# Machine Learning
## Activation Functions

Curtis Larsen

Utah Tech University—Computing

Spring 2025

# Objectives

**Objectives:**

▶ Review - Activation Functions in Neural Networks

▶ Review - Back-propagation

▶ Understand - Vanishing and Exploding Gradients

▶ Identify - Common Activation Functions

▶ Initialize - Best Practices per Activation Function

Review

# Activation Functions

$$\phi_w(x) = \sum_{i=0}^{n} w_i x_i$$

# Loss and Gradient

Loss function

$$J(\boldsymbol{\theta})$$

Gradient w.r.t $\boldsymbol{\theta}$

$$\boldsymbol{\nabla}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

# Back-propagation

- ▶ Forward pass (compute)
- ▶ Backward pass (error propagation)
- ▶ Vanishing gradients
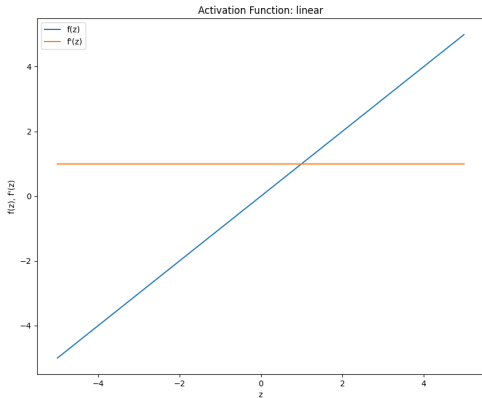- ▶ Exploding gradients

# Common Activation Functions

# Linear
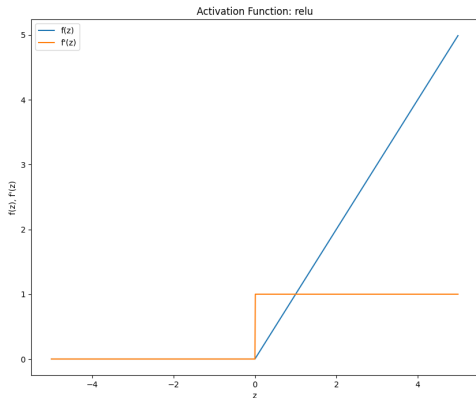
$$f(x) = x$$

**Notes:**

▶ Smooth

▶ Fast

▶ `linear`

# Rectified Linear

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \le 0 \end{cases}$$

**Notes:**

▶ Not smooth

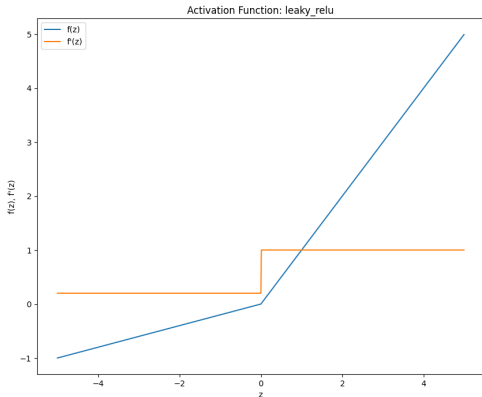▶ Fast

▶ No gradient on left side (neurons "die")

▶ `relu`

# Leaky ReLU

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha x & \text{if } x \leq 0 \end{cases}$$

**Notes:**

▶ Not smooth

▶ Fast

▶ Small gradient on left side (neurons don't "die")
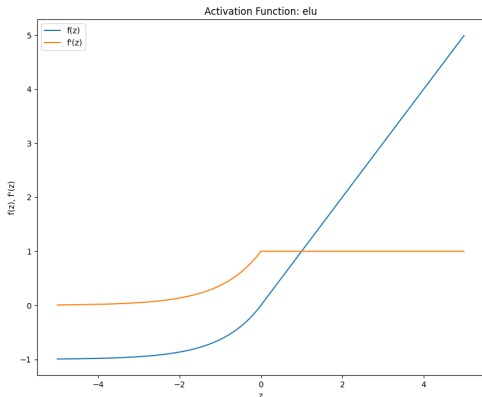
▶ `leaky_relu`



