

- ReLU

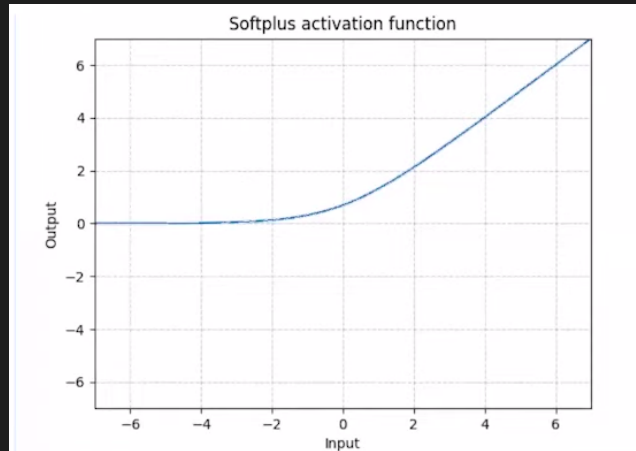
↳ Variations : Leaky, etc.

↳ 1 kink : equivariant to scalar product (scale invariant)  
esquisse, corner

- Softplus

↙ Large  $\beta \approx \text{ReLU}$

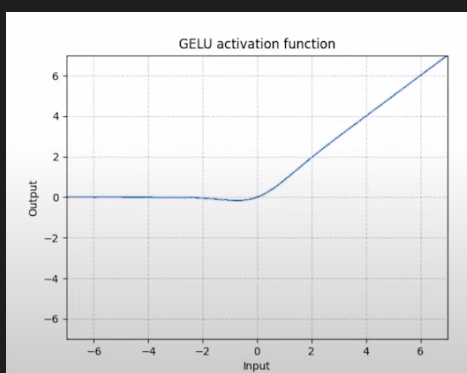
$$\text{Softplus}(x) = \frac{1}{\beta} * \log(1 + \exp(\beta * x))$$



- ELU

- CELU

- GELU : non-monotonic



$$\text{GELU}(x) = x * \Phi(x)$$

where  $\Phi(x)$  is the Cumulative Distribution Function for Gaussian Distribution.

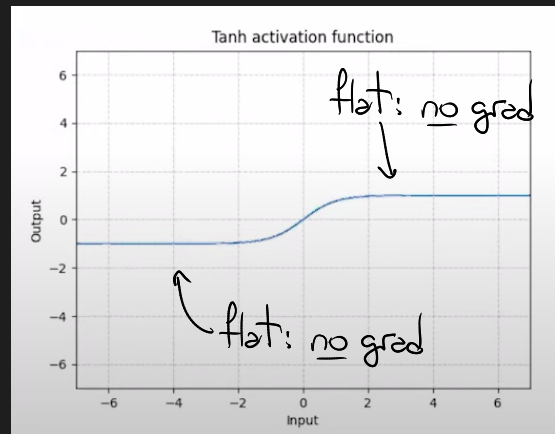
• ReLU 6  2 kinks

• Good ol' Sigmoid


• Tanh :


$$\frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Problem:



→ ReLUs are better for deeper networks (not clear why)

• Soft sign 

• Hard tanh, 

• Tanh shrink 

• Hard shrink 

• Log Sigmoid:  $\log(\text{sigmoid}(x))$

$\text{Softmax} : \text{"soft arg max"}$   
 $\text{Softmin} : \text{"soft arg min"}$

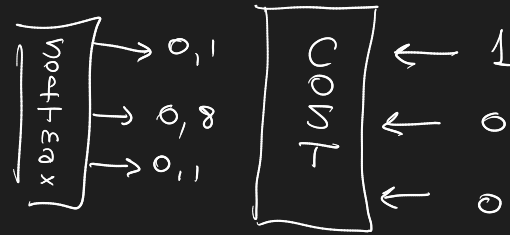
Invariant to  
 translations.

$\hookrightarrow \text{softmax}(-x)$

Cost Function for Classification.

Preds.  $\bar{y}$

Correct labels  $y$



• MSE

• Binary cross entropy

$\hookrightarrow$  Negative log likelihood (main way)

$\hookrightarrow -\log(\bar{y}_c)$

$\uparrow$   
 $\bar{y}$  of correct category.

