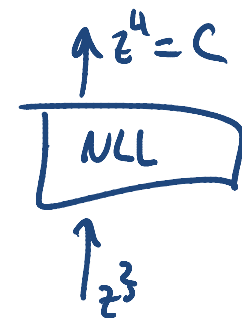
# Layer specification

$$\delta^4 = 1$$

$$C = z^4 = f_3(z^3)$$

$$z^3 = f_2(z^2)$$

$$\delta^3$$

Layer 3

Layer 2

$$z^2 = f_1(z^1)$$

$$\delta^2$$

Layer 1

$$\delta^1$$

$$z^1 = x$$

$$z^{L+1}_{out}$$

$$\delta^{L+1}$$

$$\frac{\partial C}{\partial \Theta^L}$$

$$z^L_{in}$$

$$\delta^L$$
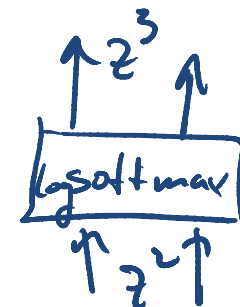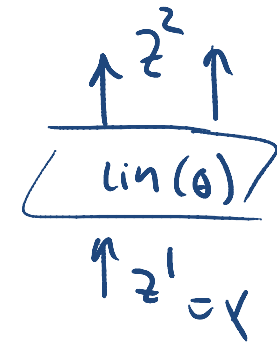
$$z^{L+1} = f(z^L)$$
forward pass

**backward pass**

$$\delta^L_i = \frac{\partial C}{\partial z^L_i} = \sum_j \frac{\partial C}{\partial z^{L+1}_j} \frac{\partial z^{L+1}_j}{\partial z^L_i} = \sum_j \delta^{L+1}_j \boxed{\frac{\partial z^{L+1}_j}{\partial z^L_i}}$$

$$\frac{\partial C}{\partial \Theta^L} = \sum_j \frac{\partial C}{\partial z^{L+1}_j} \frac{\partial z^{L+1}_j}{\partial \Theta^L} = \sum_j \delta^{L+1}_j \left( \frac{\partial z^{L+1}_j}{\partial \Theta^L} \right)$$

# Derivative via layer-specification

$$\frac{\partial c}{\partial \theta_1} = \sum_j \frac{\partial c}{\partial z_j^2} \frac{\partial z_j^2}{\partial \theta_1}$$

$$= \sum_j \left( \sum_k \frac{\partial c}{\partial z_k^3} \frac{\partial z_k^3}{\partial z_j^2} \right) \frac{\partial z_j^2}{\partial \theta_1}$$

$$= \sum_{j=1}^{2} \sum_{k=1}^{2} \frac{\partial c}{\partial z^4} \left( \frac{\partial z^4}{\partial z_k^3} \frac{\partial z_k^3}{\partial z_j^2} \frac{\partial z_j^2}{\partial \theta_1} \right)$$
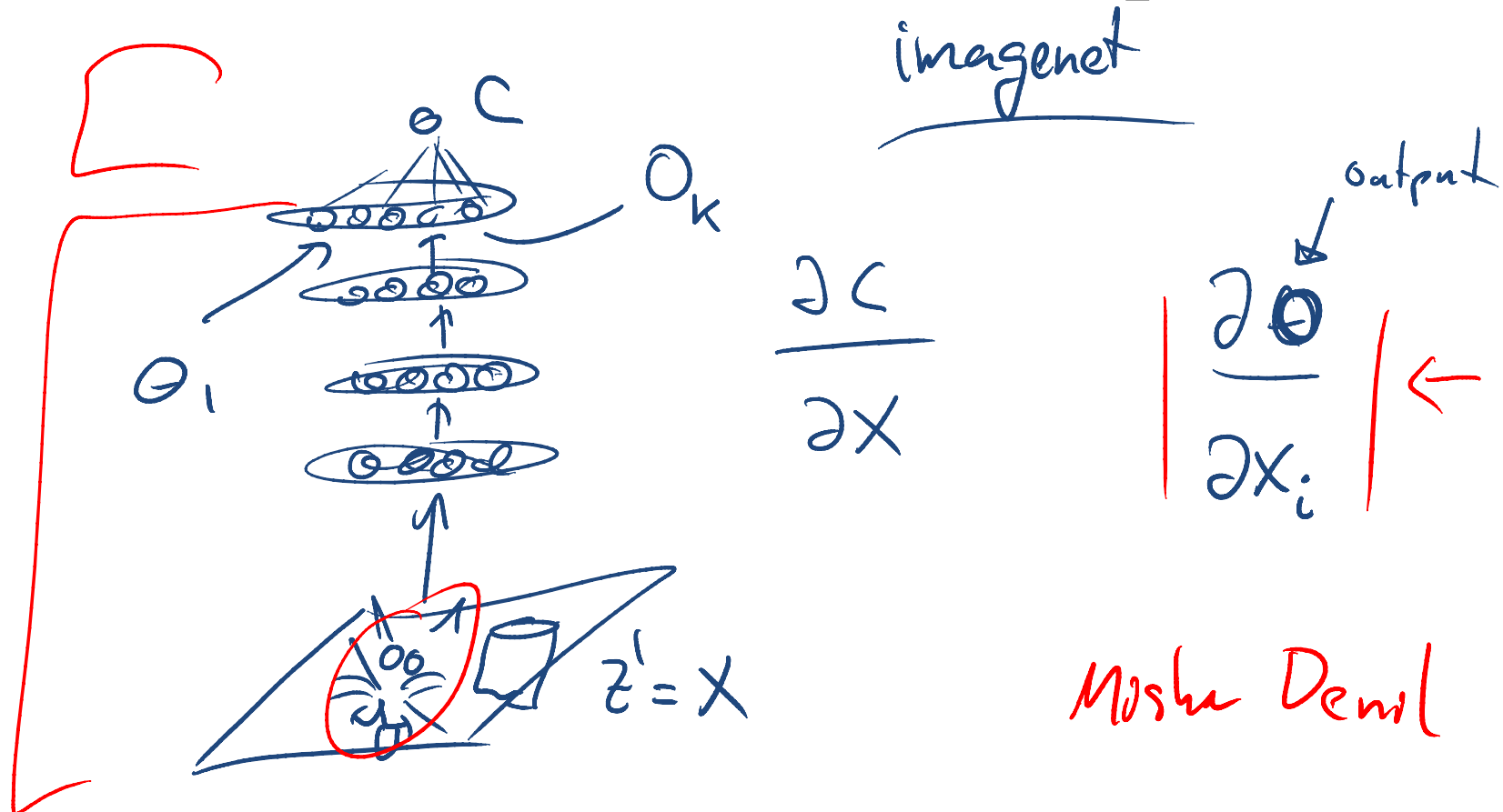
$$= \text{as before.}$$

