

# Neural Network Training

Neural networks are universal function approximators capable of learning complex, non-linear relationships from data. Unlike linear models or tree-based algorithms, neural networks build hierarchical representations through layers of interconnected neurons, enabling them to excel on large datasets with intricate patterns.

However, this flexibility comes with complexity. Neural networks have many hyperparameters — layer counts, neuron counts, activation functions, learning rates, regularization strategies — all of which must be configured appropriately. This guide provides practical, actionable guidance for training neural networks on structured tabular data, balancing performance, generalization, and computational efficiency.

1. Network Architecture .....	3
1.1 Hidden Layer Count.....	3
1.2 Hidden Unit Count (Neurons per Layer).....	3
2. Activation Functions .....	4
2.1 Hidden Layer Activation Functions.....	5
2.2 Output Layer Activation Functions .....	6
3. Weight Initialization .....	7
3.1 He / Kaiming Initialization .....	7
3.2 Xavier / Glorot Initialization .....	8
3.3 Uniform and Normal Initialization .....	8
4. Batch Normalization .....	8
5. Loss Functions.....	9
5.1 Regression Loss Functions.....	9
5.2 Classification Loss Functions.....	10
6. Regularization .....	11
6.1 Norm-Based Regularization (L1 / L2).....	11
6.2 Dropout Regularization .....	12
7. Optimization Techniques.....	13
7.1 Gradient Descent.....	13
7.2 Gradient Descent with Momentum.....	13
7.3 Adam Optimizer .....	14
8. Learning Rate Decay .....	14
8.1 Exponential Decay.....	15
8.2 Step-wise Decay.....	15
8.3 Loss-Deviation Decay (Adaptive).....	15
9. Convergence Criteria .....	15
9.1 Loss Threshold .....	16

---

9.2 Iteration Count (Fixed Epochs) .....	16
9.3 Loss Gradient Threshold .....	16
9.4 Early Stopping (Validation-Based) .....	16
10. Status Notification & Monitoring .....	17
10.1 What to Monitor .....	17
10.2 Logging Frequency .....	18
11. Troubleshooting Guide for Neural Networks .....	18
11.1 Model Not Learning (Loss Not Decreasing) .....	19
11.2 Overfitting (Train Accuracy >> Validation Accuracy) .....	19
11.3 Underfitting (Both Train and Validation Accuracy Low) .....	20
11.4 Training Instability (Loss Oscillating).....	20
11.5 Slow Convergence .....	20
12. Practical Recommendations.....	21
Architecture .....	21
Activation Functions .....	21
Weight Initialization.....	21
Regularization .....	21
Optimization .....	21
Loss Functions .....	21
Convergence & Monitoring .....	22
Preprocessing .....	22
13. Underfitting – Overfitting Trade-off .....	22

# 1. Network Architecture

The architecture of a neural network — the number and configuration of its layers and neurons — determines its capacity to learn patterns from data. Architecture design depends on the modeling objective, feature characteristics, dataset size, and computational constraints.

## 1.1 Hidden Layer Count

The number of hidden layers determines the depth of the network and directly influences its ability to learn complex, non-linear relationships.

Layer Count	Typical Use Case	Advantages	Considerations
1–3 layers	Structured tabular data	Fast training; stable; good baseline	Limited capacity for complex patterns
4–6 layers	High-dimensional or complex relationships	Better feature extraction; manageable complexity	Requires more data and tuning
> 6 layers	Unstructured data	Hierarchical pattern learning	Can be an overkill for tabular data; prone to overfitting

### Detailed Guidance

- **Shallow networks (1–3 layers):** Generally sufficient for structured tabular data where feature interactions are simpler and well-defined. A network with one hidden layer is theoretically capable of approximating any continuous function (Universal Approximation Theorem), but in practice, 2–3 layers offer a better balance between expressiveness and stable training. This range serves as the recommended baseline for most business analytics tasks.
- **Medium depth (4–6 layers):** More complex relationships arising from high-dimensional interactions or large datasets may justify deeper architectures. These networks provide improved hierarchical feature extraction while remaining computationally manageable.
- **Deep networks (> 6 layers):** Rarely applied to traditional structured data because they tend to overfit, offer diminishing performance gains, and require significantly more computational resources. Deep architectures are typically reserved for unstructured domains like image, text, and audio processing, where hierarchical patterns naturally benefit from greater depth.

