

# **DeepHack.CISS**

## **Neural Networks basics**

**Valentin Malykh**

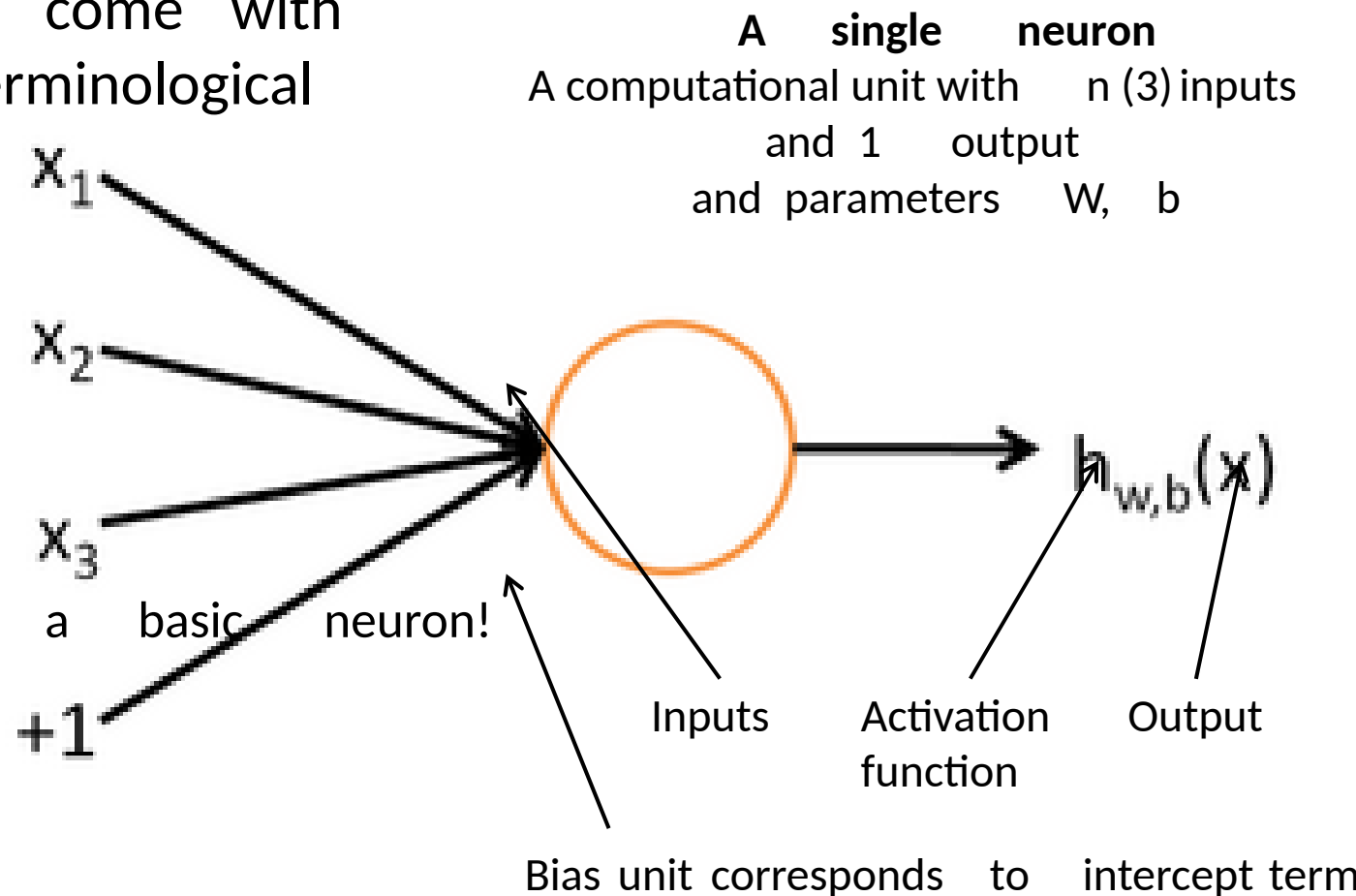
**Based on Stanford CS224d**

# Demystifying neural networks

Neural networks come with their own terminological baggage

But if you understand softmax models

Then **you** already understand a basic neuron!



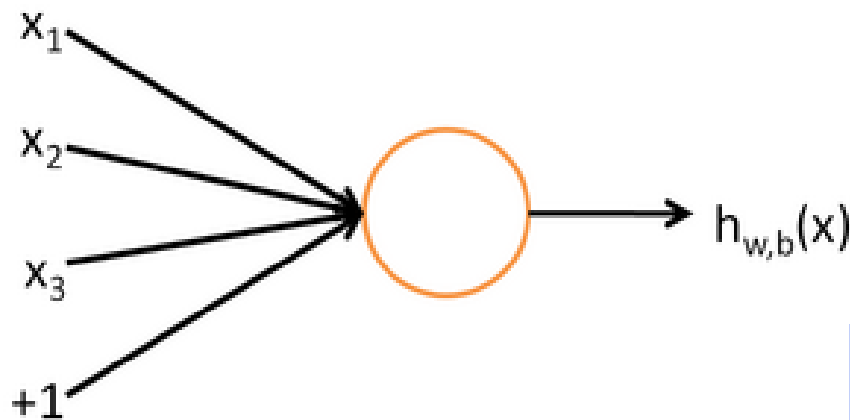
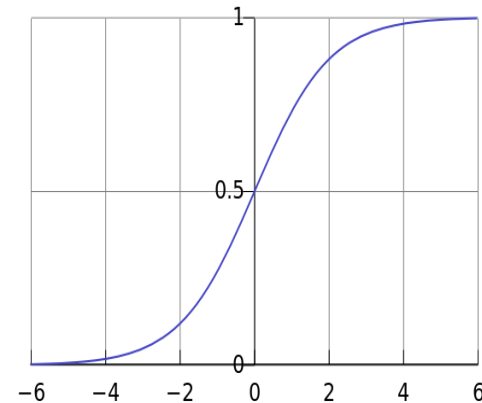
# A neuron is essentially a binary logistic regression