$z^2 \uparrow \quad \uparrow$

$\boxed{\text{Lin} (\theta)}$

$\uparrow z^1 = x$

$z^3 \uparrow \quad \uparrow$

$\boxed{\text{logsoftmax}}$

$\uparrow z^2 \uparrow$

$\uparrow z^4 = C$

$\boxed{\text{NLL}}$

$\uparrow z^3$

# Back-propagation algorithm

$$z^1 = x_i \xrightarrow{\underbrace{x_i \theta_1}} z^2(x_i) \xrightarrow{\log\left(\dfrac{e^{x_i \theta_1}}{\Sigma}\right)} z^3(x_i) \longrightarrow z^4_{(x_i)} = c$$

$$\delta^1 \longleftarrow \delta^2 \longleftarrow \delta^3 \longleftarrow \delta^4 \leftarrow 1$$

# Derivatives wrt to the input



imagenet

$\in C$

$O_k$

$\Theta_1$

$z' = X$

$\dfrac{\partial C}{\partial X}$

output

$\left| \dfrac{\partial \Theta}{\partial X_i} \right| \leftarrow$

Karen Simonyan

Moshe Demil

# Logit Regression Model in Torch

```
1    model = nn.Sequential()
2    model:add( nn.Linear(2,2) )
3    model:add( nn.LogSoftMax() )
```
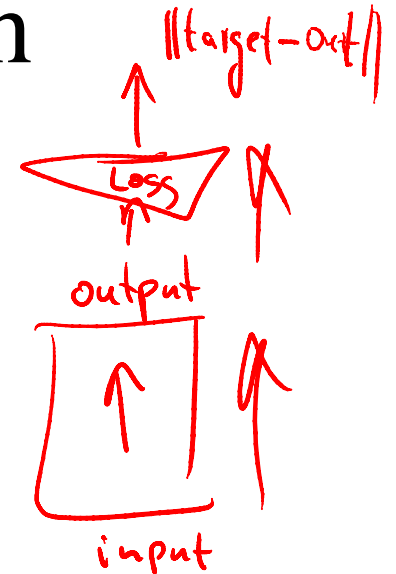
# Loss criterion in Torch

```
1 | criterion = nn.ClassNLLCriterion()
```

# Derivatives closure in Torch

||target−out||

```lua
-- params/gradients
x, dl_dx = model:getParameters()
```

output

input

loss

```lua
local loss_x = criterion:forward(model:forward(inputs), target)
model:backward(inputs, criterion:backward(model.output, target))
```

# Optimization in Torch

```
_,fs = optim.sgd(feval,x,sgd_params)
```

```
-- Functions in optim all return two things:
--    + the new x, found by the optimization method (here SGD)
--    + the value of the loss functions at all points that were used by
--      the algorithm. SGD only estimates the function once, so
--      that list just contains one value.
```

# Next lecture

In the next lecture, we consider a generalization of logistic regression, with many logistic units, called multi-layer perceptron (MLP) or feed-forward neural network.