Note: the others  $\bar{y}_i$  are not used.

$\hookrightarrow$  We need the  $\log \text{softmax}$  for this.

•  $L_1$  norm: not used very often

# Architectures

↳ Different ways of arranging modules to build Neural Networks

## Multiplicative Modules

- Quadratic modules

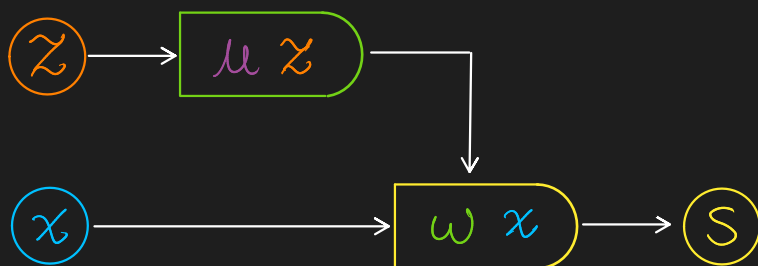
output  $\swarrow$   $\nwarrow$  weighted sum of inputs (linear function)

$$S_i = \sum_j w_{ij} x_j$$

with each  $w_{ij}$  given by

$$w_{ij} = \sum_k \mu_{ijk} z_k$$

$\nwarrow$  linear function of inputs  $z_k$



equiv:

$\nwarrow$  if  $z_k = x_j \Rightarrow$  "quadratic form"

$$S_i = \sum_{jk} \mu_{ijk} z_k x_j$$

$\nwarrow$  2<sup>nd</sup> degree monomial.

# • Attention module

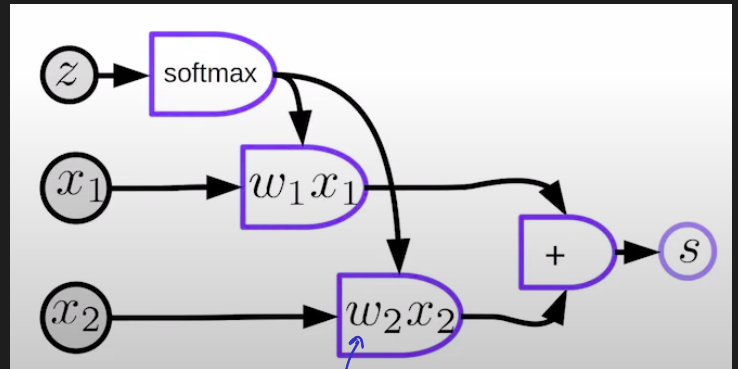
$$S_i = \sum_j w_j x_{ij}$$

with

$$w_j = \frac{e^{z_j}}{\sum_k e^{z_k}}$$

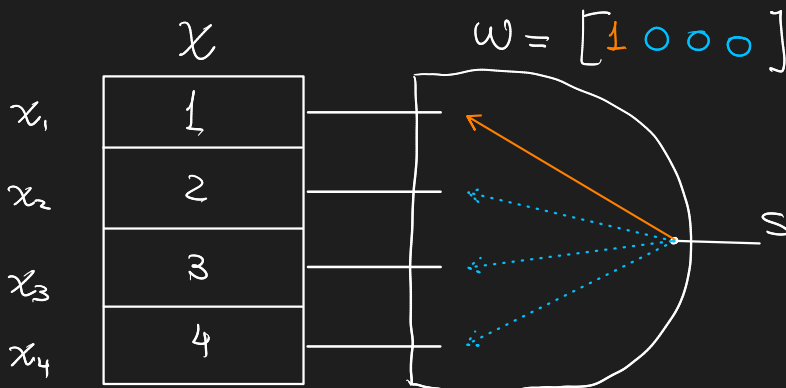
↖ pays attention to these part of the input

i	w	→	x	
0	0,7	→	1	← important
1	0	→	2	← <u>non</u> important
2	0,2	→	3	← a little
3	0,1	→	4	← less than a little



$w_i$  : scalar  
 $x_i$  : scalar

If  $w$  is 1-hot vector  $\Rightarrow$  it is a switch  
(a hard switch)