F = nonlinear activation fct. (e.g. sigmoid),  $w$  = weights,  $b$  = bias,  $h$  = hidden,  $x$  = inputs

$$h_{w,b}(x) = f(w^T x + b)$$

$b$ : We can have an “always on” Feature, which gives a class prior, or separate it out, as a bias term

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



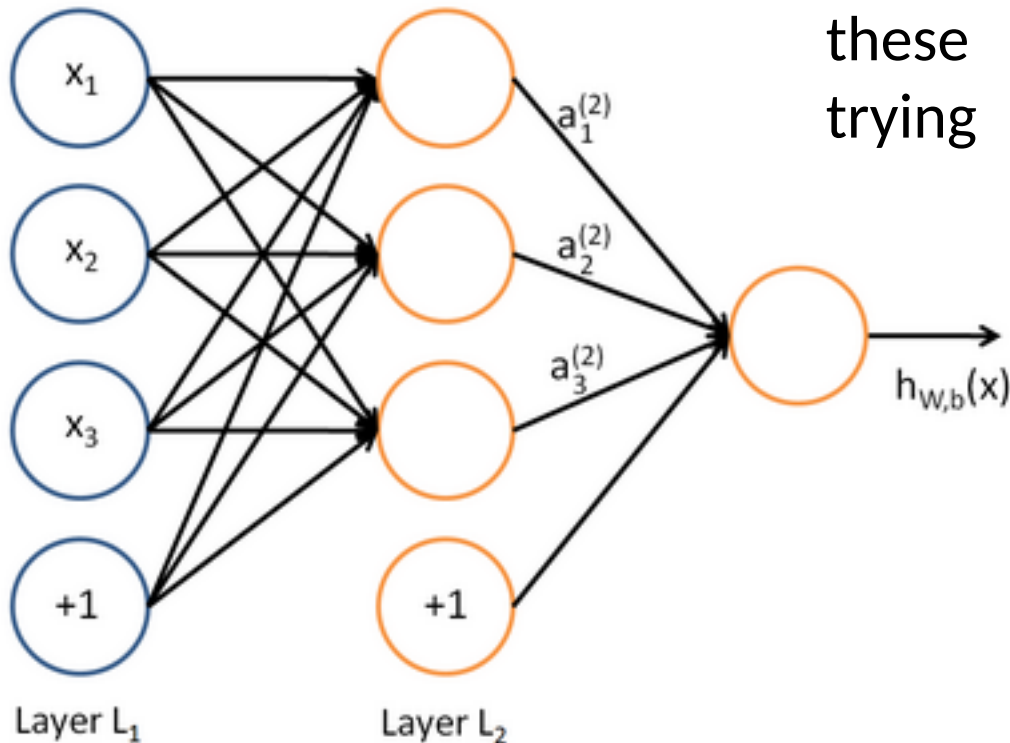
$w, b$  are the parameters of this neuron i.e., this logistic regression model

# A neural network

= running several logistic regressions at the same time

If we feed a vector of inputs through a bunch of logistic regression functions, then we get a vector of outputs ...

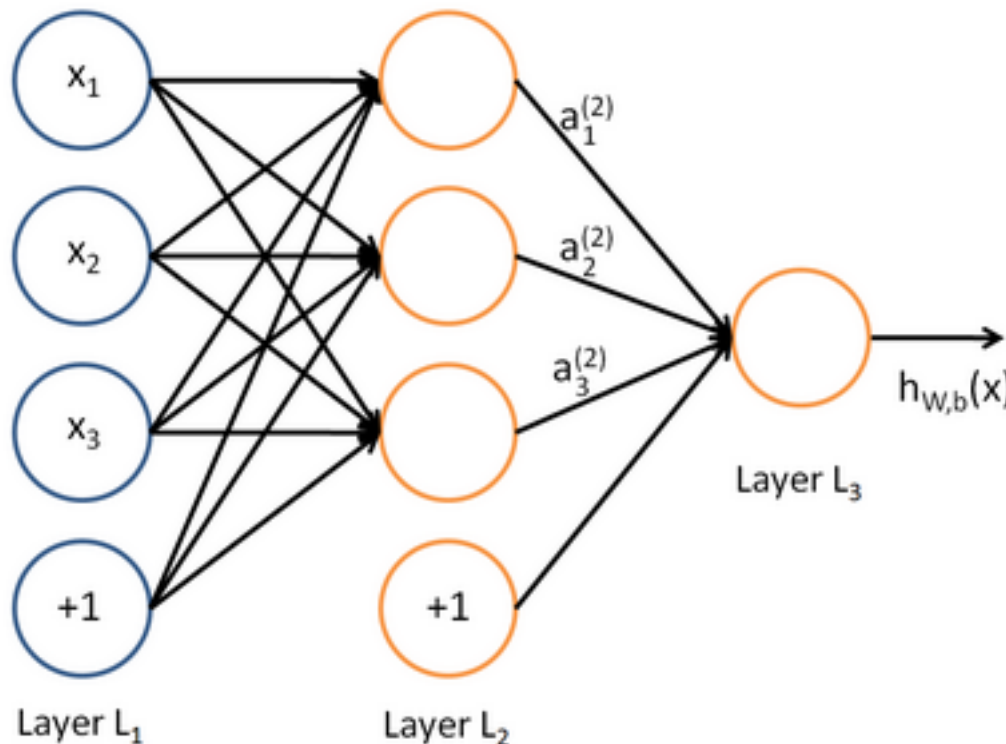
But we don't have to decide ahead of time what variables these logistic regressions are trying to predict!



