# Library for Deep Learning

- PyTorch
  - Facebook's AI Research Lab (FAIR)
- TensorFlow
  - Google Brain Team



https://medium.com/ml-cheat-sheet/machine-learning-curriculum-from-rookie-to-mastery-636bdbf74856

## **Activation Functions**

- determines if a neuron should be activated or not activated
  - use some simple mathematical operations to determine if the neuron's input to the network is relevant or not relevant in the prediction process
- introduce non-linearity to an artificial neural network
- generate output from a collection of input values fed to a layer

	Human brain	Neural networks
Input	Receive from outside world	Provide as images, sounds, numbers, etc.
Process	Neuron	Artificial neuron
Activate	Neuron tail	Algorithm

- universal approximation: a neural network can approximate any continuous function that maps inputs (X) to outputs (y)
  - The ability to represent any function is what makes the neural networks so powerful and widely-used
  - To be able to approximate any function, we need non-linearity: activation functions
  - Without activation functions, neural networks can be considered as a collection of linear models



$$y = f(\sum(weight * input) + bias)$$

# **Activation Functions**

- Linear activation function
  - proportional to the input,  $(-\infty, \infty)$
  - simply adds up the weighted total of the inputs and returns the result
  - Pros and Cons
    - not a binary activation because it only delivers a range of activations
    - Constant derivative: unrelated to x



- Binary step function
  - A threshold value determines whether a neuron should be activated or not activated
  - Pros and Cons
    - cannot provide multi-value outputs
    - gradient is zero



Mechanistic Data Science

f(x) = x

# Non-linear Activation Functions: Sigmoid

- simple to use, all the desirable qualities: nonlinearity, continuous differentiation, monotonicity, and a set output range, mainly used in binary classification problems
- Pros and Cons
  - non-linear in nature, smooth gradient, good for a classifier type problem
  - output is always going to be in the range  $(0,1) \rightarrow$  define a range for our activations
  - problem of "Vanishing gradients"
  - Its output isn't zero centered  $\rightarrow$  gradient updates go too far in different directions
  - output value is between zero and one → optimization harder (network either refuses to learn more or is extremely slow)



# Non-linear Activation Functions: Tanh (Hyperbolic Tangent)

- real-valued number to the range [-1, 1]
- non-linear, its output is zero-centered (different from Sigmoid)
- negative inputs will be mapped strongly to the negative and zero inputs will be mapped to almost zero in the graph of TanH
- Pros and Cons
  - vanishing gradient problem, but the gradient is stronger for TanH than sigmoid (derivatives are steeper)



Mechanistic Data Science

Deep Learning - 41

# Non-linear Activation Functions: ReLU (Rectified Linear Unit)

- one of the most commonly used in the applications
- maximum value of the gradient is one  $\rightarrow$  solve the problem of vanishing gradient:
- slope is never zero  $\rightarrow$  solved the problem of saturating neuron
- Pros and Cons
  - only a certain number of neurons are activated → far more computationally efficient
  - accelerates the convergence of gradient descent towards the global minimum of the loss function due to its linear, non-saturating property

0.5 -

0.0

-5.0

-2.5

2.5

5.0

0.0

- It should only be used within hidden layers
- Some gradients can be fragile during training  $\rightarrow$  dying ReLu problem

$$f(x) = max (0, x)$$



Deep Learning - 42

# Non-linear Activation Functions: Leaky ReLU

- upgraded version of the ReLU to solve the dying ReLU problem
- consistency of the benefit across tasks is presently ambiguous
- Pros and Cons

f(x)

- ReLU + enable back propagation, even for negative input values
- Making minor modification of negative input values, the gradient of the left side of the graph comes out to be a real (non-zero) value. As a result, there would be no more dead neurons in that area
- predictions may not be steady for negative input values

$$= max (0.1x, x)$$

Deep Learning - 43

Mechanistic Data Science

## Non-linear Activation Functions: ELU (Exponential Linear Units)

- one of the variations of ReLU which also solves the dead ReLU problem
- very similar to ReLU except negative inputs, slightly more computationally expensive than leaky ReLU
- Pros and Cons
  - strong alternative to ReLU: can produce negative outputs
  - Exponential operations  $\rightarrow$  increase the computational time
  - No learning about the 'a' value takes place, and exploding gradient problem

$$\begin{cases} x & for \ x \ge 0 \\ \alpha(e^x - 1) & for \ x < 0 \\ & & \\ -1 & & \\ \end{cases}$$

-2.5

-5.0

0.0

2.5

5.0

Mechanistic Data Science

Deep Learning - 44

# Non-linear Activation Functions: Softmax

- A combination of many sigmoids, determines relative probability, returns the probability of each class/labels, most commonly used in multi-class classification
- gives the probability of the current class with respect to others  $\rightarrow$  consider the possibility of other classes

2.5

5.0

Deep Learning - 45

0.0

-2.5

-5.0

- Pros and Cons
  - It mimics the one encoded label better than the absolute values
  - We would lose information if we used absolute (modulus) values, but the exponential takes care of this on its own
  - should be used for multi-label classification and regression task as well

Softmax
$$(x_i) = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$



## Non-linear Activation Functions: Swish

- allows for the propagation of a few numbers of negative weights
- crucial property that determines the success of non-monotonic smooth functions
- self-gated activation function created by Google researchers
- Pros and Cons
  - smooth activation function: it does not suddenly change direction like ReLU does near x equal to zero
  - large negative numbers are wiped out, resulting in a win-win situation
  - non-monotonous function enhances the term of input data and weight to be learnt
  - Slightly more computationally expensive