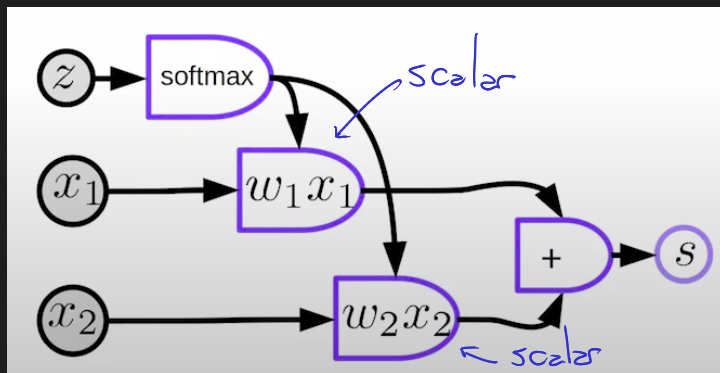


Backprop is trivial for switches,

↳ same gradient from the top goes through the only way possible.

no gradient through the non-selected inputs.

Soft switch



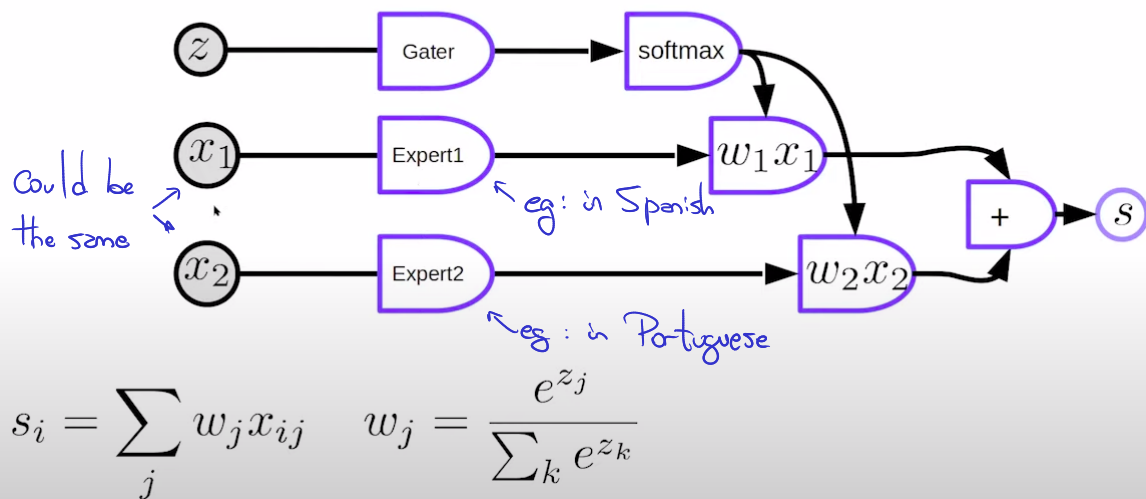
$$S = w_1 \cdot x_1 + w_2 \cdot x_2$$

$$\text{with } w_1 + w_2 = 1$$

$z$  = Attention the system pays to  $x_1$  and  $x_2$ .