Activation Function: leaky_relu

# Exponential Linear Unit

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x \leq 0 \end{cases}$$

**Notes:**

▶ Smooth if $\alpha = 1$

▶ Slower to compute

▶ Gradient descent may converge faster

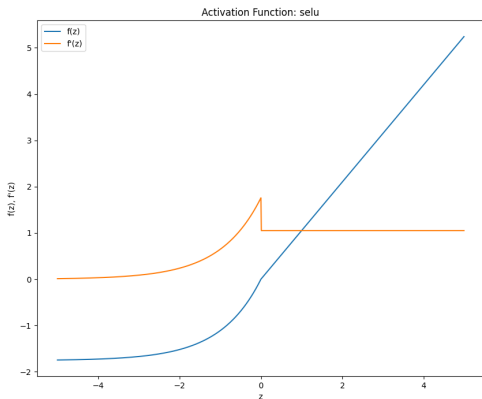▶ Small gradient on left side (neurons don't "die")

▶ `elu`



Activation Function: elu

# Scaled ELU

$$f(x) = \begin{cases} sx & \text{if } x > \\ s\alpha(e^x - 1) & \text{if } x \leq \end{cases}$$

**Notes:**

- ▶ Only useful on stacks of dense layers

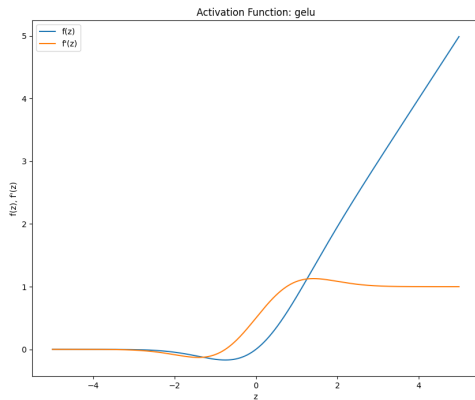- ▶ `selu`



Activation Function: selu

# Gaussian ELU

$$f(x) = x\Phi(x)$$

**Notes:**

- ▶ $\Phi(x)$ is the Gaussian cumulative distribution function

- ▶ More computationally expensive
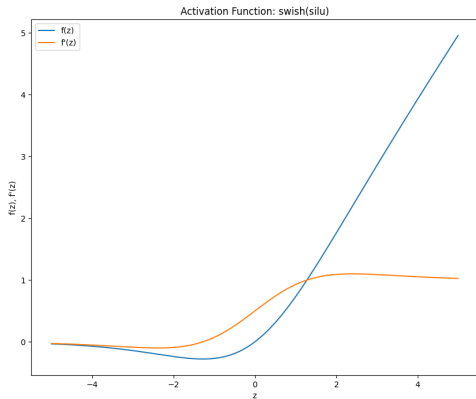
- ▶ Maybe faster convergence

- ▶ `gelu`

# Sigmoid Linear Unit (Swish)

$$f(x) = x\sigma(x)$$

**Notes:**

▶ $\sigma(x)$ is the sigmoid function
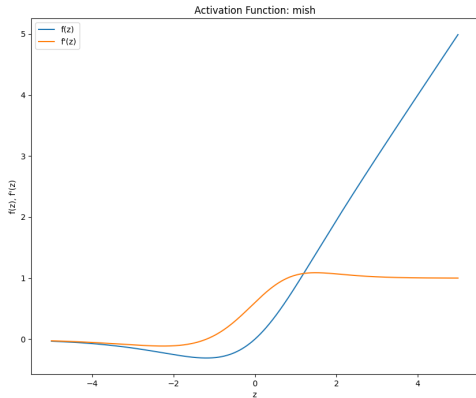
▶ Sometimes better than GELU

▶ `silu`, `swish`



Activation Function: swish(silu)