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



# **Important Considerations**

- Vanishing gradient is a common problem encountered during neural network training
  - sigmoid activation function: small output range (0 to 1), huge change in the input
     → small modification in the output, derivative also becomes small
  - These are only used for shallow networks with only a few layers, For a multi-layer network, the gradient may become too small for expected training
- Exploding gradients are situations in which massive incorrect gradients build during training, resulting in huge updates to neural network model weights
  - An unstable network might form, and training cannot be completed
  - weights' values can potentially grow to the point where they overflow, resulting in loss in NaN values

# Final Takeaways

- All hidden layers generally use the same activation functions. ReLU activation function should only be used in the hidden layer for better results.
- Sigmoid and TanH activation functions should not be utilized in hidden layers due to the vanishing gradient, since they make the model more susceptible to problems during training.
- Swish function is used in artificial neural networks having a depth more than 40 layers.
- Regression problems should use linear activation functions
- Binary classification problems should use the sigmoid activation function
- Multiclass classification problems should use the softmax activation function
- Neural network architecture and their usable activation functions,
  - Convolutional Neural Network (CNN): ReLU activation function
  - Recurrent Neural Network (RNN): TanH or sigmoid activation functions

# **Computing Gradients**

$$\nabla f = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n}\right)$$

$$F(x) = F_3\left(F_2\left(F_1(x)\right)\right) \rightarrow \frac{dF}{dx} = \frac{dF_3}{dx}\left(F_2\left(F_1(x)\right)\right) \frac{dF_2}{dx}(F_1(x)) \frac{dF_1}{dx}(x)$$

$$derivatives \begin{cases} sum rule: \frac{\partial}{\partial a}(a+b) = \frac{\partial a}{\partial a} + \frac{\partial b}{\partial a} = 1 \\ product rule: \frac{\partial}{\partial u}(uv) = \frac{\partial u}{\partial u}v + u \frac{\partial v}{\partial u} = v \end{cases}$$

$$[computational graph]$$

$$e = (a+b)(b+1) \rightarrow \begin{cases} c = a+b \\ d = b+1 \\ e = cd \end{cases}$$

$$\left\{ \frac{\partial e}{\partial a} = \frac{\partial e}{\partial c} \frac{\partial c}{\partial a} = \\ \frac{\partial e}{\partial b} = \frac{\partial e}{\partial c} \frac{\partial c}{\partial b} + \frac{\partial e}{\partial d} \frac{\partial d}{\partial b} = \end{array} \right\}$$

Mechanistic Data Science

#### Multivariate Chain Rule



Deep Learning: one loss function with many inputs (weights)  $\rightarrow$  backward!



# Patterns in Gradient Flow

- Add gate: gradient distributor
- Mul gate: swap multiplier
- Copy gate: gradient adder
- Max gate: gradient router



## Example



#### **Vector Derivatives**

Scalar to Scalar 
$$(x \in \mathbb{R}, y \in \mathbb{R}) \rightarrow \frac{\partial y}{\partial x} \in \mathbb{R}$$
  
Vector to Scalar  $(\mathbf{x} \in \mathbb{R}^N, y \in \mathbb{R}) \rightarrow \frac{\partial y}{\partial \mathbf{x}} \in \mathbb{R}^N, \left(\frac{\partial y}{\partial \mathbf{x}}\right)_n = \frac{\partial y}{\partial x_n}$ : gradient  
Vector to Vector  $(\mathbf{x} \in \mathbb{R}^N, \mathbf{y} \in \mathbb{R}^M) \rightarrow \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \in \mathbb{R}^{N \times M}, \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}}\right)_{n,m} = \frac{\partial y_m}{\partial x_n}$ : Jacobian



#### **Back Propagation with Matrices**



What part of **y** are affected by one element of **x**?  $x_{n,d}$  affects the whole row  $y_m$ How much does  $x_{n,d}$  affect  $y_{n,m}$ ?  $w_{d,m}$ 

$$\frac{\partial L}{\partial x_{n,d}} = \sum_{m} \frac{\partial L}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}} = \sum_{m} \frac{\partial L}{\partial y_{n,m}} w_{d,m} \rightarrow \begin{cases} \frac{\partial L}{\partial \mathbf{x}} = \left(\frac{\partial L}{\partial \mathbf{y}}\right)_{M \times D} \\ N \times D \\ N \times D \\ N \times D \\ N \times M \end{cases} \\ \frac{\partial L}{\partial \mathbf{w}} = \mathbf{x}_{D \times N}^{T} \left(\frac{\partial L}{\partial \mathbf{y}}\right)_{N \times M} \end{cases}$$

Mechanistic Data Science

Deep Learning - 55

$$\mathbf{A} \begin{pmatrix} \mathbf{B} & \mathbf{C} \\ n \times p & p \times q \end{pmatrix} = (\mathbf{AB})\mathbf{C}$$
  

$$\mathbf{BC} : npq \rightarrow \mathbf{A} \begin{pmatrix} \mathbf{BC} \end{pmatrix} : mnq \rightarrow \text{cost} = npq + mnq = nq(m+p)$$
  

$$\mathbf{AB} : mnp \rightarrow (\mathbf{AB}) \mathbf{C} : mpq \rightarrow \text{cost} = mnp + mpq = mp(n+q)$$
  

$$_{m \times p} \stackrel{p \times q}{\xrightarrow{p \times q}}$$
  
specialize to  $\mathbf{C} = \text{column vector } (q=1) : \mathbf{A}(\mathbf{BC}) \text{ vs. } (\mathbf{AB}) \mathbf{C}$   

$$\underbrace{\overset{\text{vector}}{\xrightarrow{n(m+p)}} \quad \underbrace{\overset{\text{big number}}{\xrightarrow{p \times q}}}_{mp(n+1)}$$