# Mixture of Experts

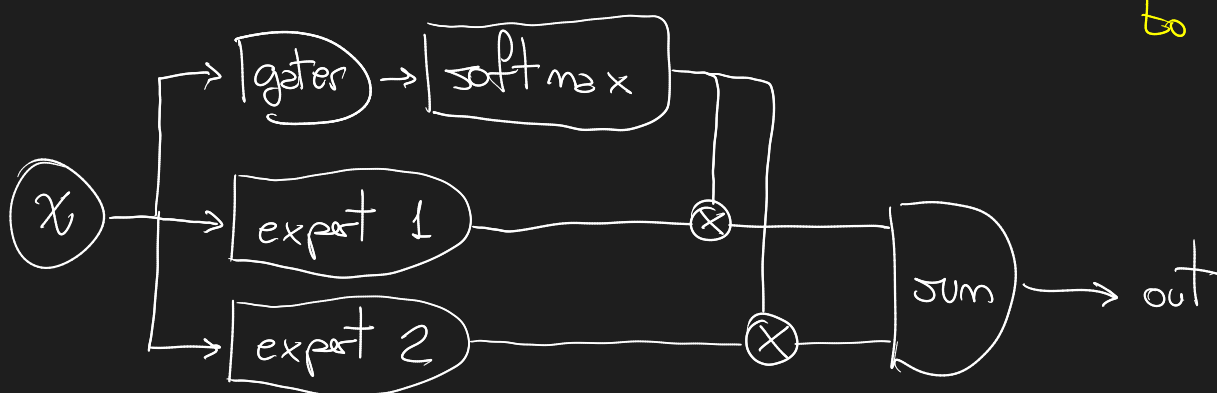
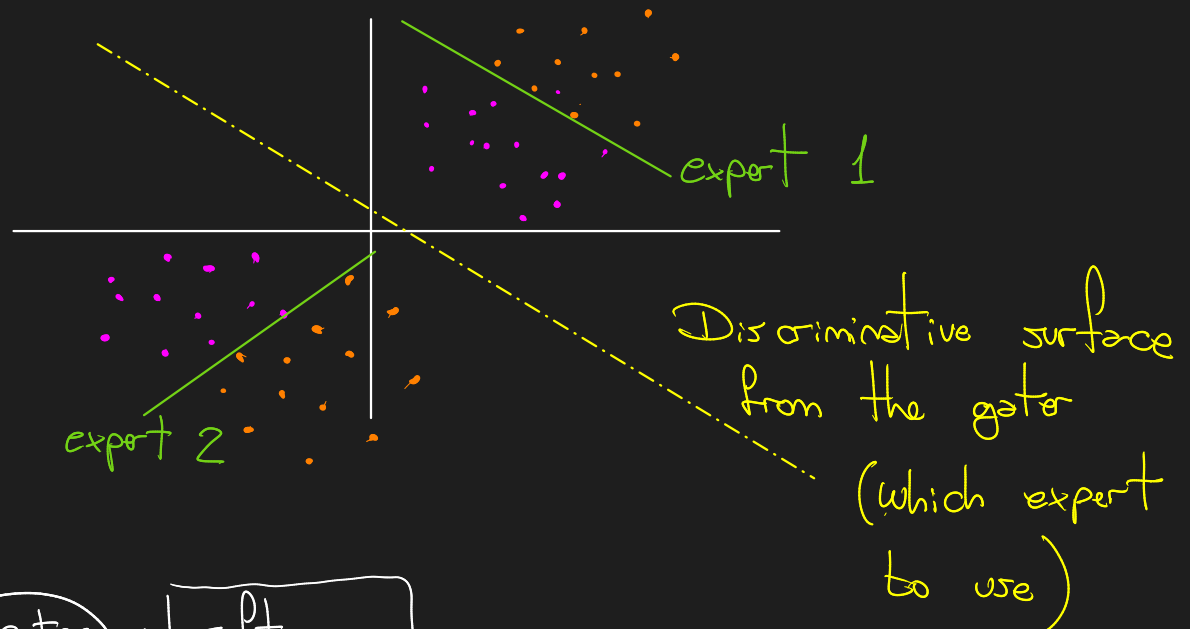
## ► Attention module “Switches” expert networks



If the input is in some mixed dialect (eg: Catalan), the attention module will switch between experts in a smart way.

Example: Linear classifiers to Nonlinear one

Imagine one of the experts could linearly classify a portion of the input space, and the other expert, classify another portion of the input space.