# MISH

$$f(x) = x \tanh(\log(1+e^x))$$

**Notes:**

▶ Sometimes better than GELU, Swish
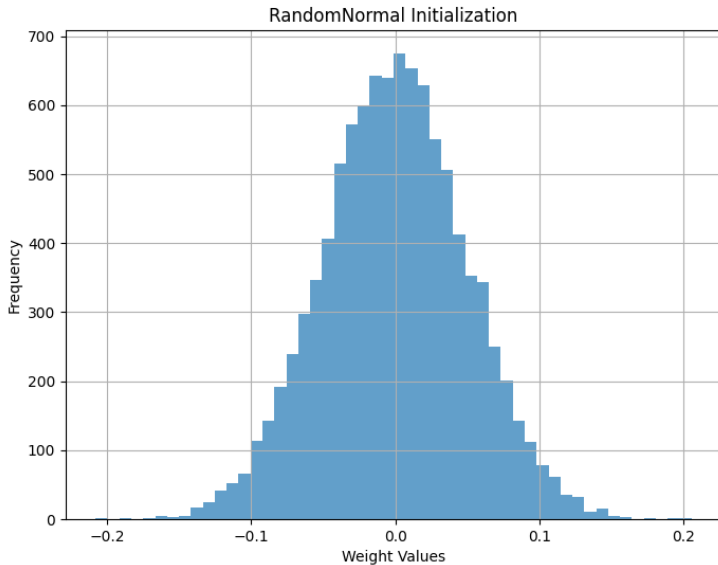
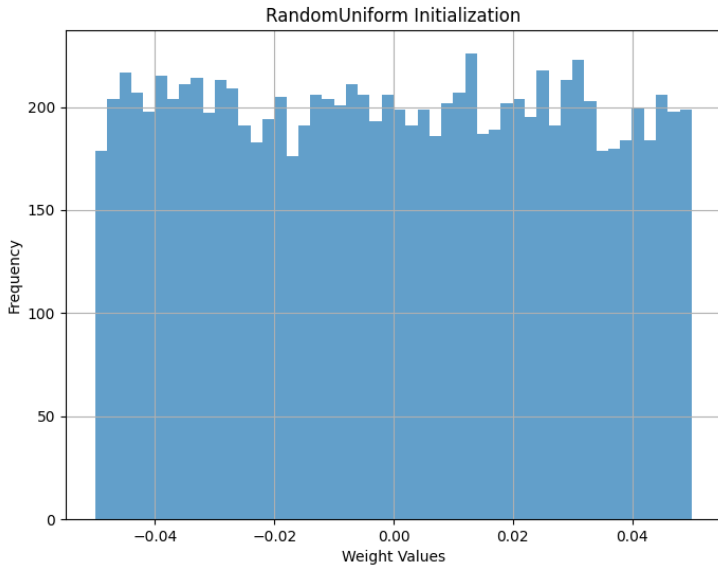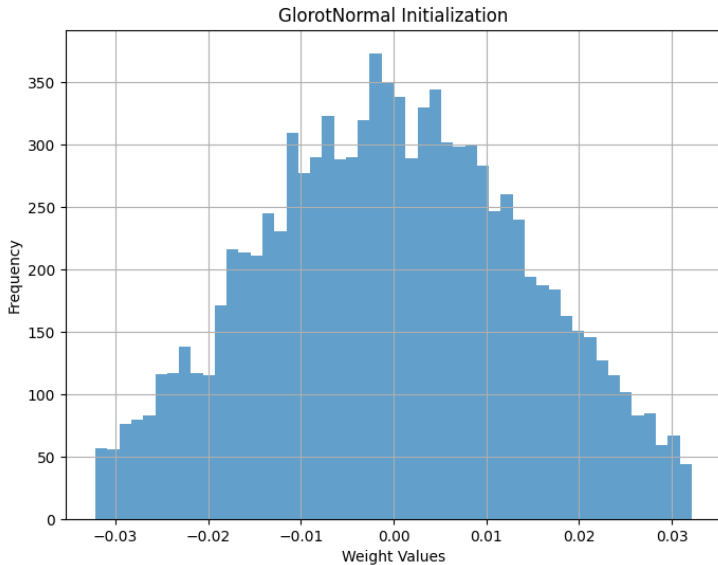▶ `mish`



