

A glowing brain with neural network connections and data points, symbolizing deep learning. The brain is rendered in a wireframe style with blue and red highlights, set against a dark background with a complex network of light blue and red lines and dots, suggesting a neural network or data flow. The text "DEEP LEARNING" is centered over the brain in a bold, white, sans-serif font.

DEEP LEARNING

Table of Contents

Introduction	1
Purpose	1
Scope	1
Procedure	1
Brief Problem	1
What is Deep Learning?	2
Input Layer:	2
Hidden Layers:	2
Output Layer:	2
Feature Extraction and Classification: ML vs DL	3
Performance vs Amount of Data	3
Forward Propagation	4
Backward Propagation (Backpropagation)	4
Optimizers	5
The Core Idea	5
Momentum	5
Idea	5
Benefits	5
Nesterov Accelerated Gradient (NAG)	6
Key difference	6
Adagrad (Adaptive Gradient)	7
Idea	7
Strengths	7
Weakness	7
RMSProp	8
Idea	8
Advantages	8
Adam (Adaptive Moment Estimation)	9
Key Idea	9
Why It Works	9
When to Use	9
AdamW (Adam + Weight Decay Done Right)	10
Why It Matters	10

Use When.....	10
Regularization in Deep Learning	11
1. L1 and L2 Regularization	11
2. Dropout	12
3. Early Stopping	12
4. Batch Normalization.....	12
5. Data Augmentation	12
Activation Functions	13
What Are Activation Functions?	13
Types of Activation Functions	13
1. Sigmoid (Logistic)	13
2. Tanh (Hyperbolic Tangent)	14
3. ReLU (Rectified Linear Unit)	15
4. Leaky ReLU	16
5. Softmax	17
Choosing an Activation Function	19
Vanishing & Exploding Gradients	20
1. Vanishing Gradient.....	20
2. Exploding Gradient.....	20
How Activation Functions Handle Them.....	20
Loss Functions	21
Regression Losses.....	21
Binary Classification Losses.....	21
Multi-Class Classification Losses	22
Probabilistic / Distribution Losses.....	22
Data Preprocessing in DL	23
Loading Data	23
1. ImageDataGenerator	23
2. Using tf.keras.utils.image_dataset_from_directory.....	24
3. Loading from File Paths and Custom Pipelines	25
4. Loading from NumPy Arrays or Pandas DataFrames	26
Data Augmentation.....	27
A. Using ImageDataGenerator	27
B. Using tf.keras.layers (preferred modern method)	28

Other Important Preprocessing Steps.....	29
1. Normalization / Standardization	29
2. One-Hot or Label Encoding	29
3. Balancing Classes	29
Visualization	30
A. Visualize 5–10 Samples from Each Class	30
B. Visualize After Augmentation	31
Perceptron (Single Neuron Model)	32
Working Mechanism	32
Example (Binary Classification)	32
Limitation	33
ANN (Artificial Neural Network).....	34
Structure (Common to All ANNs)	34
ANN = The Family of Neural Architectures	34
Multi-Layer Perceptron (MLP)	35
Mathematically	35
Advantages.....	35
Deep Neural Network (DNN).....	37
Representation Learning.....	37
Architecture Variations	37
Transfer Learning.....	39
What Is Transfer Learning?	39
Why Do We Need Transfer Learning?	39
The Core Idea (How It Works).....	40
What Does a Pretrained Model Contain?	40
Early Layers:	40
Middle Layers:.....	40
Final Layers:.....	40
Types of Transfer Learning	41
Feature Extraction.....	41
Fine-Tuning.....	41
When Does Transfer Learning Work Best?	42
Case 1: Very Similar Tasks	42
Case 2: Moderately Similar Tasks.....	42

Case 3: Very Different Tasks	42
Benefits of Transfer Learning	43
Transfer Learning Workflow	43
When NOT to Use Transfer Learning.....	43
Key Hyperparameters in Transfer Learning.....	44
Learning Rate Rule	44
TensorFlow	45
1. The Model (The Container)	45
A. Sequential API (Easiest)	45
B. Functional API (Most Flexible)	45
2. Layers (The Building Blocks)	46
3. Activation Functions (The Non-Linearity)	47
4. Loss Functions (The Goal)	47
5. Optimizers (The Driver)	48
6. Metrics (The Scoreboard).....	48
7. Putting It All Together: A "Template" Code.....	49
8. Essential "Ease of Use" Features.....	50
Callbacks:.....	50
Model Saving/Loading:	50
Conclusion	51
Feedback & Contribution	51
Copyright & Usage.....	51
Acknowledgments	52

Introduction

Purpose

The purpose of this document is to provide a comprehensive overview of **Deep Learning**, one of the most transformative areas within Artificial Intelligence.

Scope

Includes understanding the theoretical foundations of deep neural networks, their architecture, learning mechanisms, and the wide range of applications they enable.

Procedure

Involves breaking down the core components of deep learning—such as neurons, layers, activation functions, backpropagation, optimization, and regularization—and examining how these contribute to building powerful models capable of learning complex data patterns. This document also explores different types of deep learning networks, including Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Transformers, along with their practical use cases.

Brief Problem

Deep learning addresses the limitations of traditional machine learning methods when handling high-dimensional, unstructured, or large-scale data such as images, audio, and natural language. Deep learning provides the solution through hierarchical feature learning and representation, enabling systems to automatically discover intricate relationships without manual intervention.

What is Deep Learning?

Deep Learning (DL) is a subfield of **Machine Learning (ML)** that focuses on algorithms inspired by the structure and function of the human brain — **Artificial Neural Networks (ANNs)**.

While ML relies heavily on **manual feature engineering**, Deep Learning automatically **learns feature hierarchies** directly from raw data through multiple layers of representation.

In simpler terms:

- **Machine Learning:** You give the model both the **data** and the **rules/features** to learn from.
- **Deep Learning:** You give the model just the **data**, and it learns the **rules/features** on its own through neural networks.

How Deep Learning Handles Inputs and Outputs

Input Layer:

- Accepts raw data (images, text, audio, etc.).
- Each neuron in the input layer represents one feature or pixel/value from the dataset.

Hidden Layers:

- These layers transform input data step by step.
- Each neuron computes a weighted sum of its inputs, applies an **activation function**, and passes it forward.
- Multiple hidden layers allow the model to learn **complex, abstract representations** (e.g., edges → shapes → objects in images).

Output Layer:

- Provides the final prediction or classification.
- For example:
 - In regression → single continuous output (e.g., price).
 - In classification → multiple categories with probabilities (using softmax).

Feature Extraction and Classification: ML vs DL

Aspect	Machine Learning	Deep Learning
Feature Extraction	Done manually by humans (e.g., extracting edges, texture, or specific statistics).	Done automatically by the neural network through hierarchical layers.
Model Example	Logistic Regression, SVM, Decision Tree.	CNN, RNN, Transformers.
Input Data	Requires structured, preprocessed features.	Works directly on raw data (images, text, sound).
Computation	Faster to train, less data-hungry.	Requires more data and compute, but yields higher performance.
Interpretability	Easier to interpret.	Acts as a black box (less interpretable).

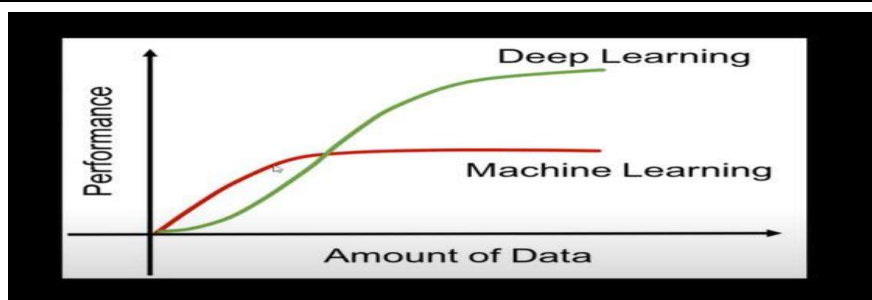
👉 Example:

In ML, you might extract color histograms or texture manually to classify images.

In DL, a **CNN** automatically detects edges → corners → objects through training.

Performance vs Amount of Data

Model Type	Behavior with Data
Machine Learning	Performance increases initially, then plateaus when more data doesn't significantly improve results (since features are manually designed).
Deep Learning	Performance continues improving as data increases — because the model can learn deeper, more complex patterns given enough examples.



Forward Propagation

Forward propagation is how the neural network makes predictions.

1. Data enters the input layer.
2. Each neuron computes:

$$z = w_1x_1 + w_2x_2 + \dots + b$$

3. The neuron applies an activation function:

$$a = f(z)$$

4. Output of each layer becomes the input of the next layer.
5. The final layer produces the prediction (\hat{y}).

So, **forward propagation = computing outputs based on current weights.**

Backward Propagation (Backpropagation)

After making predictions, the model checks **how wrong** it was and **adjusts** the weights accordingly — this is **backpropagation**.

1. Compute **loss (error)** between predicted output (\hat{y}) and true output (y):

$$L = \text{Loss}(y, \hat{y})$$

2. Compute **gradients** (partial derivatives of loss with respect to weights).
3. Propagate the gradients backward from the output layer to the input layer.
4. Update weights using **gradient descent**:

$$w := w - \eta \frac{\partial L}{\partial w}$$

Where:

- η = learning rate
- $\frac{\partial L}{\partial w}$ = gradient (how much each weight affects the loss)

Optimizers

The Core Idea

Training a neural network means **adjusting weights** to minimize a **loss function** (e.g., MSE, cross-entropy).