# A neural network

= running several logistic regressions at the same time

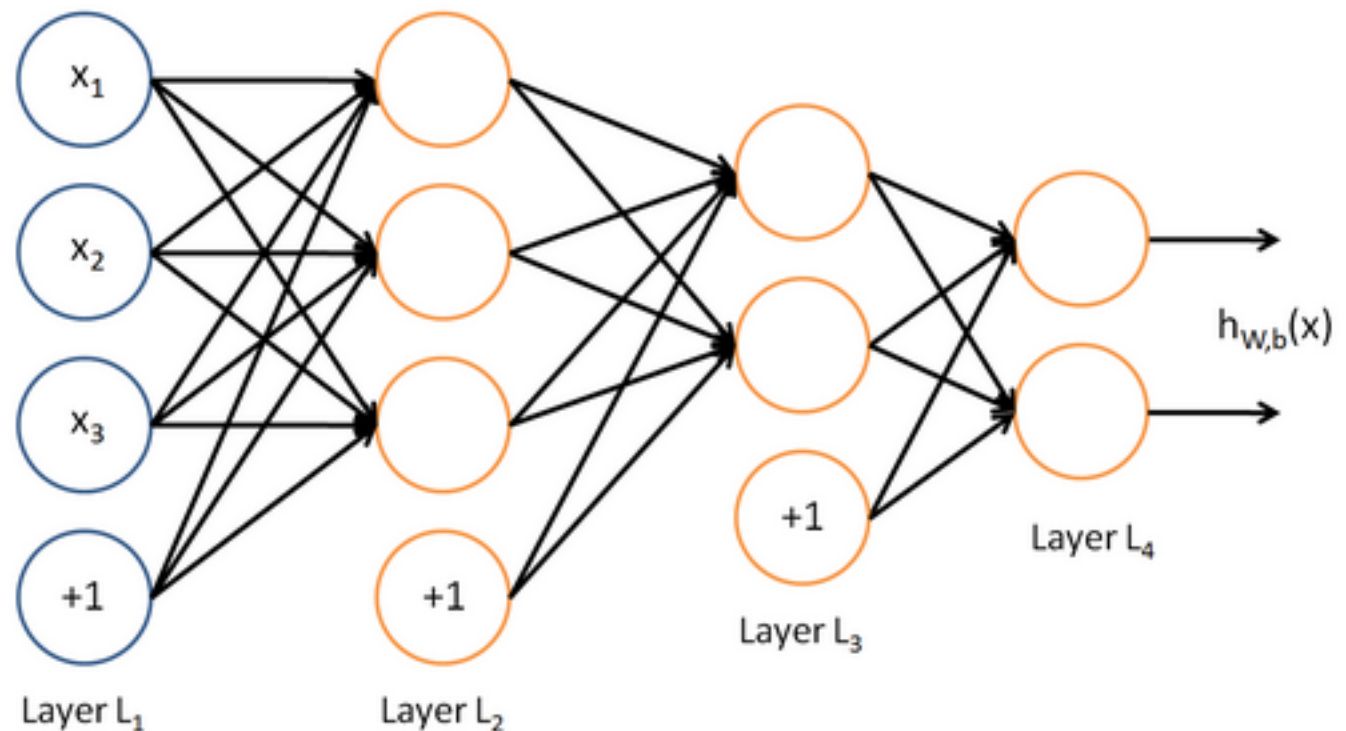
... which we can feed into another logistic regression function



It is the loss function that will direct what the intermediate hidden variables should be, so as to do a good job at predicting the targets for the next layer, etc.

**A neural network**  
**= running several logistic regressions at the same time**

Before we know it, we have a multilayer neural network....



# Matrix notation for a layer

We have

$$a_1 = f(W_{11}x_1 + W_{12}x_2 + W_{13}x_3 + b_1)$$

$$a_2 = f(W_{21}x_1 + W_{22}x_2 + W_{23}x_3 + b_2)$$

etc.

