

A black and white photograph of a long, arched stone corridor. The corridor is formed by a series of large, square stone pillars supporting a vaulted ceiling. A single ornate lantern hangs from the ceiling in the foreground. The perspective leads the eye down the length of the corridor, which recedes into the distance. The lighting is dramatic, with strong shadows and highlights on the stone surfaces.

# Backpropagation: A modular approach (Torch NN)

Nando de Freitas



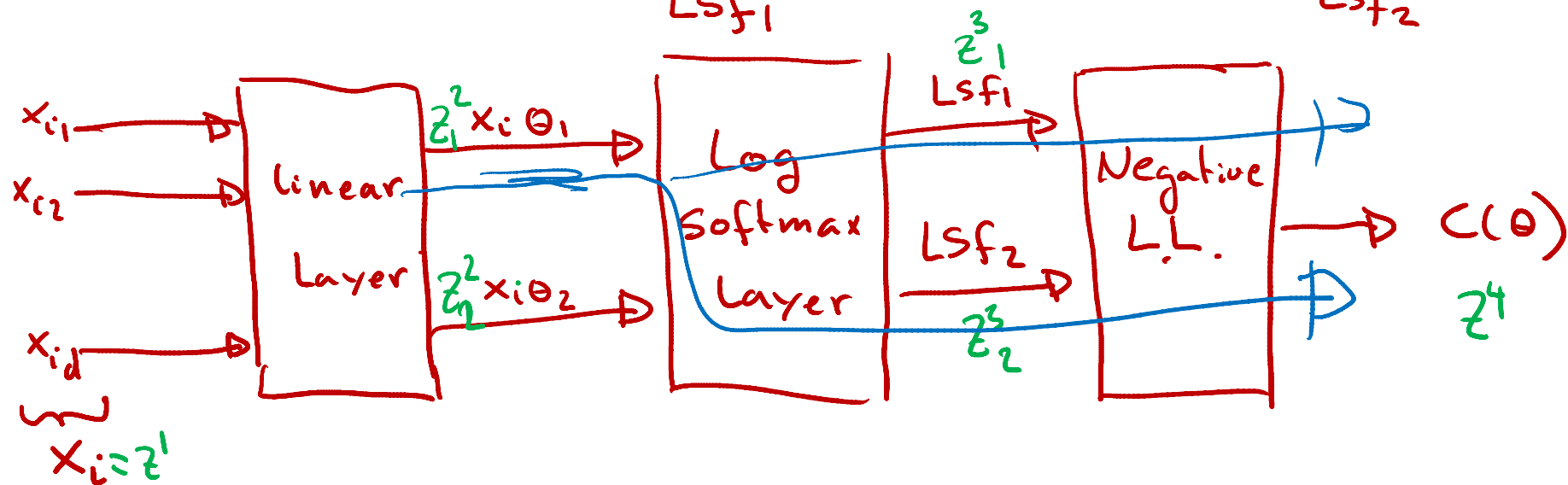
UNIVERSITY OF  
OXFORD

# 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 **local** rules results in correct **global** rules.
- ❑ How back-propagation works.
- ❑ How this manifests itself in **Torch**.

$$C(\theta) = - \sum_{i=1}^n \underbrace{\Pi_0(y_i) \log \left( \frac{e^{x_i \theta_1}}{e^{x_i \theta_1} + e^{x_i \theta_2}} \right)}_{Lsf_1} + \underbrace{\Pi_1(y_i) \log \left( \frac{e^{x_i \theta_2}}{e^{x_i \theta_1} + e^{x_i \theta_2}} \right)}_{Lsf_2}$$