To do that, we compute:

- The **gradient** → tells us how to change weights to reduce error
 - The **optimizer** → decides *how* to apply that gradient
-

Momentum

SGD sometimes:

- Oscillates back and forth
 - Gets stuck in flat regions
-

Idea

Carry “velocity” from previous updates—like pushing a ball downhill:

$$v_t = \beta v_{t-1} + \alpha \nabla L$$
$$w_{t+1} = w_t - v_t$$

Where:

- β = momentum coefficient (e.g., 0.9)
-

Benefits

- Reduces oscillations
- Faster convergence
- Can escape shallow minima

Nesterov Accelerated Gradient (NAG)

A smarter momentum.

Key difference

Instead of calculating gradient at current position:

- Look **ahead first**, then compute gradient

$$v_t = \beta v_{t-1} + \alpha \nabla L(w - \beta v_{t-1})$$

When to Use

- If training is unstable with momentum
- Classic for older deep nets

Adagrad (Adaptive Gradient)

Idea

Every parameter gets its **own adaptive learning rate**, based on past gradients:

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{G_t + \epsilon}} \cdot \nabla L_t$$

Where G_t stores accumulated squared gradients.

Strengths

- Great for **sparse features** (text, NLP)
 - Frequently-updated weights get lower LR
 - Rarely-updated weights get higher LR
-

Weakness

- Learning rate **keeps shrinking**
- Eventually becomes too small → stops learning

RMSProp

Created to fix Adagrad's shrinking LR problem.

Idea

Use **exponentially decaying average** instead of accumulating forever:

$$E[g^2]_t = \beta E[g^2]_{t-1} + (1 - \beta)g_t^2$$

Then update:

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

Advantages

- Maintains adaptive learning rates
- Works extremely well for:
 - RNNs
 - Reinforcement learning
 - Online learning

Adam (Adaptive Moment Estimation)

Adam is the modern default used everywhere. It has a learning rate for each parameter.

Key Idea

Combines:

- **Momentum** (moving average of gradients)
- **RMSProp** (moving average of squared gradients)

Two moving averages:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$
$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

Bias-corrected:

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

Update:

$$w_{t+1} = w_t - \frac{\alpha \hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

Why It Works

- Handles sparse data
 - Handles noisy gradients
 - Works well even with minimal tuning
-

When to Use

- Deep learning in general
- CNNs, RNNs, Transformers
- If you're not sure → **start with Adam**

AdamW (Adam + Weight Decay Done Right)

In Adam, L2 regularization is implemented incorrectly.

AdamW fixes it by decoupling weight decay:

$$w_{t+1} = w_t - \alpha(\text{Adam update}) - \alpha\lambda w_t$$

Why It Matters

Modern models like Transformers **train much better with AdamW**.

Use When

- Large deep networks
- NLP models
- Vision models

Regularization in Deep Learning

Regularization prevents **overfitting**, which happens when your model learns the training data *too well* and fails to generalize.

1. L1 and L2 Regularization

These add penalties on **large weights** to make the model simpler.

L2 Regularization (Ridge)

Adds a penalty proportional to the **square of the weights**:

$$L_{total} = L_{original} + \lambda \sum w^2$$

→ Pushes weights toward smaller values (but rarely zero).

Used by default in many optimizers (e.g., weight decay in AdamW).

L1 Regularization (Lasso)

Adds a penalty proportional to the **absolute value** of the weights:

$$L_{total} = L_{original} + \lambda \sum |w|$$

→ Encourages **sparse models** (some weights become exactly 0).

Keras Example:

```
Dense(64, activation='relu', kernel_regularizer=tf.keras.regularizers.l2(0.001))
```

2. Dropout

Randomly “drops” (sets to 0) a fraction of neurons during training.

This prevents co-dependence between neurons and improves generalization.

```
Dropout(0.5)
```

means 50% of neurons are dropped each batch.

3. Early Stopping

Stop training when validation loss stops improving — prevents overfitting.

```
EarlyStopping(monitor='val_loss', patience=5, restore_best_weights=True)
```

4. Batch Normalization

Normalizes the activations between layers, stabilizing learning and providing slight regularization.

```
BatchNormalization()
```

5. Data Augmentation

Artificially increase dataset diversity (e.g., rotating, flipping, or scaling images).

```
ImageDataGenerator(rotation_range=20, horizontal_flip=True)
```

Activation Functions

What Are Activation Functions?

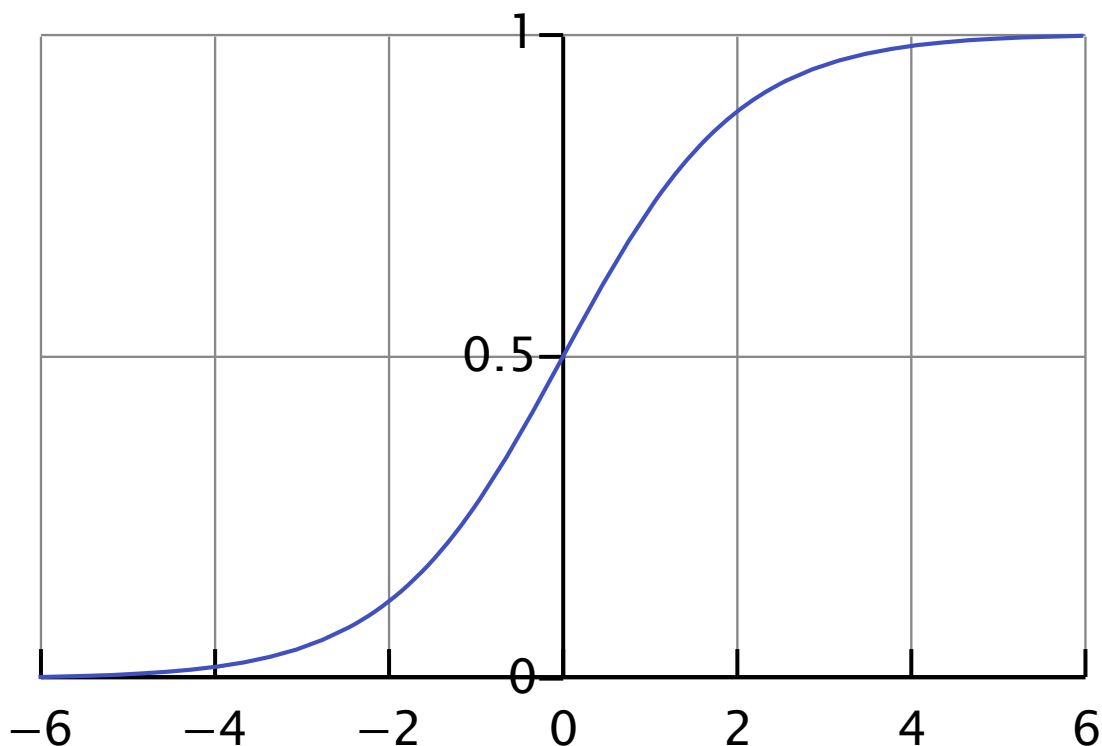
- Activation functions decide **whether a neuron should be activated or not**, by introducing **non-linearity** into the network.
 - Without them, neural networks would just be linear models → unable to learn complex patterns.
-

Types of Activation Functions

1. Sigmoid (Logistic)

$$\frac{1}{1 + e^{-x}}$$

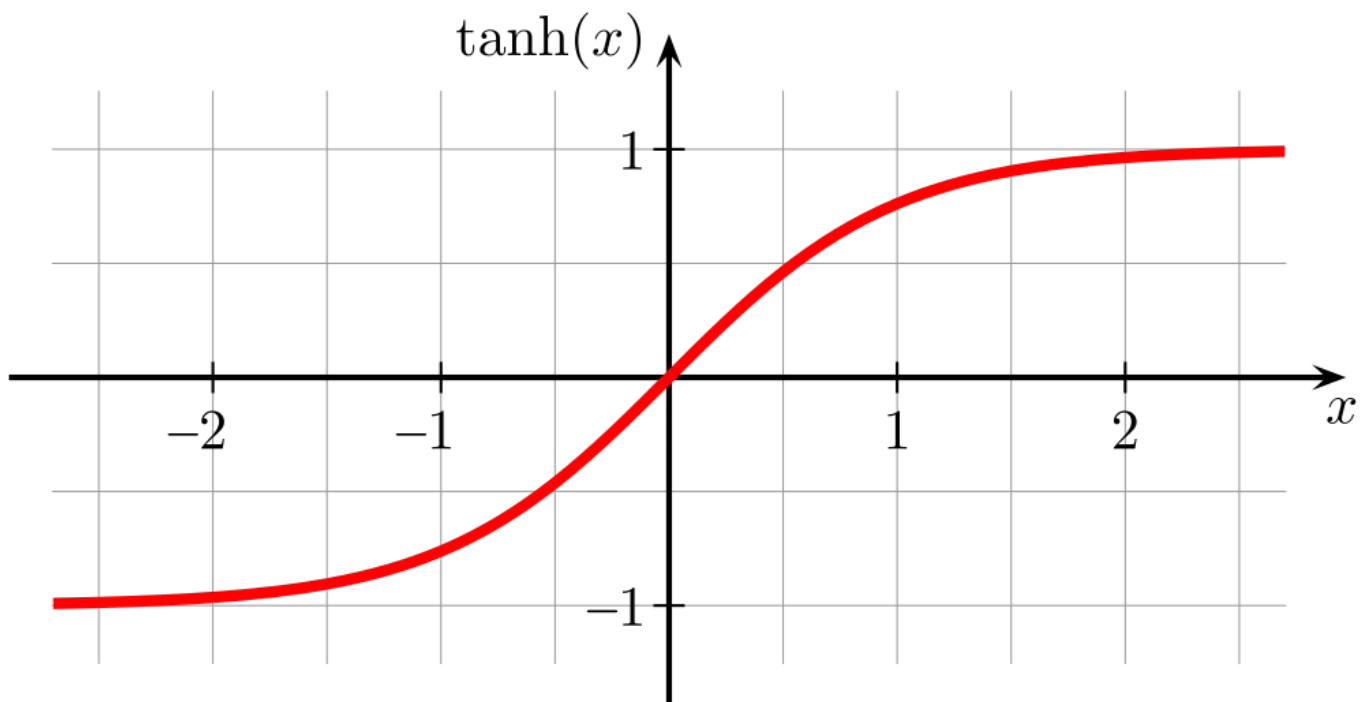
- Used in Output Layer
- Outputs between **0 and 1**.
- Historically popular for binary classification.
- **Problems:** Slow convergence