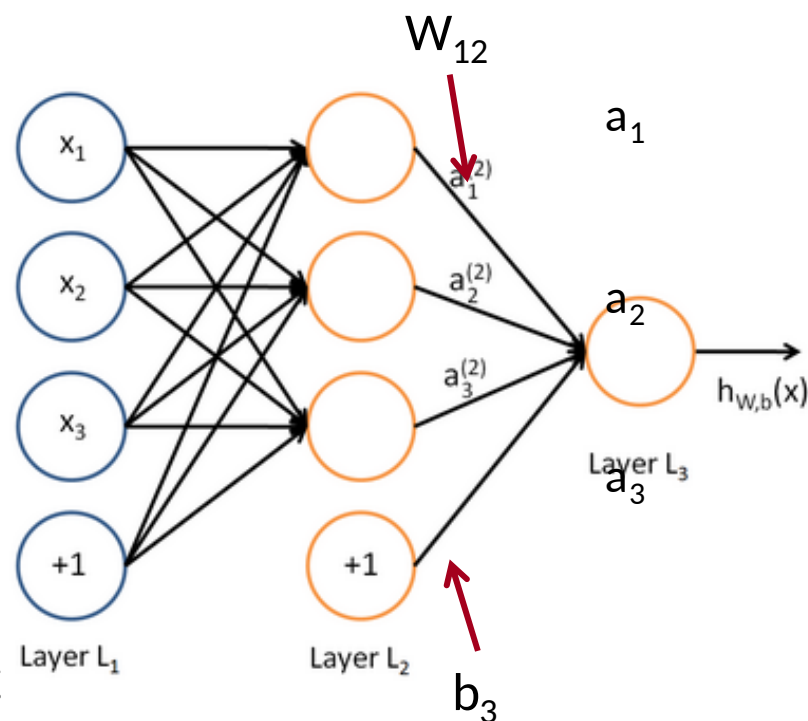
In matrix notation

$$z = Wx + b$$

$$a = f(z)$$

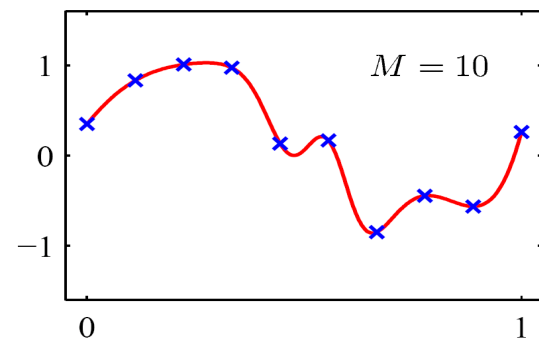
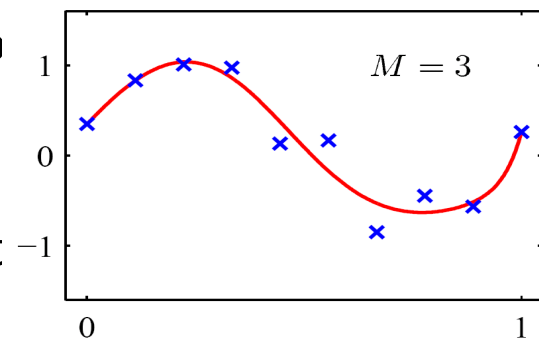
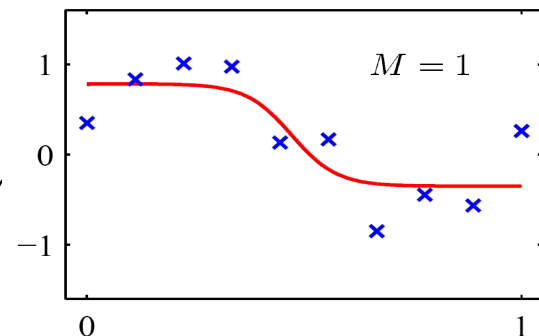
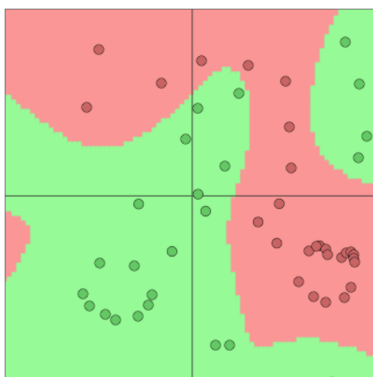
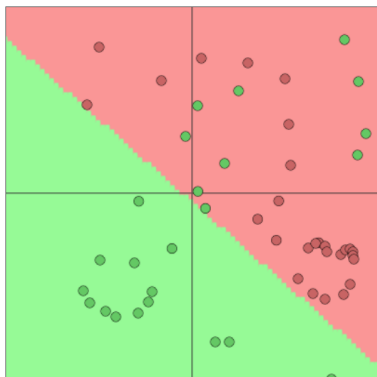
where  $f$  is applied element-wise:

$$f([z_1, z_2, z_3]) = [f(z_1), f(z_2), f(z_3)]$$



# Non-linearities (aka “f”): Why they’re needed

- Example: function approximation, e.g., regression or classification
  - Without non-linearities, deep neural network can’t do anything more than a linear transform
  - Extra layers could just be compiled into a single linear transform:
$$W_1 W_2 x = Wx$$
  - With more layers, they can approximate complex functions!





# Binary classification with unnormalized scores

- Revisiting our previous example:

$$X_{\text{window}} = \begin{bmatrix} x_{\text{museums}} & x_{\text{in}} & x_{\text{Paris}} & x_{\text{are}} & x_{\text{amazing}} \end{bmatrix}$$

- Assume we want to classify whether the center word is a Location (Named Entity Recognition)
- Similar to word2vec, we will go over all positions in a Corpus. But this time, it will be supervised and only some positions should get a high score.
- The positions that have an actual NER location in their center are called “true” positions.