$$z_1^2 = \underline{x_{i1} \theta_1}$$

$$z_2^2 = x_{i2} \theta_2$$

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

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

$$z^4 = \sum_i \Pi_0(y_i) z_1^3 + \Pi_1(y_i) z_2^3$$

# Derivative using the chain rule

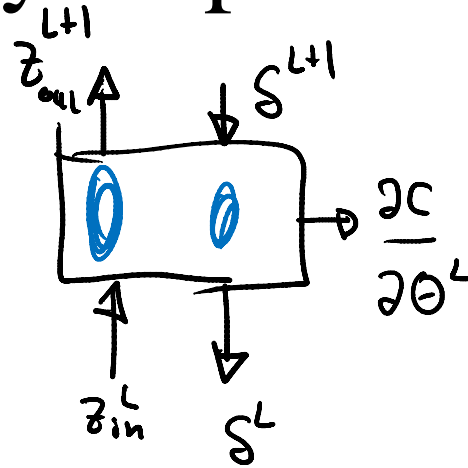
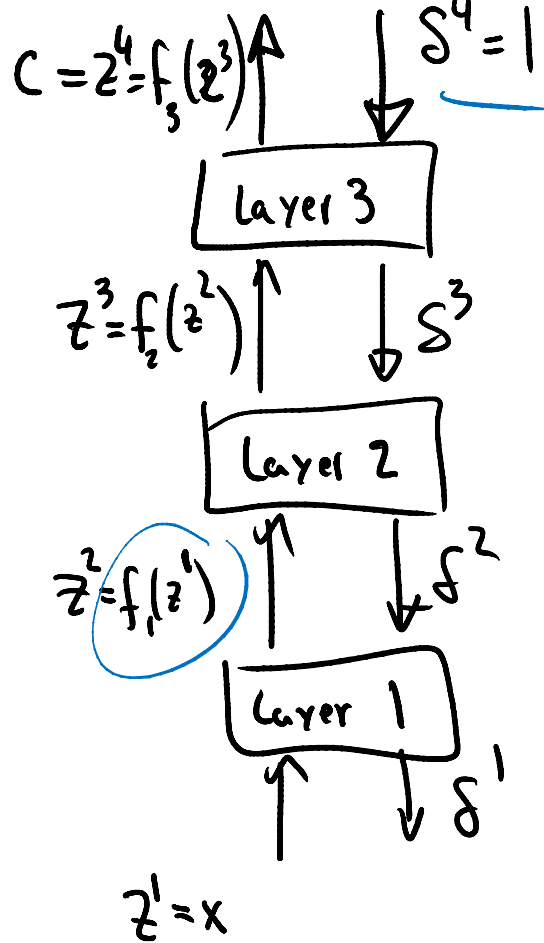
$$C(\theta) = - \sum_{i=1}^n \mathbb{I}_0(y_i) \log \left( \frac{e^{\underline{x_i \theta_1}} \overline{z_{i2}}}{e^{x_i \theta_1} + e^{x_i \theta_2}} \right) + \mathbb{I}_1(y_i) \log \left( \frac{e^{\overline{x_i \theta_2}}}{e^{x_i \theta_1} + e^{x_i \theta_2}} \right)$$

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

$$\begin{aligned} \frac{\partial C(\theta)}{\partial \theta_1} = & \underbrace{\frac{\overline{\partial z^4}}{\partial z_1^3} \frac{\partial z_1^3}{\partial z_1^2} \frac{\partial z_1^2}{\partial \theta_1}}_{\text{Term 1}} + \underbrace{\frac{\overline{\partial z^4}}{\partial z_1^3} \frac{\partial z_1^3}{\partial z_2^2} \frac{\partial \overline{z_2^2}}{\partial \theta_1}}_{\text{Term 2}} \\ & + \underbrace{\frac{\partial z^4}{\partial z_2^3} \frac{\partial z_2^3}{\partial z_1^2} \frac{\partial z_1^2}{\partial \theta_1}}_{\text{Term 3}} + \underbrace{\frac{\partial z^4}{\partial z_2^3} \frac{\partial z_2^3}{\partial z_2^2} \frac{\partial \overline{z_2^2}}{\partial \theta_1}}_{\text{Term 4}} \end{aligned}$$