2. Tanh (Hyperbolic Tangent)

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

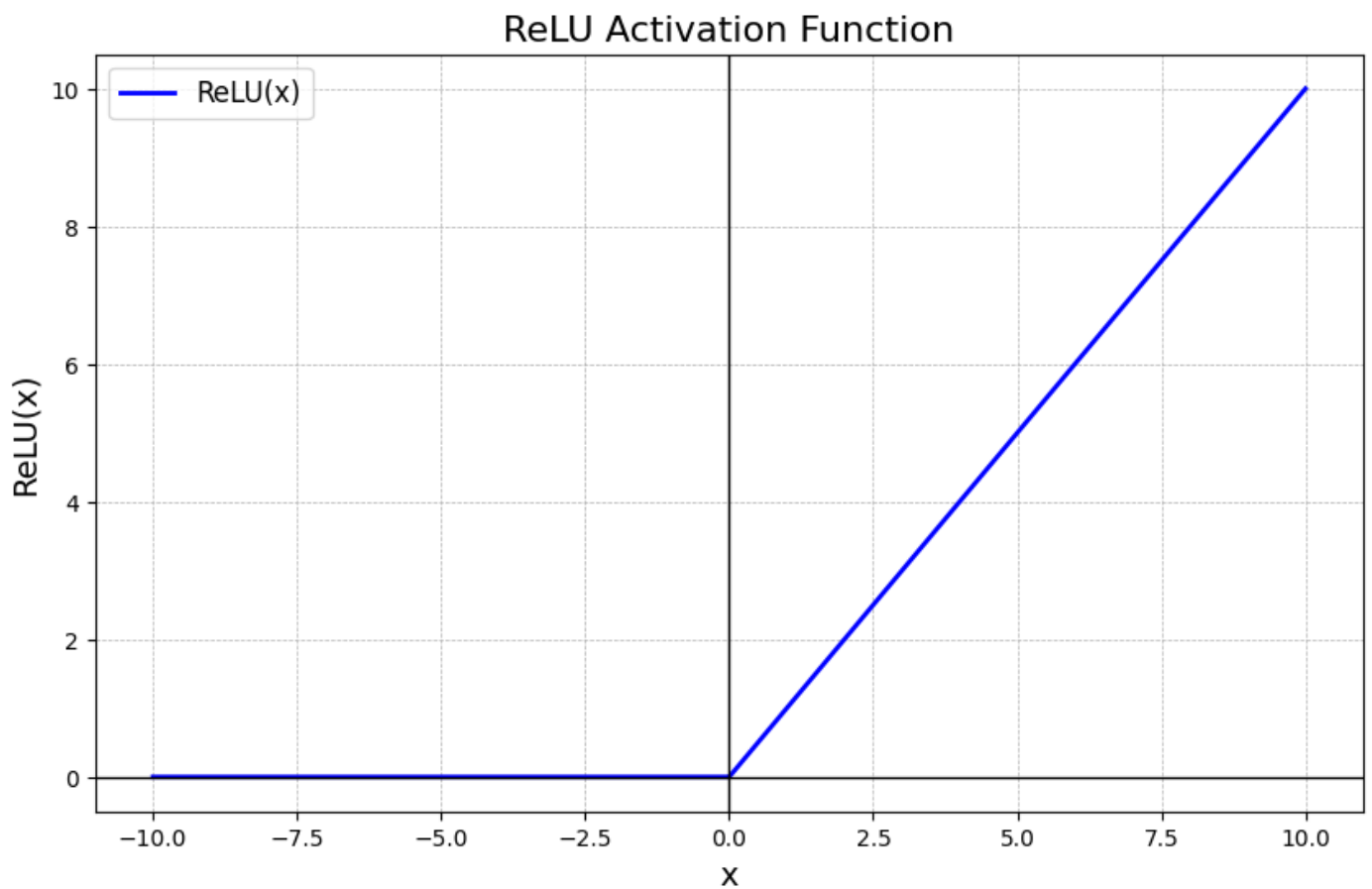
- Used in Output Layer
- Outputs between **-1 and 1**.
- Zero-centered → better than sigmoid.
- **Problem:** Still suffers from vanishing gradients.



3. ReLU (Rectified Linear Unit)

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

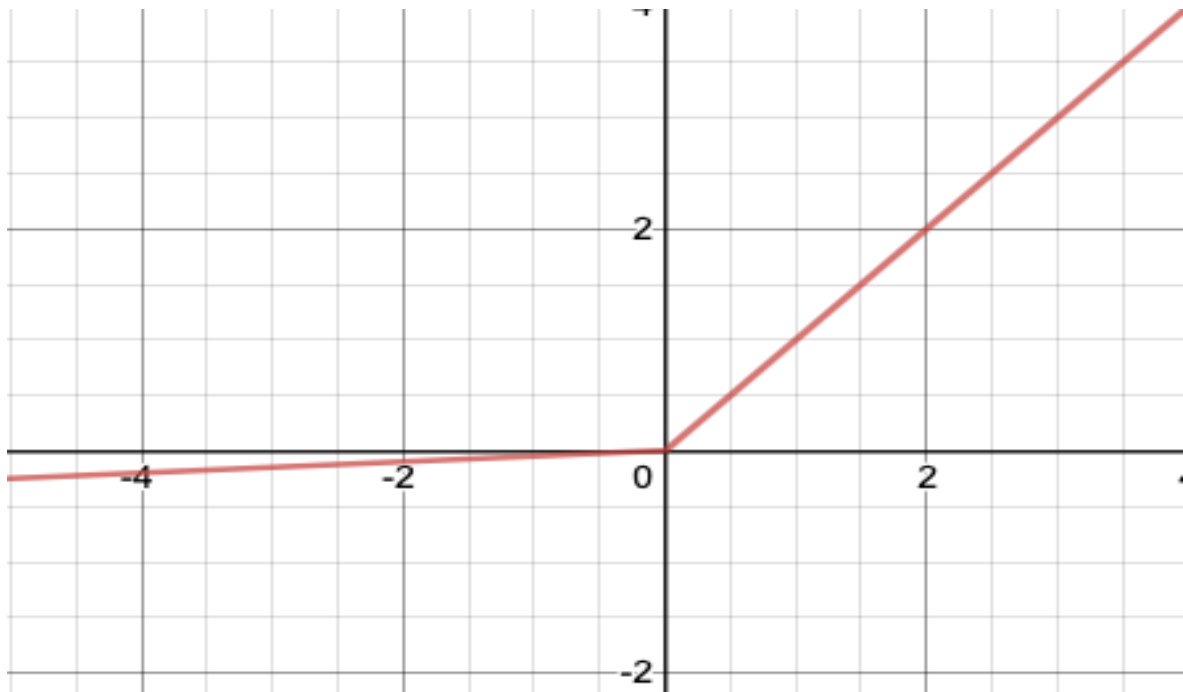
- Most widely used in deep networks.
- Simple, fast, reduces vanishing gradient issue.
- Used in Hidden Layer
- **Problems:** “Dead neurons” (weights stop updating if stuck at 0).