# Binary classification for NER

- Example: Not all museums in Paris are amazing.
- Here: one true window, the one with Paris in its center and all other windows are “corrupt” in terms of not having a named entity location in their center.
- “Corrupt” windows are easy to find and there are Many: Any window that isn’t specifically labeled as NER location in our corpus

# A Single Layer Neural Network

- A single layer is a combination of a **linear layer** and **nonlinearity**:

$$z = Wx + b$$

$$a = f(z)$$

- The neural **activations**  $a$  can then be used to compute some output.

- For instance, a probability via softmax:

$$p(y|x) = \text{softmax}(Wa)$$

- Or an unnormalized score (even simpler):

$$\text{score}(x) = U^T a \in \mathbb{R}$$

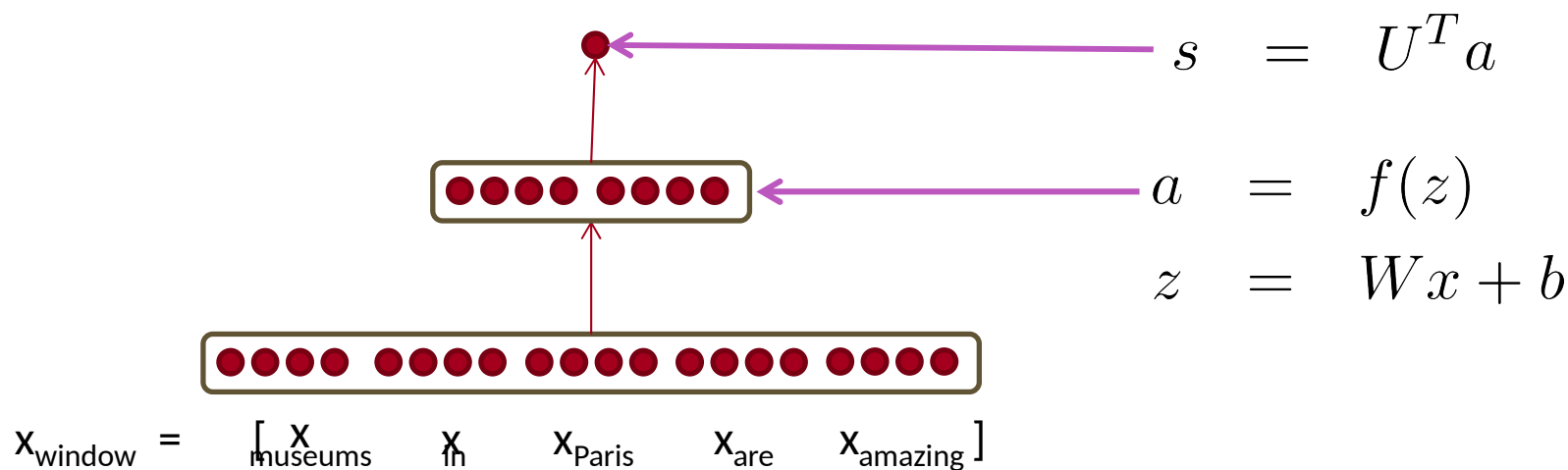
# Summary: Feed-forward Computation

We compute a window's score with a 3-layer neural

- $s = \text{score}(\text{"museums in Paris are amazing"})$

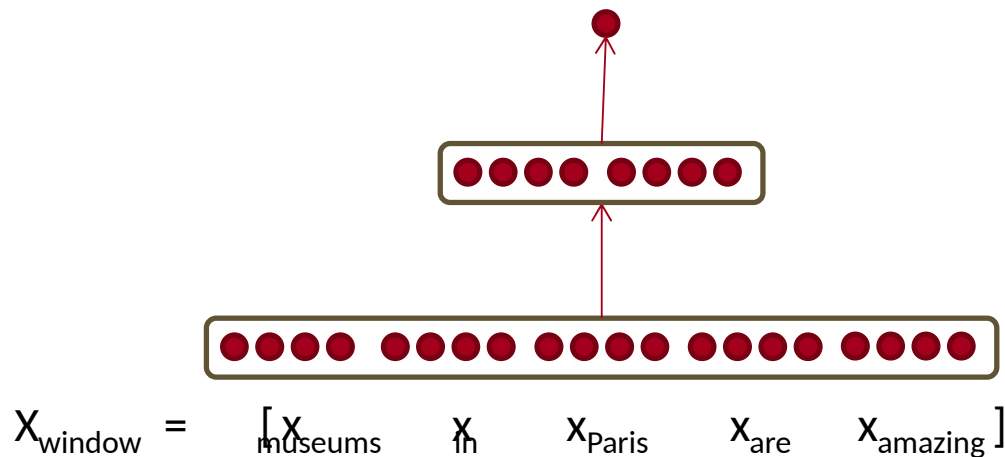
$$s = U^T f(Wx + b)$$

$$x \in \mathbb{R}^{20 \times 1}, W \in \mathbb{R}^{8 \times 20}, U \in \mathbb{R}^{8 \times 1}$$



# Main intuition for extra layer

The layer learns **non-linear interactions** between the input word vectors.