Activation Function: mish

# Implementation

```
model.add(keras.layers.Dense(units, activation="leaky_relu"))
```
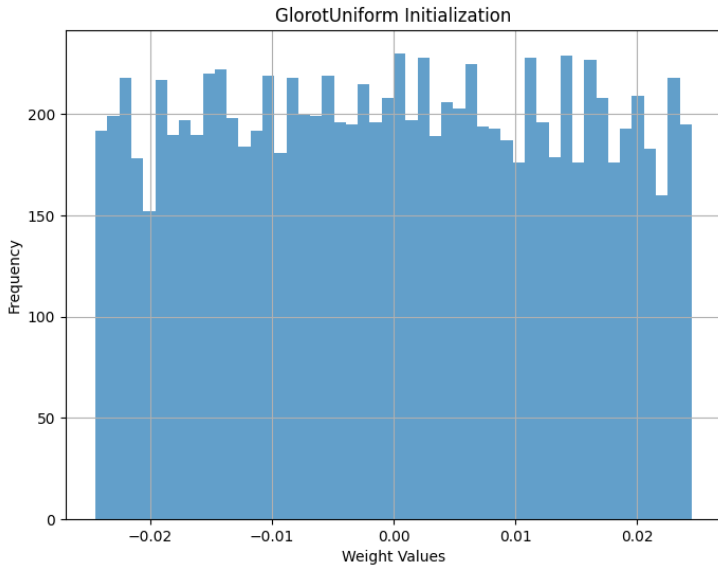
# Weight Initialization

RandomNormal Initialization

RandomUniform Initialization

GlorotNormal Initialization

GlorotUniform Initialization

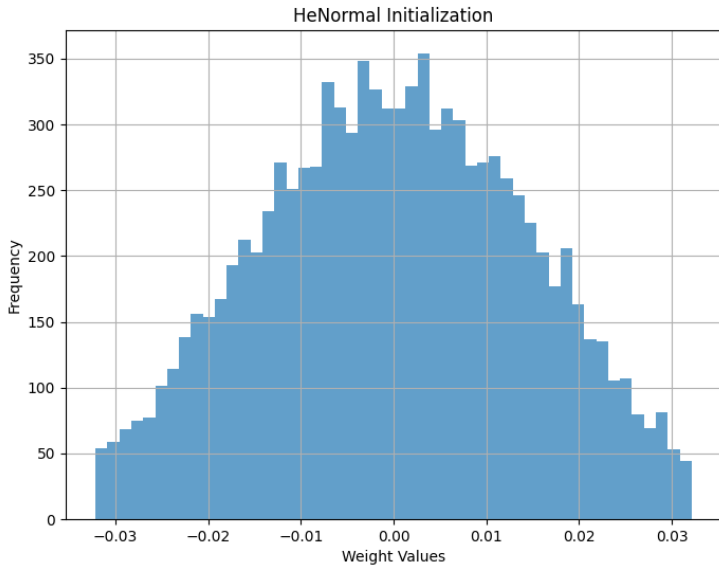HeNormal Initialization

HeUniform Initialization

LecunNormal Initialization

LecunUniform Initialization

| Activation | Recommended Initialization |
|---|---|
| Linear | Glorot |
| Sigmoid | Glorot |
| Hyperbolic Tangent | Glorot |
| Rectified Linear Unit | He |
| Leaky ReLU | He |
| Exponential Linear Unit | He |
| Scaled Exponential Linear Unit | LeCun |
| Gaussian Error Linear Unit | He |
| Sigmoid Linear Unit, Swish | He |
| Mish | He |

# Implementation

```
model.add(keras.layers.Dense(units, activation="leaky_relu",
      kernel_intializer="he_normal"))
```

# Summary

# Topics

► Activation Functions

► Initializations

## Activation Functions

| Function | Keras Name |
|---|---|
| Linear | linear |
| Sigmoid | sigmoid |
| Hyperbolic Tangent | tanh |
| Rectified Linear Unit | relu |
| Leaky ReLU | leak_relu |
| Exponential Linear Unit | elu |
| Scaled Exponential Linear Unit | selu |
| Gaussian Error Linear Unit | gelu |
| Sigmoid Linear Unit, Swish | swish,silu |
| Mish | mish |

# Initializations

| Function | Keras Name |
|---|---|
| Zeros | `zeros` |
| Ones | `ones` |
| Random Normal | `random_normal` |
| Random Uniform | `random_uniform` |
| Glorot (Xavier) Normal | `glorot_normal` |
| Glorot (Xavier) Uniform | `glorot_uniform` |
| He Normal | `he_normal` |
| He Uniform | `he_uniform` |
| Lecun Normal | `lecun_normal` |
| Lecun Uniform | `lecun_uniform` |