4. Leaky ReLU

$$f(x) = \max(\alpha x, x)$$

- Fixes dead neuron problem by allowing a small negative slope.
- Used in Hidden Layer



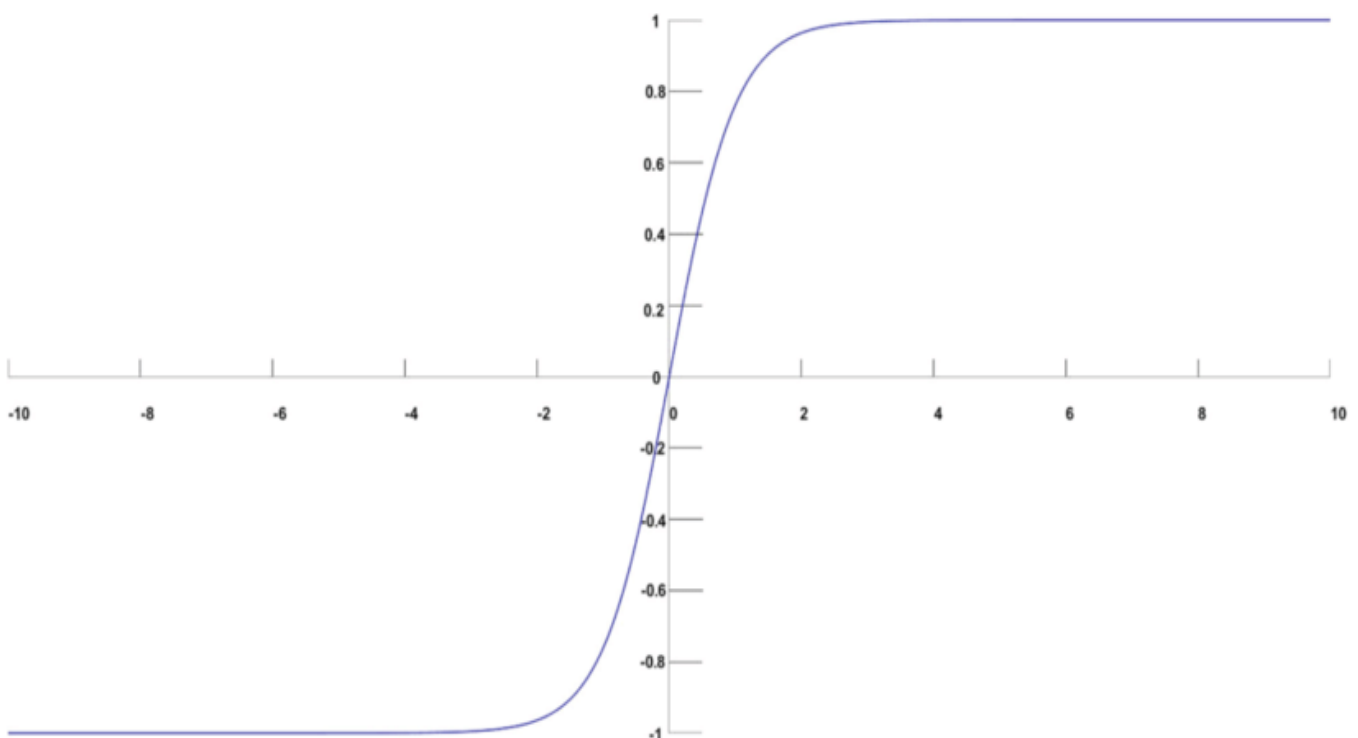
5. Softmax

The Core Idea: Scores to Probabilities

In a multi-class problem (say, 3 classes: Cat, Dog, Bird), your network's final layer will output a raw score (logit) for each class. These scores can be anything: negative, positive, or zero.

The **Softmax Function** squashes these scores so that:

- All values are between 0 and 1.
- The sum of all values equals 1 (creating a valid probability distribution).



TensorFlow Note

Using `from_logits = True` + linear activation function in output layer is the same as using sigmoid + being more stable and more numerically accurate as an intermediate value, but we need to apply sigmoid on these logits to get probabilities (logits are output of neurons in output layer without using an activation function).

The Mathematics

Given a vector of raw scores (logits)

$$z = [z_1, z_2, \dots, z_K]$$

for K classes, the Softmax function $\sigma(z)$ for the i -th class is:

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} = p(y=i | x)$$

Why use e (Euler's number)?

Exponentials make the largest input stand out significantly more (hence "Soft-max"). Even if one score is slightly higher than the others, the exponential function magnifies that difference, making the model more confident.

Example:

Raw Logits (z): [2.0, 1.0, 0.1]

Exponentials (e^z): [7.39, 2.72, 1.10]

Sum: 11.21

Softmax Probabilities:

- Class 1: $7.39/11.21 \approx \mathbf{0.66}$ (66%)
- Class 2: $2.72/11.21 \approx \mathbf{0.24}$ (24%)
- Class 3: $1.10/11.21 \approx \mathbf{0.10}$ (10%)

The Companion: Categorical Cross-Entropy

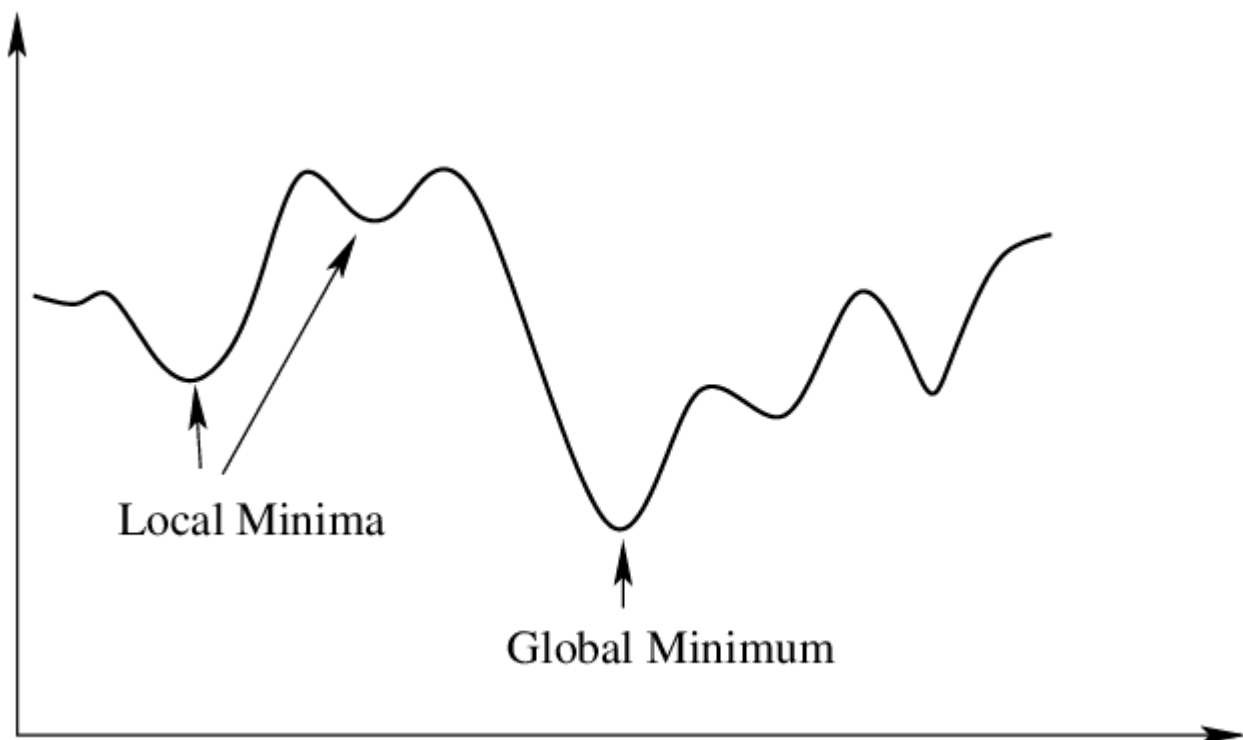
Softmax Regression is almost always paired with **Categorical Cross-Entropy Loss**.

- Softmax calculates the probabilities.
- Cross-Entropy calculates how different those probabilities are from the true label (which is 1.0 for the correct class and 0.0 for others).

If your target is "Class 1" (One-hot: [1, 0, 0]) and your prediction is [0.66, 0.24, 0.10], the loss focuses entirely on maximizing that 0.66.

Choosing an Activation Function

- **Hidden layers** → ReLU or Leaky ReLU (default choice) because it goes flat in one part of the graph not two like the sigmoid, when $g(z)$ is flat it becomes inefficient for gradient descent because of multiple local minima.
- **Output layer:**
 - Binary classification → Sigmoid.
 - Multi-class classification → Softmax.
 - Regression → Linear (+ve or -ve values), relu (+ve values only)