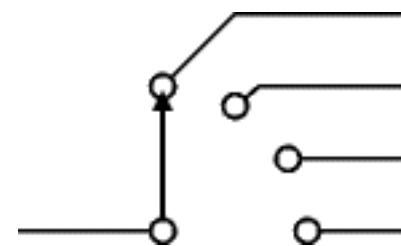
Example: only if “museums” is first vector should it matter that “in” is in the second position

# The max-margin loss

- Idea for training objective: Make true window's score larger and corrupt window's score lower (until they're good enough): minimize
- $s = \text{score}(\text{museums in Paris are amazing})$
- $s_c = \text{score}(\text{Not all museums in Paris})$

$$J = \max(0, 1 - s + s_c)$$

- This is not differentiable but it is continuous --> we can use SGD.



# Max-margin loss

- Objective for a single window:

$$J = \max(0, 1 - s + s_c)$$

- Each window with an NER location at its center should have a score +1 higher than any window without a location at its center
- For full objective function: Sample several corrupt windows per true one. Sum over all training windows.
- Similar to negative sampling in word2vec

# Deriving gradients for backprop

$$J = \max(0, 1 - s + s_c)$$

$$s = U^T f(Wx + b)$$

$$s_c = U^T f(Wx_c + b)$$

Assuming cost  $J$  is  $> 0$ ,  
compute the derivatives of  $s$  and  $s_c$   
involved variables:  $U, W, b, x$

$$\frac{\partial s}{\partial U} = \frac{\partial}{\partial U} U^T a \quad \Rightarrow \quad \frac{\partial s}{\partial U} = a$$



# Deriving gradients for backprop

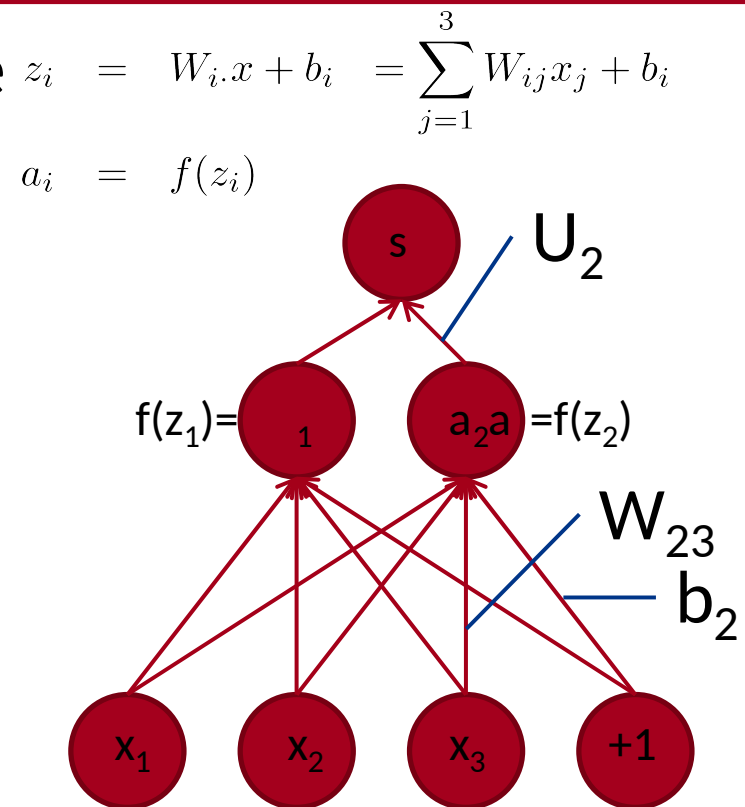
- For this function:

$$\frac{\partial s}{\partial W} = \frac{\partial}{\partial W} U^T a = \frac{\partial}{\partial W} U^T f(z) = \frac{\partial}{\partial W} U^T f(Wx + b)$$

- Let's consider the derivative of a single weight  $W_{ij}$

- $W_{ij}$  only appears inside  $a_i$

- For example:  $W_{23}$  is only used to compute  $a_2$  not  $a_1$



# Deriving gradients for backprop

Derivative of single weight<sub>ij</sub>:  $W$

$$\begin{aligned}
 \frac{\partial}{\partial W_{ij}} U^T a &\rightarrow \frac{\partial}{\partial W_{ij}} U_i a_i \\
 U_i \frac{\partial}{\partial W_{ij}} a_i &= U_i \frac{\partial a_i}{\partial z_i} \frac{\partial z_i}{\partial W_{ij}} \\
 &= U_i \frac{\partial f(z_i)}{\partial z_i} \frac{\partial z_i}{\partial W_{ij}} \\
 &= U_i f'(z_i) \frac{\partial z_i}{\partial W_{ij}} \\
 &= U_i f'(z_i) \frac{\partial W_{i \cdot} x + b_i}{\partial W_{ij}}
 \end{aligned}$$