←

# Layer specification



$$z^{L+1} = f(z^L)$$

forward pass

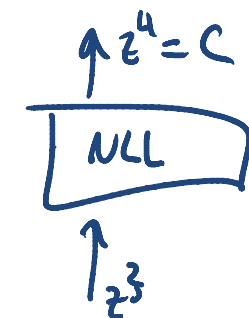
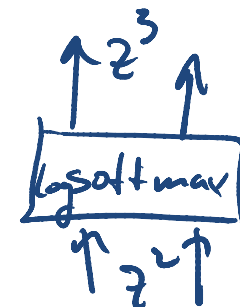
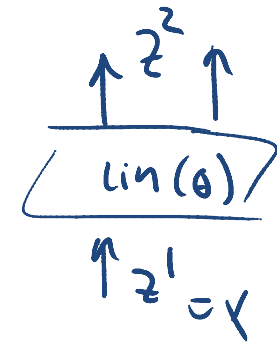
backward pass

$$\delta_i^L = \frac{\partial c}{\partial z_i^L} = \sum_j \frac{\partial c}{\partial z_j^{L+1}} \frac{\partial z_j^{L+1}}{\partial z_i^L} = \sum_j \delta_j^{L+1} \left[ \frac{\partial z_j^{L+1}}{\partial z_i^L} \right]$$

$$\frac{\partial c}{\partial \theta^L} = \sum_j \frac{\partial c}{\partial z_j^{L+1}} \frac{\partial z_j^{L+1}}{\partial \theta^L} = \sum_j \delta_j^{L+1} \left( \frac{\partial z_j^{L+1}}{\partial \theta^L} \right)$$

# Derivative via layer-specification

$$\begin{aligned}
 \frac{\partial C}{\partial \theta_1} &= \sum_j \underbrace{\frac{\partial C}{\partial z_j^2}}_{\substack{\uparrow \\ \text{ReLU}}} \frac{\partial z_j^2}{\partial \theta_1} \\
 &= \sum_j \left( \sum_k \left( \frac{\partial C}{\partial z_k^3} \frac{\partial z_k^3}{\partial z_j^2} \right) \right) \frac{\partial z_j^2}{\partial \theta_1} \\
 &= \sum_{j=1}^2 \sum_{k=1}^2 \frac{\partial C}{\partial z_k^4} \left( \frac{\partial z_k^4}{\partial z_k^3} \frac{\partial z_k^3}{\partial z_j^2} \right) \frac{\partial z_j^2}{\partial \theta_1} \\
 &= \text{as before.}
 \end{aligned}$$