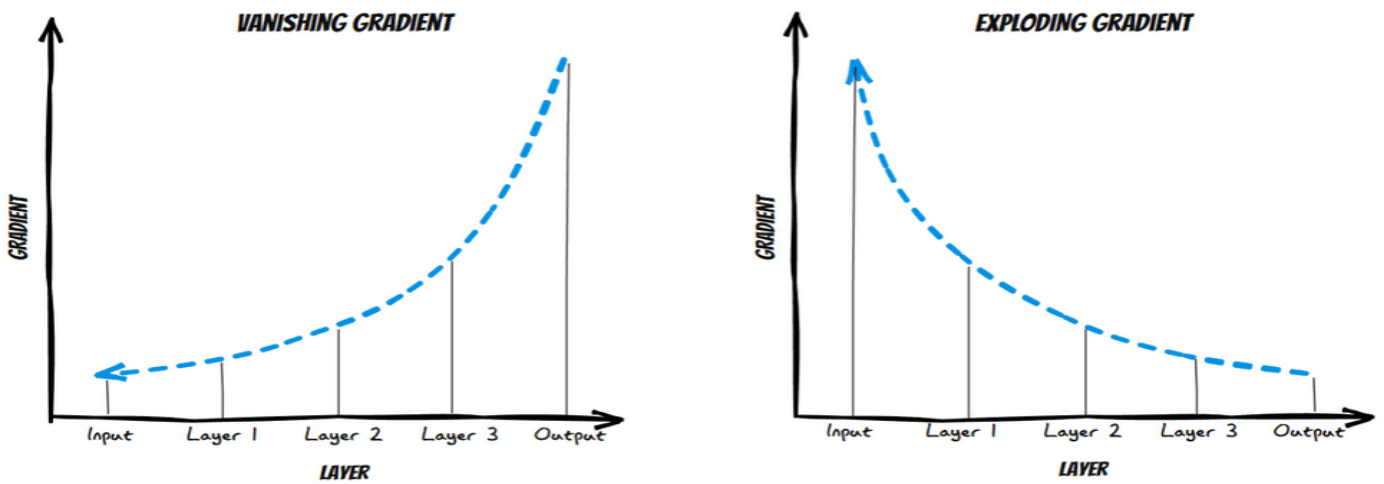
Vanishing & Exploding Gradients

1. Vanishing Gradient

- Happens when gradients become **very small** during backpropagation.
- Leads to **slow or no learning** in earlier layers.
- Common in **deep networks with Sigmoid/Tanh** because their gradients shrink for large $|x|$.

2. Exploding Gradient

- Happens when gradients grow **too large** during backpropagation.
- Leads to **unstable training** (loss oscillates or diverges).
- Common in very deep networks without proper normalization.



How Activation Functions Handle Them

Activation Function	Range	Vanishing Gradient?	Exploding Gradient?
Sigmoid	(0, 1)	✗ Severe (gradients vanish for large	✗
Tanh	(-1, 1)	✗ Still vanishes for large	✗
ReLU	$[0, \infty)$	✓ Avoids vanishing (gradient = 1 for $x > 0$)	⚠ Can explode if inputs/weights unscaled
Leaky ReLU	$(-\infty, \infty)$	✓ Avoids vanishing (small slope for $x < 0$)	⚠ Can explode if not regularized
Softmax	(0, 1)	✗ Vanishing if inputs are large	Not typical

Loss Functions

Loss (or cost) functions tell the model **how wrong its predictions are**.

During training, the model adjusts its weights to **minimize this loss**.

Different loss functions are used for different kinds of tasks.

Regression Losses

Loss Function	Mathematical Rule	When Used	TensorFlow Name
Mean Squared Error (MSE)	$L = \frac{1}{n} \sum (y - \hat{y})^2$	Continuous values (price, temperature)	"mse" or tf.keras.losses.MeanSquaredError()
Mean Absolute Error (MAE)	$L = \frac{1}{n} \sum y - \hat{y} $	A loss that is robust to outliers because it penalizes errors linearly instead of squaring them.	"mae" or tf.keras.losses.MeanAbsoluteError()
Huber Loss	Quadratic if small error, linear if large	Handles outliers better than MSE	tf.keras.losses.Huber()
Log-Cosh Loss	$\sum \log (\cosh (\hat{y} - y))$	Smooth alternative to MSE	tf.keras.losses.LogCosh()

Binary Classification Losses

Loss Function	Rule	When Used	TensorFlow Name
Binary Crossentropy	$L = -(y \log (\hat{y}) + (1 - y) \log (1 - \hat{y}))$	Yes/No prediction	"binary_crossentropy"
Hinge Loss	$L = \max (0, 1 - y\hat{y})$	SVM-style classification	tf.keras.losses.Hinge()
Squared Hinge	$L = \max (0, 1 - y\hat{y})^2$	Stronger penalty than hinge	SquaredHinge()

Multi-Class Classification Losses

Loss Function	Rule	When Used	TensorFlow Name
Categorical Crossentropy	$L = -\sum y_i \log(\hat{y}_i)$	One-hot encoded labels	"categorical_crossentropy"
Sparse Categorical Crossentropy	Same as above but labels are integers	Most common for multi-class	"sparse_categorical_crossentropy"
KLDivergence	$\sum y \log(y/\hat{y})$	Comparing distributions	tf.keras.losses.KLDivergence()

Probabilistic / Distribution Losses

Loss Function	Rule	When Used	TensorFlow Name
Poisson Loss	$L = \hat{y} - y \log(\hat{y})$	Count predictions	tf.keras.losses.Poisson()

Data Preprocessing in DL

Loading Data

1. ImageDataGenerator

This class allows **loading + preprocessing + augmentation** on the fly.

```
from tensorflow.keras.preprocessing.image import ImageDataGenerator

datagen = ImageDataGenerator(
    rescale=1./255,      # Normalize pixel values (0-1)
    validation_split=0.2 # Split into training/validation
)

train_data = datagen.flow_from_directory(
    'dataset/train',
    target_size=(224, 224),
    batch_size=32,
    class_mode='categorical',
    subset='training'
)

val_data = datagen.flow_from_directory(
    'dataset/train',
    target_size=(224, 224),
    batch_size=32,
    class_mode='categorical',
    subset='validation'
)
```

✓ Key features:

- Automatically labels data based on folder names.
- Supports augmentation and normalization.
- Efficiently loads batches of data during training.

2. Using `tf.keras.utils.image_dataset_from_directory`

This is the **new recommended API** for TensorFlow ≥ 2.5 .

```
from tensorflow.keras.utils import image_dataset_from_directory

train_ds = image_dataset_from_directory(
    "dataset/train",
    validation_split=0.2,
    subset="training",
    seed=123,
    image_size=(224, 224),
    batch_size=32
)

val_ds = image_dataset_from_directory(
    "dataset/train",
    validation_split=0.2,
    subset="validation",
    seed=123,
    image_size=(224, 224),
    batch_size=32
)
```

✔ Advantages:

- Automatically infers labels from folder names.
- Integrates directly with `tf.data` pipelines.
- Efficient prefetching, caching, shuffling built-in.

3. Loading from File Paths and Custom Pipelines

When your data isn't neatly arranged in directories, you can manually load it.

```
import tensorflow as tf
import pandas as pd

df = pd.read_csv('labels.csv') # contains columns: filepath, label

file_paths = df['filepath'].values
labels = df['label'].values

def process_image(path, label):
    img = tf.io.read_file(path)
    img = tf.image.decode_jpeg(img, channels=3)
    img = tf.image.resize(img, (224, 224))
    img = tf.cast(img, tf.float32) / 255.0
    return img, label

dataset = tf.data.Dataset.from_tensor_slices((file_paths, labels))
dataset = dataset.map(process_image).batch(32).shuffle(1000)
```

✓ Used when:

- Data is custom-organized (not directory-based)
- You have metadata (e.g., CSV of labels)
- You want full control over transformations

4. Loading from NumPy Arrays or Pandas DataFrames

If data is already loaded in memory.

```
X_train = np.load('x_train.npy')
y_train = np.load('y_train.npy')

dataset = tf.data.Dataset.from_tensor_slices((X_train, y_train)).batch(32)
```

✓ Common in:

- Preprocessed datasets
- Tabular or numerical features

Data Augmentation

Data augmentation increases dataset **diversity artificially** — so the model generalizes better and doesn't overfit.

✓ Purpose:

- Teach model invariance to transformations (rotation, brightness, etc.)
 - Make small datasets appear “bigger”
 - Improve robustness
-

A. Using ImageDataGenerator

```
datagen = ImageDataGenerator(  
    rescale=1./255,  
    rotation_range=20,  
    width_shift_range=0.2,  
    height_shift_range=0.2,  
    shear_range=0.15,  
    zoom_range=0.2,  
    horizontal_flip=True,  
    fill_mode='nearest'  
)
```

B. Using `tf.keras.layers` (preferred modern method)

Create augmentations directly in the model (GPU accelerated).

```
from tensorflow.keras import layers, Sequential

data_augmentation = Sequential([
    layers.RandomFlip("horizontal"),
    layers.RandomRotation(0.2),
    layers.RandomZoom(0.2),
    layers.RandomBrightness(factor=0.2)
])
Then include it in your model:
model = Sequential([
    data_augmentation,
    layers.Conv2D(32, 3, activation='relu', input_shape=(224,224,3)),
    ...
])
```

✓ Advantages:

- Fast (runs on GPU)
- Works during training automatically
- Integrates seamlessly with datasets

Other Important Preprocessing Steps

1. Normalization / Standardization

Ensures data values are in the same scale for stable training.

- **Image data:** divide by 255 → range [0, 1]
 - **Numeric data:** standardize using mean and std
-

2. One-Hot or Label Encoding

Convert categorical labels into numeric form.

```
from tensorflow.keras.utils import to_categorical  
  
y_train = to_categorical(y_train, num_classes=10)
```

3. Balancing Classes

If your dataset is imbalanced (e.g., 90% class A, 10% class B):

- Use **class weights**
- Or **augment minority classes more**

```
class_weights = {0: 1.0, 1: 5.0}  
model.fit(train_ds, class_weight=class_weights)
```

Visualization

Visualization ensures that:

- Data is loaded correctly
- Labels are matched to the right images
- Augmentation works as expected

A. Visualize 5-10 Samples from Each Class

If using image_dataset_from_directory:

```
import matplotlib.pyplot as plt
import numpy as np

class_names = train_ds.class_names

plt.figure(figsize=(10, 10))
for images, labels in train_ds.take(1): # Take one batch
    for i in range(10):
        ax = plt.subplot(2, 5, i + 1)
        plt.imshow(images[i].numpy().astype("uint8"))
        plt.title(class_names[labels[i]])
        plt.axis("off")
```

- ✓ This shows 10 random images with their labels — perfect for sanity checking your dataset.