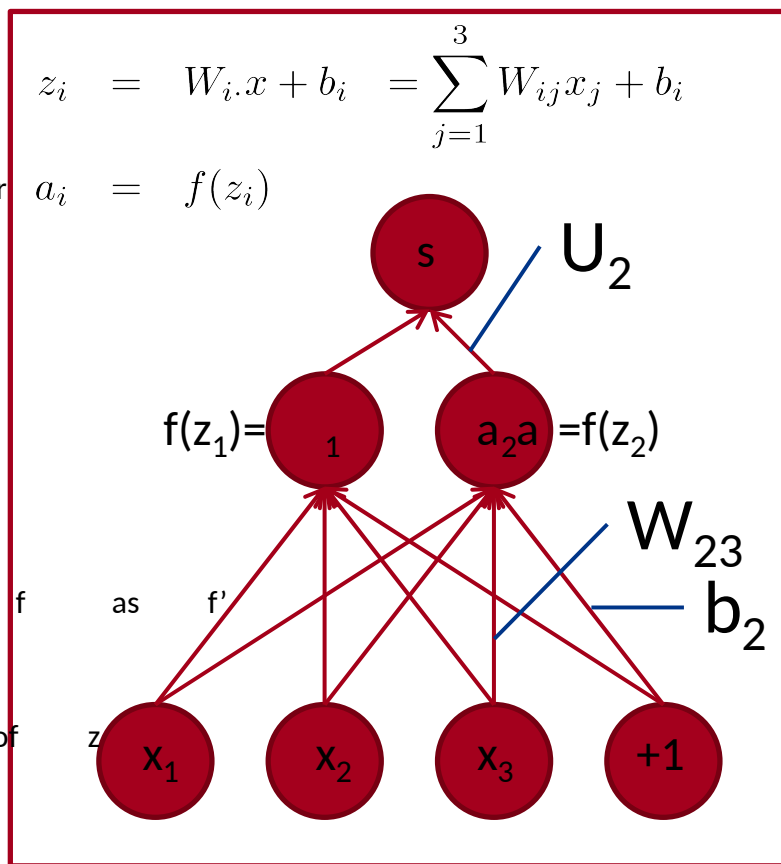
Ignore constants in which  $W$  doesn't appear

Pull out since it's constant, apply chain rule

Apply definition of  $a$

Just define derivative of  $f$  as  $f'$

Plug in definition of  $z$

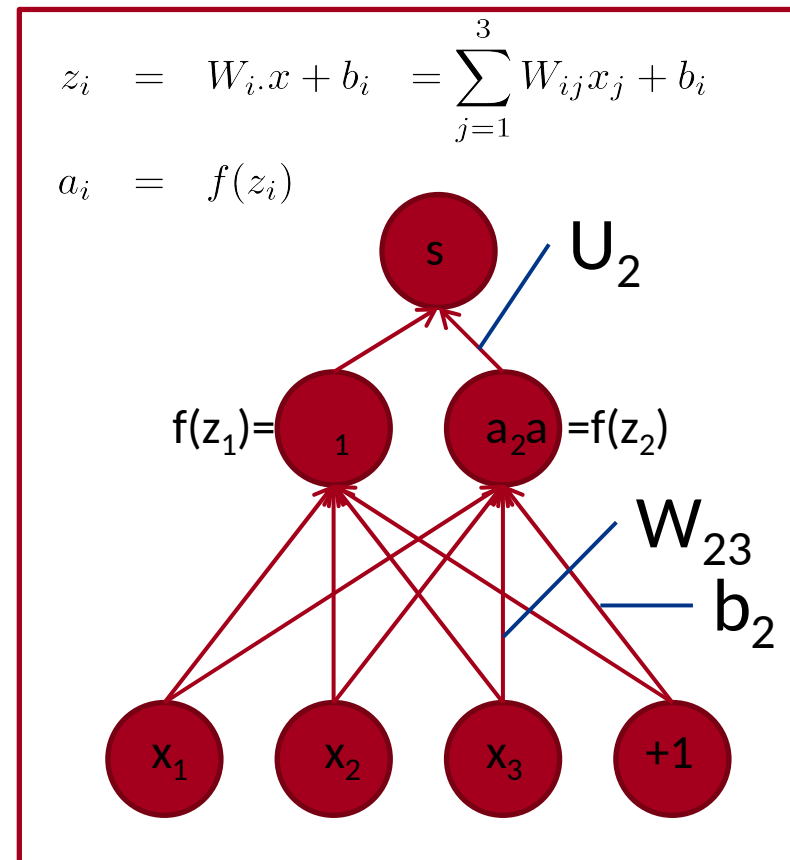


# Deriving gradients for backprop

Derivative of single weight  $w_{ij}$  continued:

$$\begin{aligned}
 U_i \frac{\partial}{\partial W_{ij}} a_i &= U_i f'(z_i) \frac{\partial W_{i \cdot} x + b_i}{\partial W_{ij}} \\
 &= U_i f'(z_i) \frac{\partial}{\partial W_{ij}} \sum_k W_{ik} x_k \\
 &= \underbrace{U_i f'(z_i)}_{\delta_i} x_j \\
 &= \underbrace{\delta_i}_{\text{Local signal}} \underbrace{x_j}_{\text{error Local signal}} \underbrace{\text{input}}
 \end{aligned}$$

where  $f'(z) = f(z)(1 - f(z))$



# Deriving gradients for backprop

- So far, derivative of single  $y_j$  only, but we want gradient for full  $W$ .

$$\begin{aligned}\frac{\partial s}{\partial W_{ij}} &= \underbrace{U_i f'(z_i)}_{\delta_i} x_j \\ &= \delta_i x_j\end{aligned}$$