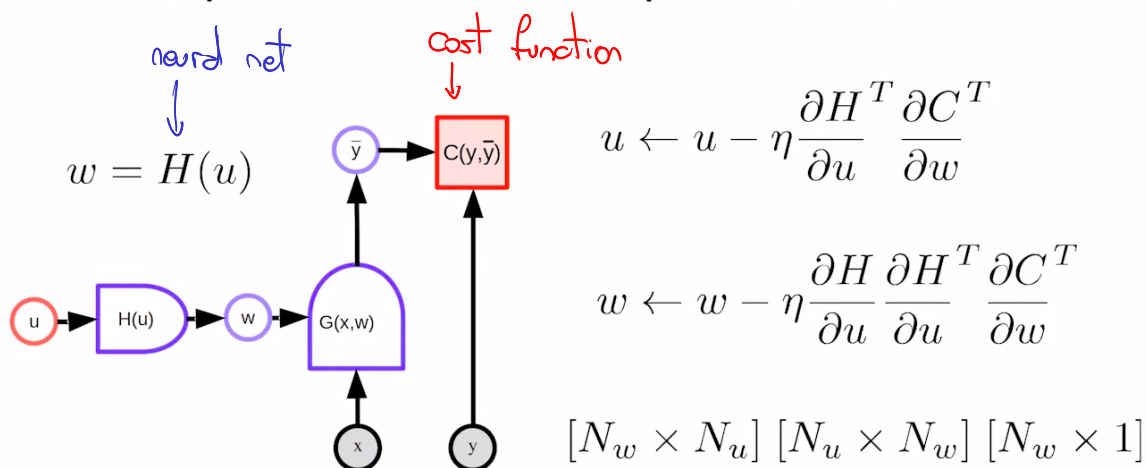


Training: Simple Backprop.

Generalization.

## Parameter transformations

- ▶ When the parameter vector is the output of a function



## Simple parameter transform: weight sharing

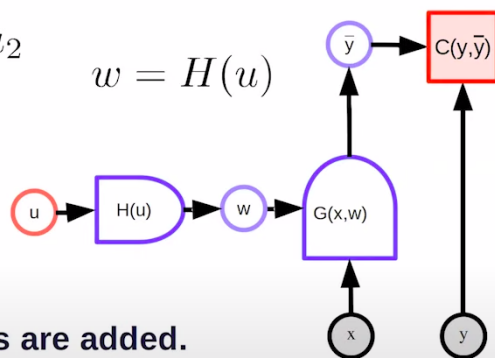
- ▶ Function  $H(u)$  replicates one component of  $u$  into multiple components of  $w$

$$w_1, = w_2 = u_1 \quad w_3 = w_4 = u_2$$

- ▶  $H$  is like a "Y" branch.

- ▶ Gradients are summed in the backprop

- ▶ The gradients w.r.t. shared parameters are added.



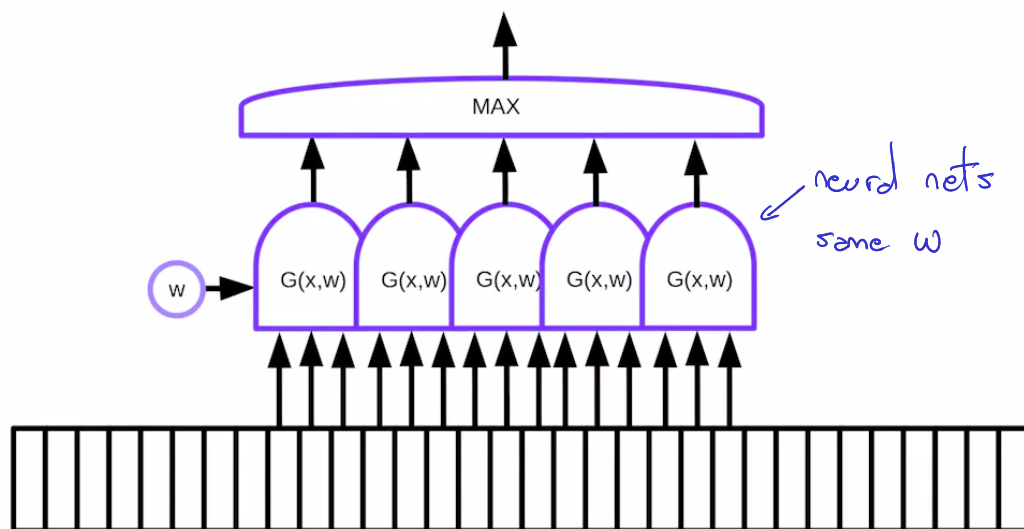


Special Case.