If using ImageDataGenerator

```
from tensorflow.keras.preprocessing.image import ImageDataGenerator
import matplotlib.pyplot as plt
import numpy as np

# Create the generator
datagen = ImageDataGenerator(rescale=1./255)

train_gen = datagen.flow_from_directory(
    'dataset/train',
    target_size=(224, 224),
    batch_size=32,
    class_mode='categorical'
)

# Get one batch
images, labels = next(train_gen)

# Plot first 10 images
plt.figure(figsize=(10, 10))
for i in range(10):
    ax = plt.subplot(2, 5, i + 1)
    plt.imshow(images[i])
    label_index = np.argmax(labels[i])
    label_name = list(train_gen.class_indices.keys())[label_index]
    plt.title(label_name)
    plt.axis("off")
plt.show()
```

`flow_from_directory()` automatically creates `class_indices` mapping: { 'cats': 0, 'dogs': 1 }

`next(train_gen)` gives one batch of 32 images + their one-hot encoded labels.

You take 10 samples and display them with their label names.

B. Visualize After Augmentation

```
for images, _ in train_ds.take(1):
    augmented_images = data_augmentation(images)
    plt.figure(figsize=(10,10))
    for i in range(9):
        ax = plt.subplot(3, 3, i + 1)
        plt.imshow(augmented_images[i].numpy().astype("uint8"))
        plt.axis("off")
```

Perceptron (Single Neuron Model)

The **Perceptron** is the **simplest neural network**, proposed by **Frank Rosenblatt (1958)**. It's basically **one neuron** that takes several inputs, applies weights, and produces one output.

$$y = f(w_1x_1 + w_2x_2 + \dots + w_nx_n + b)$$

Where:

- x_i : input features
- w_i : learned weights
- b : bias
- f : activation function (e.g., step, sigmoid, ReLU)

Working Mechanism

1. Each input x_i is multiplied by its weight w_i .
2. The weighted sum + bias is computed.
3. The activation function decides the output (e.g., 0 or 1).

Example (Binary Classification)

Let's classify whether a point lies **above or below a line**:

$$y = f(w_1x_1 + w_2x_2 + b)$$

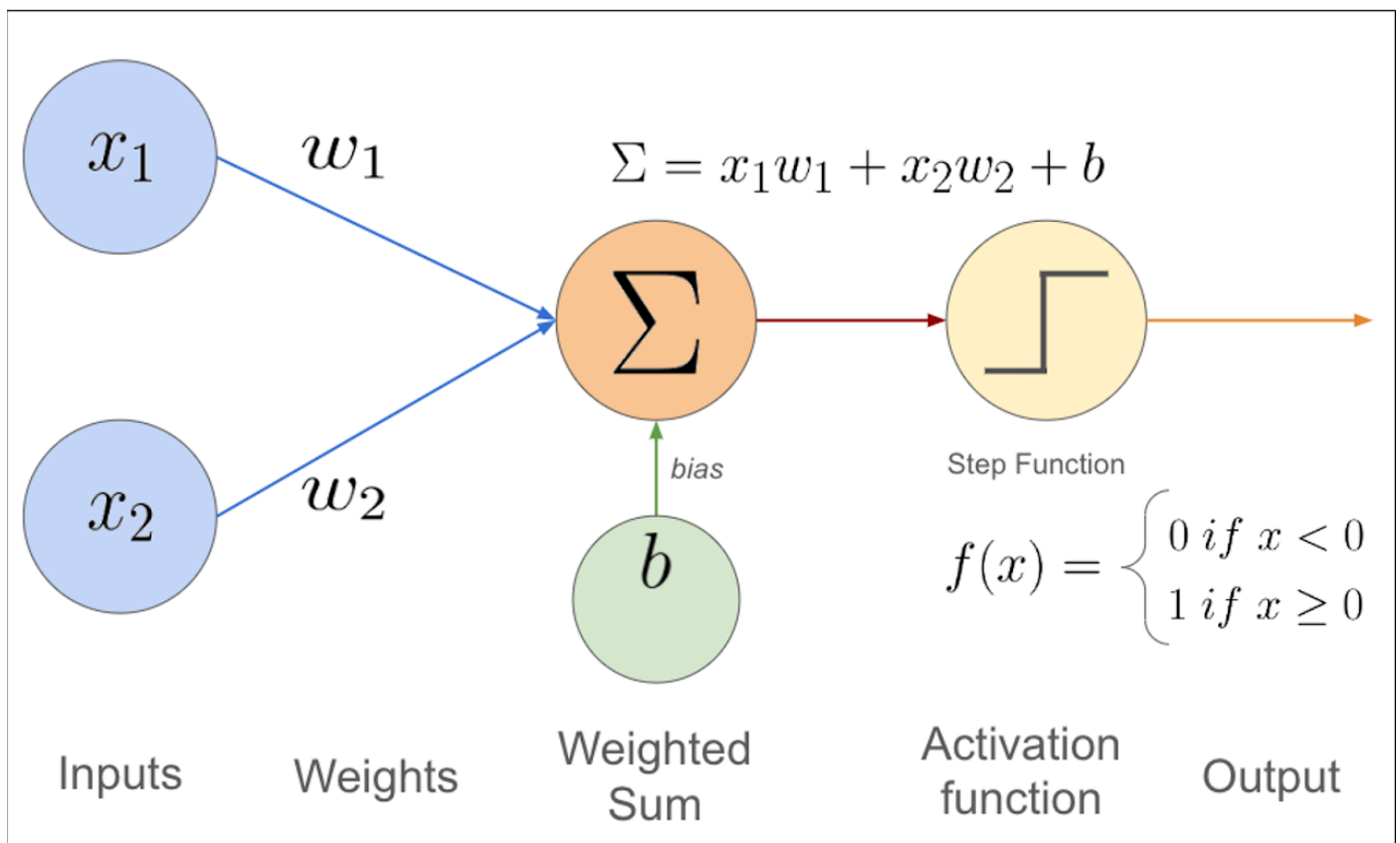
If result $> 0 \rightarrow$ Class 1

Else \rightarrow Class 0

So a **single perceptron** can only learn **linearly separable data** (a straight-line boundary in 2D).

Limitation

- Cannot handle **nonlinear problems** like XOR.
- Only outputs a single value (usually binary).
- Only one layer → very limited capacity.



ANN (Artificial Neural Network)

An **Artificial Neural Network (ANN)** is the **general term** for *any computational model* that mimics how biological neurons work — by passing weighted signals through layers of interconnected “neurons.”

It’s the **broad family name** that includes **all types of neural networks**, no matter how simple or complex.

Structure (Common to All ANNs)

All neural networks share these basic components:

1. **Input layer** — receives the data
2. **Hidden layers** — process information (apply weights, activations)
3. **Output layer** — produces predictions
4. **Weights & biases** — learned during training
5. **Activation functions** — introduce nonlinearity

ANN = The Family of Neural Architectures

Category	Model	Description
Basic ANN	MLP (Multi-Layer Perceptron)	Fully connected layers, general-purpose network
Deep ANN	DNN (Deep Neural Network)	Same as MLP but with many hidden layers
Vision-based ANN	CNN (Convolutional Neural Network)	Specialized for images using convolutional layers
Sequence-based ANN	RNN (Recurrent Neural Network)	Processes sequences (text, time series) step by step
Improved RNNs	LSTM / GRU	Handle long-term dependencies in sequences
Attention-based ANN	Transformer	Uses attention mechanism for global context
Generative ANN	Autoencoder / GAN / Diffusion	Learn to generate or reconstruct data
Graph-based ANN	GNN (Graph Neural Network)	Operates on graph structures (nodes + edges)

Multi-Layer Perceptron (MLP)

A **Multi-Layer Perceptron (MLP)** is an ANN with **multiple perceptrons stacked in layers**.

Structure:

Input layer → Hidden layer(s) → Output layer

Each neuron in one layer connects to **every neuron in the next layer** → that's why it's also called a **Fully Connected Neural Network (FCNN)** or **Dense Network**.