- We want all combinations of  $i = 1, 2$  and  $j = 1, 2$ ,

- Solution: Outer product:  $\frac{\partial s}{\partial W} = \delta x^T$

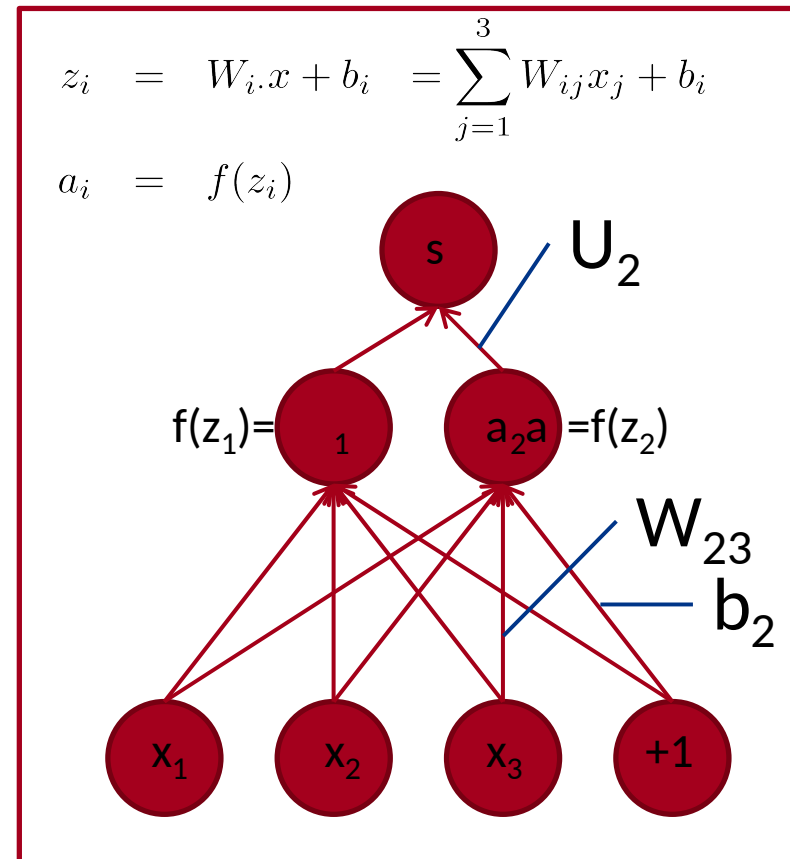
where  $\delta \in \mathbb{R}^{2 \times 1}$  is the “responsibility” or error signal coming from each activation  $a$

# Deriving gradients for backprop

- How to derive gradients for biases  $b$ ?

$$\frac{\partial s}{\partial b} = a$$

$$\begin{aligned} & U_i \frac{\partial}{\partial b_i} a_i \\ &= U_i f'(z_i) \frac{\partial W_i \cdot x + b_i}{\partial b_i} \\ &= \delta_i \end{aligned}$$



# Training with Backpropagation

That's almost backpropagation

It's taking derivatives and using the chain rule

Remaining trick: we can **re-use** derivatives computed for Higher layers in computing derivatives for lower layers!

Example: last derivatives of model, the word vectors in  $x$

# Training with Backpropagation

- Take derivative of score with respect to single element of word vector
- Now, we cannot just take into consideration one  $a$  because each is connected to all the neurons above and hence  $x$  influences the overall score through all of these, hence:

$$\begin{aligned}
 \frac{\partial s}{\partial x_j} &= \sum_{i=1}^2 \frac{\partial s}{\partial a_i} \frac{\partial a_i}{\partial x_j} \\
 &= \sum_{i=1}^2 \frac{\partial U^T a}{\partial a_i} \frac{\partial a_i}{\partial x_j} \\
 &= \sum_{i=1}^2 U_i \frac{\partial f(W_{i \cdot} x + b)}{\partial x_j} \\
 &= \sum_{i=1}^2 \underbrace{U_i f'(W_{i \cdot} x + b)} \frac{\partial W_{i \cdot} x}{\partial x_j} \\
 &= \sum_{i=1}^2 \delta_i W_{ij} \\
 &= W_{\cdot j}^T \delta'_{\cdot j}
 \end{aligned}$$

Re-used part of previous derivative

# Training with Backpropagation

- With  $\frac{\partial s}{\partial x_j} = W_{\cdot j}^T \delta$

$$\frac{\partial s}{\partial x} = W^T \delta$$

- Observations: The error message that arrives at a hidden layer has the same dimensionality as that hidden layer



## Putting all gradients together:

- Remember: Full objective function for each window was:

$$J = \max(0, 1 - s + s_c) \quad \begin{aligned} s &= U^T f(Wx + b) \\ s_c &= U^T f(Wx_c + b) \end{aligned}$$

- For example: gradient for just  $U$ :

$$\frac{\partial J}{\partial U} = 1\{1 - s + s_c > 0\} (-f(Wx + b) + f(Wx_c + b))$$

$$\frac{\partial J}{\partial U} = 1\{1 - s + s_c > 0\} (-a + a_c)$$

# Summary

Congrats! Super useful basic components and real model

- Word vectortraining
- Windows
- Softmax andcross entropy error
- Scores and max-margin loss
- Neural network