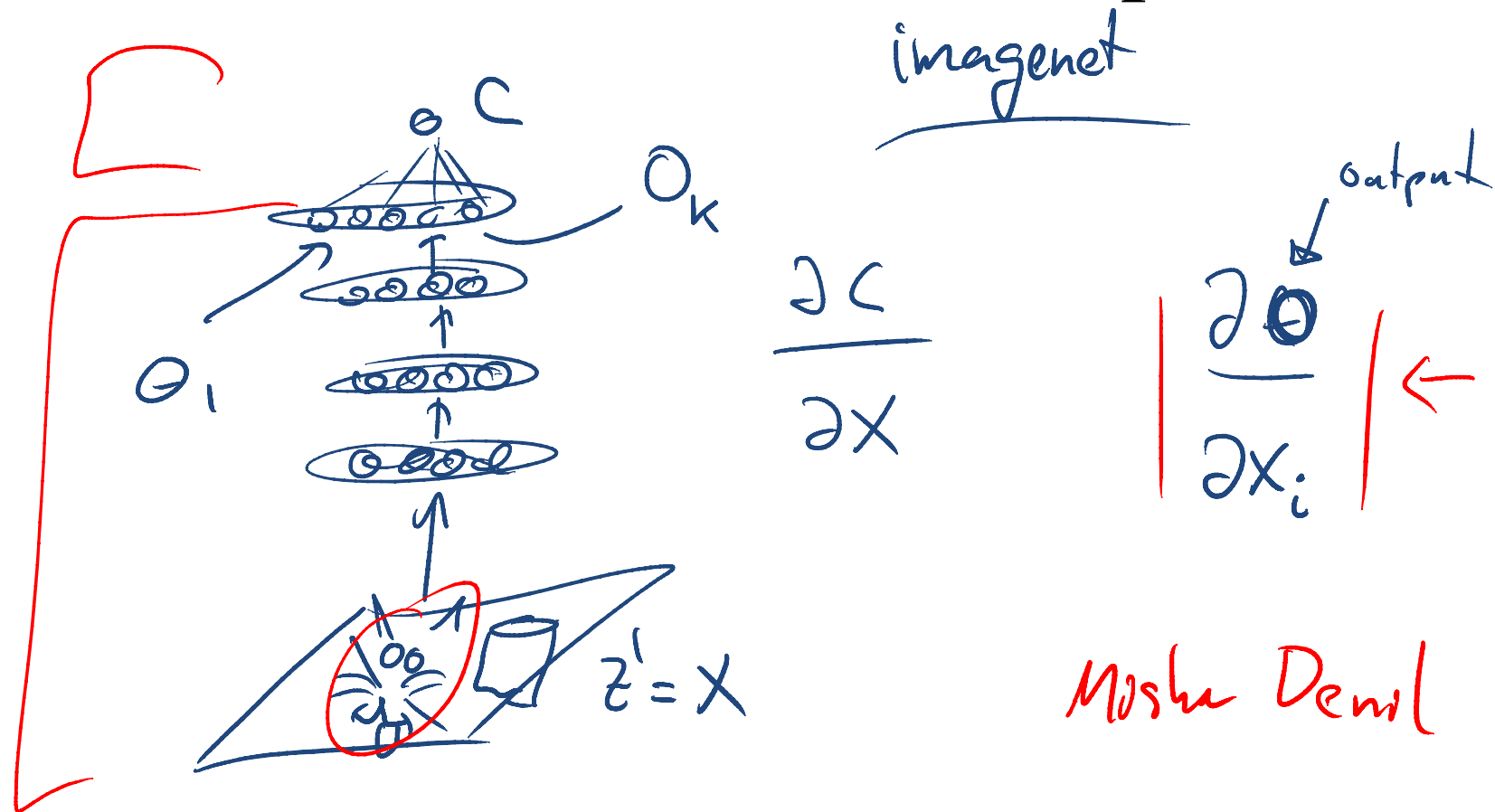


# Back-propagation algorithm

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

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

# Derivatives wrt to the input



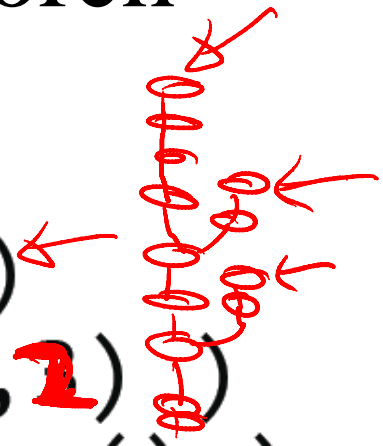
Misha Denil

Karen Simonyan



# Logit Regression Model in Torch

```
1 model = nn.Sequential()  
2 model.add( nn.Linear(2, 1) )  
3 model.add( nn.LogSoftMax() )
```

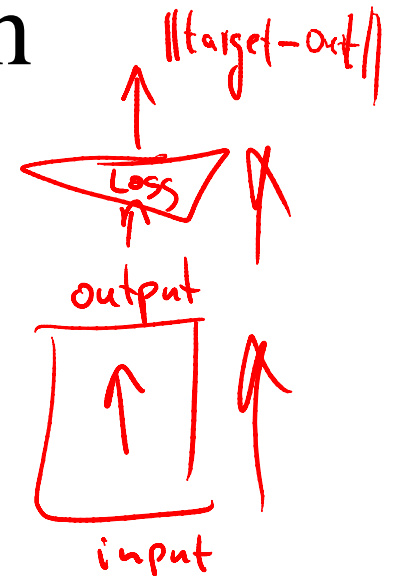


# Loss criterion in Torch

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

# Derivatives closure in Torch

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



```
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.