Mathematically

For one hidden layer:

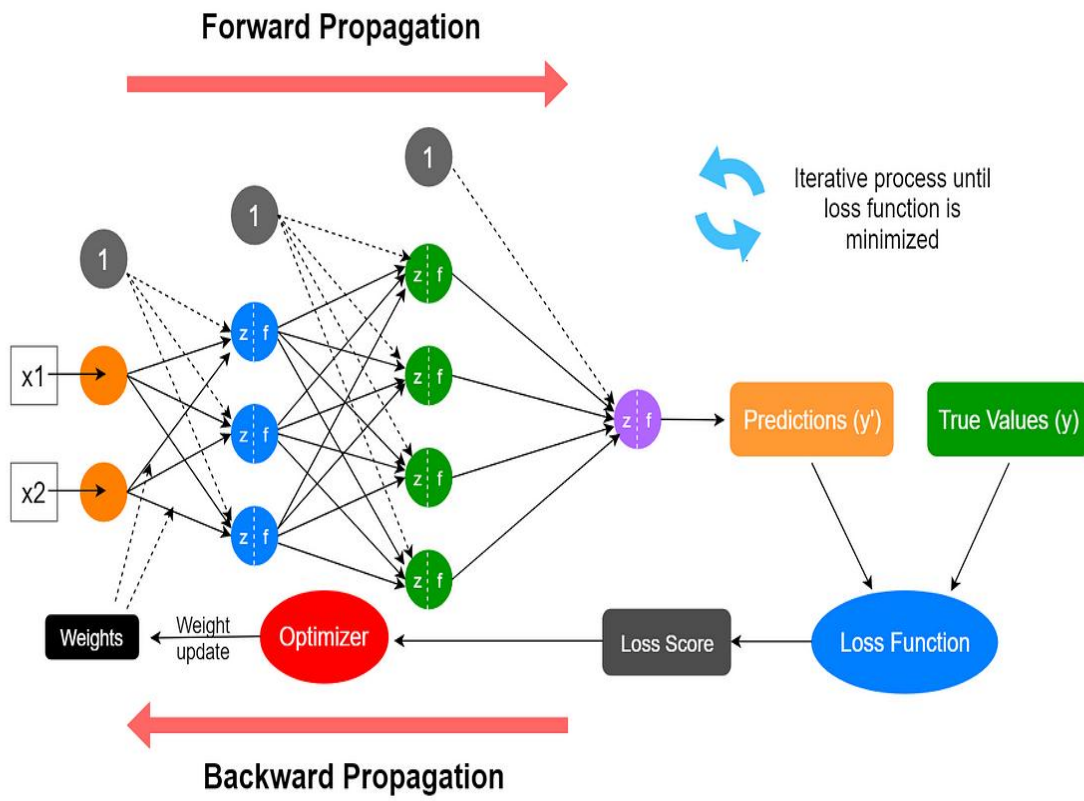
$$h = f(W_1x + b_1)$$
$$y = g(W_2h + b_2)$$

Where:

- f, g : activation functions (ReLU, sigmoid, etc.)
 - W_1, W_2 : weight matrices
 - h : hidden layer output
-

Advantages

- Can learn **nonlinear relationships** using activation functions.
- Multiple hidden layers allow the model to learn **complex patterns**.
- The foundation of **ANNs and DNNs**.



Deep Neural Network (DNN)

A DNN is an ANN with **many hidden layers** — sometimes tens or hundreds.

The key idea:

Each layer learns **higher-level abstractions** from the previous one.

For example, in an image:

- Layer 1 → detects edges
 - Layer 2 → detects shapes
 - Layer 3 → detects objects
 - Layer 4 → detects context (e.g., “cat on bed”)
-

Representation Learning

DNNs automatically learn **feature representations** from raw data, unlike traditional machine learning that needs **manual feature engineering**.

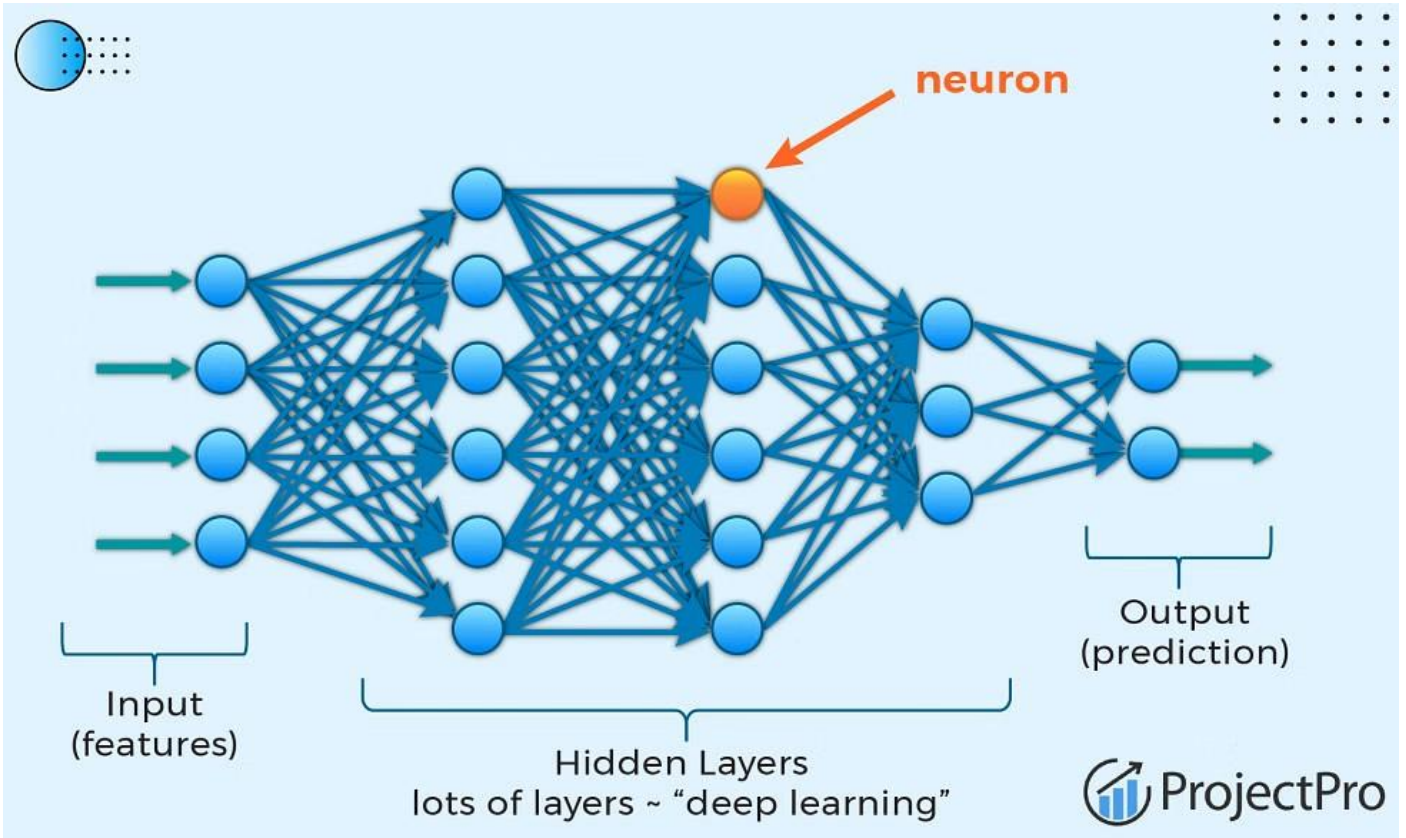
This is what made deep learning revolutionary — it can find patterns humans don’t explicitly program.

Architecture Variations

DNNs can have various architectures depending on the task:

- **Fully Connected (Feedforward)** — classic DNN for classification/regression.
- **Autoencoders** — for dimensionality reduction or denoising.
- **Recurrent Networks (RNNs, LSTMs)** — for sequences or time series.
- **Convolutional Networks (CNNs)** — for images.
- **Transformers** — for text, language, and now multimodal tasks.

All of these are **forms of DNNs** — they just have specialized structures.



Transfer Learning

What Is Transfer Learning?

Using knowledge learned from one task to help solve a different but related task.

Instead of training from scratch, we **reuse a pretrained model**.

Why Do We Need Transfer Learning?

Training deep neural networks from scratch requires:

- Huge datasets
- Massive compute
- Long training time

Example:

Training a large CNN from scratch may need **millions of images**.

Most real-world projects have:

- Small datasets
- Limited GPU power

Transfer learning solves this.

The Core Idea (How It Works)

A pretrained network has already learned:

- Edges
- Shapes
- Patterns
- High-level features

These features are useful across many tasks.

So instead of learning them again, we reuse them.

What Does a Pretrained Model Contain?

Think of layers:

Early Layers:

Learn **general features**

- Edges
- Colors
- Textures

These are reusable for almost any vision task.

Middle Layers:

Learn:

- Shapes
- Object parts

Still fairly reusable.

Final Layers:

Learn **task-specific patterns**.

These are usually replaced.

Types of Transfer Learning

Feature Extraction

Idea:

Freeze most layers and use them as a feature extractor.

Steps:

1. Load pretrained model
2. Freeze its weights
3. Replace the final layer
4. Train only the new layer

When to Use:

- Small dataset
 - Similar task
 - Limited compute
-

Fine-Tuning

Idea:

Allow some pretrained layers to update during training.

Steps:

1. Load pretrained model
2. Replace output layer
3. Unfreeze top layers
4. Train with a small learning rate

When to Use:

- Larger dataset
- Target task different from original
- Want higher accuracy

When Does Transfer Learning Work Best?

It depends on **task similarity**.

Case 1: Very Similar Tasks

Example:

- General object detection → cat vs dog classifier

Best strategy:

Freeze most layers

Case 2: Moderately Similar Tasks

Example:

- ImageNet → medical image classification

Best strategy:

Fine-tune upper layers

Case 3: Very Different Tasks

Example:

- Natural images → satellite radar images

Best strategy:

Fine-tune many layers

Benefits of Transfer Learning

Faster Training

No need to learn basic features from scratch.

