

# Activation Functions in Neural Networks

This handout presents common activation functions with their formulas, conceptual motivation, pros/cons, computational aspects, and typical use cases.

Table 1: Common Activation Functions in Neural Networks

Activation	Formula / Concept	Motivation (Conceptual)	Advantages	Disadvantages	Computational
Sigmoid	$f(x) = \frac{1}{1+e^{-x}}$	Encodes “probabilistic firing” with smooth saturation; bounded output aids probability interpretation.	<ul style="list-style-type: none"> <li>Bounded in (0, 1)</li> <li>Probabilistic interpretation</li> </ul>	<ul style="list-style-type: none"> <li>Vanishing gradients for large <math> x </math></li> <li>Not zero-centered</li> </ul>	<ul style="list-style-type: none"> <li>Uses <code>exp</code>; mod</li> </ul>
Tanh	$f(x) = \tanh(x)$	Center outputs at zero to speed optimization; antisymmetric.	<ul style="list-style-type: none"> <li>Zero-mean activations</li> <li>Smoother than sigmoid</li> </ul>	<ul style="list-style-type: none"> <li>Still saturates <math>\Rightarrow</math> vanishing gradients</li> </ul>	<ul style="list-style-type: none"> <li>Similar cost to</li> </ul>
ReLU	$f(x) = \max(0, x)$	Mimics neuron firing only for positive signals; encourages sparsity.	<ul style="list-style-type: none"> <li>Simple and very fast</li> <li>Avoids saturation for <math>x &gt; 0</math></li> </ul>	<ul style="list-style-type: none"> <li>“Dying ReLU” for persistently negative inputs</li> <li>Unbounded output</li> </ul>	<ul style="list-style-type: none"> <li>Branch/compute (cheap)</li> </ul>
Leaky ReLU	$f(x) = \max(\alpha x, x), \alpha \approx 10^{-2}$	Preserve gradient flow for negative region to avoid dead units.	<ul style="list-style-type: none"> <li>Mitigates dying ReLU</li> </ul>	<ul style="list-style-type: none"> <li>Small bias for <math>x &lt; 0</math></li> </ul>	<ul style="list-style-type: none"> <li>One multiply + ReLU</li> </ul>
PReLU	$f(x) = \max(\alpha x, x), \alpha$ learned	Learn slope for negatives from data to adapt nonlinearity.	<ul style="list-style-type: none"> <li>Adaptive capacity can improve accuracy</li> </ul>	<ul style="list-style-type: none"> <li>Extra parameters; mild overfitting risk</li> </ul>	<ul style="list-style-type: none"> <li>Slightly higher compute/memory</li> </ul>
ELU	$f(x) = \begin{cases} x, & x > 0 \\ \alpha(e^x - 1), & x \leq 0 \end{cases}$	Push activations’ mean toward zero; smooth negative branch.	<ul style="list-style-type: none"> <li>Faster convergence in some settings</li> <li>Avoids dead neurons</li> </ul>	<ul style="list-style-type: none"> <li>Requires <code>exp</code>; sensitive to <math>\alpha</math></li> </ul>	<ul style="list-style-type: none"> <li>Higher cost than</li> </ul>
SELU	$\text{SELU}(x) = \lambda \cdot \begin{cases} x, & x > 0 \\ \alpha(e^x - 1), & x \leq 0 \end{cases}$	Self-normalizing activations that keep layer activations near zero mean/unit variance.	<ul style="list-style-type: none"> <li>Stable activations without Batch-Norm</li> </ul>	<ul style="list-style-type: none"> <li>Requires LeCun normal init; architectural constraints</li> </ul>	<ul style="list-style-type: none"> <li>Slightly costlier</li> </ul>
Swish	$f(x) = x \sigma(x) = \frac{x}{1+e^{-x}}$	Smooth gating of input magnitude; continuous ReLU alternative.	<ul style="list-style-type: none"> <li>Smooth gradients; strong empirical results</li> </ul>	<ul style="list-style-type: none"> <li>Slower than ReLU; requires sigmoid</li> </ul>	<ul style="list-style-type: none"> <li>Multiply + <code>exp</code></li> </ul>
Mish	$f(x) = x \tanh(\ln(1 + e^x))$	Smooth, robust ReLU-like activation with improved gradient flow.	<ul style="list-style-type: none"> <li>Good stability and propagation</li> </ul>	<ul style="list-style-type: none"> <li>Higher cost; less standardized</li> </ul>	<ul style="list-style-type: none"> <li>Uses <code>tanh</code> and</li> </ul>
GELU	$f(x) = x \Phi(x) \approx \frac{x}{2} \left(1 + \text{erf}\left(\frac{x}{\sqrt{2}}\right)\right)$	Probabilistic “gating by magnitude”; smooth ReLU-like behavior.	<ul style="list-style-type: none"> <li>Smooth; strong results in Transformers</li> </ul>	<ul style="list-style-type: none"> <li>Slightly slower; derivative more complex</li> </ul>	<ul style="list-style-type: none"> <li>Uses <code>erf</code>/norm</li> </ul>
Softmax	$f_i(\mathbf{x}) = \frac{e^{x_i}}{\sum_j e^{x_j}}$	Convert logits to a categorical probability distribution.	<ul style="list-style-type: none"> <li>Interpretable probabilities that sum to 1</li> </ul>	<ul style="list-style-type: none"> <li>Sensitive to large logits; numerical overflow</li> </ul>	<ul style="list-style-type: none"> <li>Vector <code>exp</code> +</li> </ul>
Linear	$f(x) = x$	Preserve numeric information without distortion.	<ul style="list-style-type: none"> <li>Simple; suitable for regression</li> </ul>	<ul style="list-style-type: none"> <li>No nonlinearity <math>\Rightarrow</math> limited expressivity</li> </ul>	<ul style="list-style-type: none"> <li>Minimal cost (mapping)</li> </ul>

**Pedagogical Notes.** Nonlinearity is essential (without it, deep stacks collapse to an affine map). ReLU-family dominates CNN/MLP hidden layers; GELU/Swish are common in Transformers; tanh/sigmoid remain in RNN gates; SELU expects LeCun-normal initialization. Output activations depend on the task: Linear (regression), Sigmoid (binary), Softmax (multiclass). BatchNorm and initialization strongly interact with activation choice.