“face on image”; “alexz”/“hey google”

## Shared Weights for Motif Detection

► Detecting motifs anywhere on an input



Exercise :

- Backprop through softmax

$$\sigma_i(x) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$\text{Jacobian}(\sigma(x)) = \begin{bmatrix} \frac{\partial \sigma_1}{\partial x_1} & \frac{\partial \sigma_1}{\partial x_2} & \dots & \frac{\partial \sigma_1}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \sigma_n}{\partial x_1} & \frac{\partial \sigma_n}{\partial x_2} & \dots & \frac{\partial \sigma_n}{\partial x_n} \end{bmatrix}$$

First, we apply  $\log$  (all positive values)

$$\begin{aligned}\log(\sigma_i(x)) &= \log\left(\frac{e^{x_i}}{\sum_j e^{x_j}}\right) \\ &= \underbrace{\log e^{x_i}}_{= x_i} - \underbrace{\log\left(\sum_j e^{x_j}\right)}_{\text{same for all}}\end{aligned}$$

$$\log \sigma_i(x) = x_i - \log\left(\sum_k e^{x_k}\right)$$

Differentiating

$$\frac{\partial}{\partial x_j} \log \sigma_i(x) = \frac{\partial x_i}{\partial x_j} - \frac{\partial}{\partial x_j} \cdot \log\left(\sum_k e^{x_k}\right)$$

Where

$$\frac{\partial x_i}{\partial x_j} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

So

$$\frac{\partial x_i}{\partial x_j} = \mathbb{1}_{\{i=j\}}$$

The other part

$$\frac{\partial}{\partial x_j} \cdot \log \left( \sum_k e^{x_k} \right) = \frac{1}{\sum_k e^{x_k}} \cdot \frac{\partial}{\partial x_j} \left( \sum_k e^{x_k} \right)$$

where

$$\begin{aligned} \frac{\partial}{\partial x_j} \left( \sum_k e^{x_k} \right) &= \sum_k \frac{\partial e^{x_k}}{\partial x_j} \\ &= 0 + 0 + \frac{\partial e^{x_j}}{\partial x_j} + 0 + 0 + \dots \\ &= \frac{\partial e^{x_j}}{\partial x_j} \\ &= e^{x_j} \end{aligned}$$

All together

$$\frac{\partial}{\partial x_j} \log \sigma_i(x) = \frac{\partial x_i}{\partial x_j} - \frac{\partial}{\partial x_j} \cdot \log \left( \sum_k e^{x_k} \right)$$

$$= \mathbb{1}\{i=j\} - \underbrace{\frac{1}{\sum_k e^{x_k}}}_{\sigma_j(x)} \cdot e^{x_j}$$

$$= \mathbb{1}\{i=j\} - \sigma_j(x)$$

Multiplying both sides by  $\sigma_i(x)$ , we get ①

$$\sigma_i(x) \left( \frac{\partial}{\partial x_j} \log \sigma_i(x) \right) = \sigma_i(x) \cdot \left( \mathbb{1}\{i=j\} - \sigma_j(x) \right)$$

If we take partial derivative of  $\log(\text{softmax})$

$$\frac{\partial}{\partial x_j} \log(\sigma_i(x)) = \underbrace{\frac{1}{\sigma_i(x)}}_{\text{deriv. of log}} \cdot \frac{\partial \sigma_i(x)}{\partial x_j}$$

$$\frac{\partial \sigma_i(x)}{\partial x_j} = \sigma_i(x) \cdot \frac{\partial}{\partial x_j} \log(\sigma_i(x))$$

①

↓

$$= \sigma_i(x) \cdot \left( \mathbb{1}\{i=j\} - \sigma_j(x) \right)$$

To conclude, the Jac. matrix becomes:

$$J(\sigma(x)) = \begin{bmatrix} \sigma_1(1-\sigma_1) & \sigma_1(-\sigma_2) & 0 & 0 & 0 & \sigma_1(-\sigma_n) \\ \sigma_2(-\sigma_1) & \sigma_2(1-\sigma_2) & 0 & 0 & 0 & \sigma_2(-\sigma_n) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \sigma_n(-\sigma_1) & \sigma_n(-\sigma_1) & 0 & 0 & 0 & \sigma_n(1-\sigma_n) \end{bmatrix}$$