Requires Less Data

Works well with small datasets.

Better Performance

Pretrained features are strong.

Lower Compute Cost

Transfer Learning Workflow

1. Typical steps:
 2. Choose pretrained model.
 3. Remove final layer.
 4. Add new output layer.
 5. Freeze base layers.
 6. Train new layers.
 7. Optionally fine-tune.
-

When NOT to Use Transfer Learning

Avoid when:

- Dataset is huge
- Task is completely unrelated
- Custom architecture needed

Key Hyperparameters in Transfer Learning

Important ones:

- Learning rate (usually smaller)
 - Number of frozen layers
 - Batch size
 - Regularization strength
-

Learning Rate Rule

Always use:

Smaller learning rate when fine-tuning

Because large updates can destroy pretrained knowledge.

TensorFlow

1. The Model (The Container)

The model is the object that holds everything together. There are two main ways to build one.

A. Sequential API (Easiest)

Use this for a simple stack of layers where each layer has exactly one input and one output.

Attributes & Methods:

- **.add(layer)**: Appends a layer to the stack.
- **.summary()**: Prints a useful description of the model architecture and parameter count.
- **.compile(...)**: Configures the model for training (sets optimizer, loss, metrics).
- **.fit(x, y, ...)**: Trains the model for a fixed number of epochs.

```
model = tf.keras.Sequential([
    # Define layers here or use .add() later
])
```

B. Functional API (Most Flexible)

Use this for complex architectures (e.g., ResNet, multiple inputs/outputs). You treat layers like functions.

```
inputs = tf.keras.Input(shape=(28, 28))
x = tf.keras.layers.Flatten()(inputs) # Layer is called on the tensor
outputs = tf.keras.layers.Dense(10)(x)
model = tf.keras.Model(inputs=inputs, outputs=outputs)
```

2. Layers (The Building Blocks)

Layers contain the weights and biases. They live in `tf.keras.layers`.

Layer Type	Python Class	Key Arguments (Attributes)	What It Does
Dense	<code>Dense()</code>	<p><code>units (int)</code>: Number of neurons.</p> <p><code>activation (str)</code>: Activation function.</p>	Standard fully connected layer. $y = \text{activation}(Wx + b)$
Convolution	<code>Conv2D()</code>	<p><code>filters (int)</code>: Number of output filters.</p> <p><code>kernel_size (tuple)</code>: Window size (e.g., (3,3)).</p> <p><code>strides</code>: Step size.</p>	Spatial filtering. Used for Images/Computer Vision.
Flatten	<code>Flatten()</code>	<i>None usually.</i>	Flattens 2D/3D inputs into 1D vector. Bridge between Conv2D and Dense.
Dropout	<code>Dropout()</code>	<code>rate (float 0.0-1.0)</code> : % of neurons to drop.	Prevents overfitting by randomly setting inputs to 0 during training.
Input	<code>InputLayer()</code>	<code>input_shape (tuple)</code> : Shape of data (e.g., (28,28)).	Explicitly defines the entry point shape.

```
tf.keras.layers.Dense(units=64, activation='relu', name='hidden_layer_1')
```

3. Activation Functions (The Non-Linearity)

These decide "how much" a neuron fires. You can pass them as strings to layers or use them as standalone layers.

Activation	String Alias	Class Object	Best Used For
ReLU	'relu'	tf.keras.activations.ReLU()	Hidden Layers. Default choice. Fast and effective.
Sigmoid	'sigmoid'	tf.keras.activations.Sigmoid()	Output Layer (Binary). Squashes output between 0 and 1.
Softmax	'softmax'	tf.keras.activations.Softmax()	Output Layer (Multi-class). Probabilities summing to 1.
Linear	'linear'	<i>None (default)</i>	Output Layer (Regression). No activation; raw values.

4. Loss Functions (The Goal)

The measure of error the model tries to minimize.

- **Key Argument:** from_logits=True
 - Use this if your model output layer has **no activation**. It is more numerically stable.
 - If your model output has softmax or sigmoid, use from_logits=False.

Common Choices:

- 'binary_crossentropy' (Binary Classification)
- 'categorical_crossentropy' (Multi-class, One-hot labels)
- 'sparse_categorical_crossentropy' (Multi-class, Integer labels)
- 'mse' (Regression)

5. Optimizers (The Driver)

The algorithm that updates the weights based on gradients.

Optimizer	Class	Key Arguments	When to use
Adam	tf.keras.optimizers.Adam	learning_rate (default 0.001)	Just use this. It's the standard "good enough" default for most problems.
SGD	tf.keras.optimizers.SGD	learning_rate, momentum	specific research cases or when Adam fails.
RMSprop	tf.keras.optimizers.RMSprop	learning_rate	Good for Recurrent Neural Networks (RNNs).

```
opt = tf.keras.optimizers.Adam(learning_rate=0.01) # Tuning the speed of learning
model.compile(optimizer=opt, ...)
```

6. Metrics (The Scoreboard)

Used to judge performance *for humans*. They do **not** affect training (gradients are not calculated for metrics).

- **Regression:** 'mae', 'mse'
- **Classification:** 'accuracy', Precision(), Recall(), AUC()

```
model.compile(..., metrics=['accuracy', tf.keras.metrics.Precision()])
```

7. Putting It All Together: A "Template" Code

```
import tensorflow as tf

# 1. Define the Model
model = tf.keras.Sequential([
    tf.keras.layers.InputLayer(input_shape=(10,)), # 10 input features
    tf.keras.layers.Dense(64, activation='relu'), # Hidden layer
    tf.keras.layers.Dropout(0.2), # Regularization
    tf.keras.layers.Dense(32, activation='relu'),
    tf.keras.layers.Dense(1) # Output (Linear for regression, or add activation)
])

# 2. Compile (Configure)
# For binary classification (Fraud Detection)
model.compile(
    optimizer=tf.keras.optimizers.Adam(learning_rate=0.001),
    loss=tf.keras.losses.BinaryCrossentropy(from_logits=True), # Stable!
    metrics=['accuracy']
)

# 3. Fit (Train)
history = model.fit(
    x=X_train,
    y=y_train,
    epochs=20, # How many times to see the whole dataset
    batch_size=32, # How many samples per gradient update
    validation_split=0.2 # Use 20% of data to check performance during training
)
```

8. Essential "Ease of Use" Features

Callbacks:

These are tools that run *during* training.

- EarlyStopping: Stops training if validation loss stops improving (saves time).
- ModelCheckpoint: Saves the best version of your model automatically.

```
callback = tf.keras.callbacks.EarlyStopping(monitor='val_loss', patience=3)
model.fit(..., callbacks=[callback])
```

```
# 1. Define where you want to save the model
# Using .keras is the modern standard (replaces .h5)
checkpoint_path = "best_model.keras"

# 2. Create the Callback
checkpoint_callback = tf.keras.callbacks.ModelCheckpoint(
    filepath=checkpoint_path,
    monitor='val_loss',    # Metric to watch (e.g., 'val_loss', 'val_accuracy')
    save_best_only=True,   # KEY: Only save if this metric improves
    save_weights_only=False, # False = save entire model (arch + weights + optimizer)
    mode='min',           # 'min' for loss, 'max' for accuracy/AUC
    verbose=1            # 1 = print a message when saving
)

# 3. Train with the callback
# Assuming you have X_train, y_train, etc. ready
history = model.fit(
    X_train, y_train,
    epochs=50,
    validation_data=(X_val, y_val),
    callbacks=[checkpoint_callback] # <--- Add it here!
)

# 4. Load the best model afterwards
# (Since the last epoch might not be the best one)
best_model = tf.keras.models.load_model(checkpoint_path)
```

Model Saving/Loading:

```
model.save('my_model.keras') # Saves architecture, weights, and optimizer state.
tf.keras.models.load_model('my_model.keras') # Reloads it exactly as it was.
```

Conclusion

In conclusion, deep learning stands as a cornerstone of modern AI, pushing the boundaries of what machines can learn and achieve. By leveraging large datasets, high computational power, and layered neural architectures, deep learning models excel at recognizing patterns, understanding context, and generating intelligent predictions. This document highlights not only the technical foundations but also the significance of deep learning in shaping current and future innovations—from autonomous vehicles to language models. Understanding these concepts provides the groundwork for applying deep learning effectively across various domains and for advancing toward more intelligent, adaptive systems.

Feedback & Contribution

This is a living document. If you have feedback, suggestions, or additional insights that could improve it, I welcome your input. Sharing ideas helps keep this material accurate, relevant, and valuable for everyone.

Copyright & Usage

© 2025 [Youssef Amgad Elkhatib].

This document is intended **for learning and reference purposes only**.

You are free to read, use, and share the content for personal or educational use, but:

- **Do not copy or republish this document as your own.**
- Always provide proper credit when referencing.
- Commercial use or redistribution without permission is not allowed.

You may not reproduce this document in whole or in part without attribution. Suggestions and feedback are welcome, and contributors will be acknowledged, but copyright remains with the author.

Acknowledgments

This document is authored and copyrighted by **[Youssef Amgad Elkhatib]**.

Special thanks to the contributors who provided valuable feedback and suggestions for improvements:

Name	Contribution