#### Best Practice

For tabular data, start with 2–3 hidden layers. Only increase depth if you have > 100K training samples, high-dimensional features (100+ columns), and evidence that shallower networks are underfitting. Most business problems are well-served by shallow architectures.

## 1.2 Hidden Unit Count (Neurons per Layer)

The number of hidden units (neurons) in each layer controls the representational capacity of the network. This number generally depends on the number of input features, the number of output targets, and the amount of training data.

### Heuristics for Initial Sizing

Two commonly used starting points for selecting hidden unit count:

- Average of input and output features — e.g., if 50 inputs and 1 output, start with  $\approx 25$  neurons
- $2\times$  to  $10\times$  the number of input features — e.g., if 20 inputs, try 40–200 neurons

### Trade-offs

- **Too few units:** Underfitting — the model lacks capacity to learn underlying patterns
- **Too many units:** Increased computational cost and risk of overfitting without proper regularization

### Relationship Between Depth and Width

Depth and width are not independent choices:

- Deeper networks can use narrower layers (parameters distributed across depth)
- Shallow networks may need wider layers (parameters concentrated in fewer layers)
- Total parameter count matters more than specific configuration
- Example: 3 layers  $\times$  64 neurons  $\approx$  2 layers  $\times$  96 neurons (similar capacity)

#### Note

For structured data, a wide and shallow architecture (more neurons per layer, fewer layers) often performs better than a very deep network. Start with 2–3 layers of 64–128 neurons each, then adjust based on validation performance.

## 2. Activation Functions

Activation functions introduce non-linearity into the model. Without activation functions, a neural network would collapse into a single-layer linear model, regardless of how many hidden layers it contains.

Function	Primary Use	Strengths	Limitations
<b>ReLU</b>	Hidden layers (default)	Fast computation; effective gradient flow	Dying ReLU in deep networks
<b>Leaky ReLU</b>	Hidden layers (alternative)	Prevents dying neurons	Extra hyperparameter; minimal benefit for shallow nets
<b>Sigmoid</b>	Binary classification output	Outputs probabilities [0, 1]	Vanishing gradients in deep layers

Function	Primary Use	Strengths	Limitations
<b>ReLU</b>	Hidden layers (default)	Fast computation; effective gradient flow	Dying ReLU in deep networks
<b>Tanh</b>	Hidden/output (bounded targets)	Zero-centered; range [-1, 1]	Vanishing gradients
<b>Linear</b>	Regression output	Unbounded continuous predictions	No non-linearity (only for output)

## 2.1 Hidden Layer Activation Functions

### ReLU (Rectified Linear Unit)

ReLU is the most widely used activation function due to its simple gradient and low computational cost. It outputs the input directly when positive and zero otherwise:  $f(x) = \max(0, x)$ .

#### Advantages

- Fast computation — simple thresholding operation
- Effective gradient flow — no vanishing gradient problem for positive inputs
- Encourages sparse representations — many neurons output zero

#### Limitation: Dying ReLU

If a neuron consistently receives negative inputs, it outputs 0 forever and stops learning. However, this risk is limited in shallow networks and can be mitigated with proper weight initialization.

### Leaky ReLU

Leaky ReLU addresses the "dying ReLU" problem by allowing a small, non-zero gradient for negative inputs (typically  $0.01 \times \text{input}$ ):  $f(x) = \max(0.01x, x)$ .

#### When to Use

- Deep networks (> 6 layers) where dying ReLU is a concern
- When you observe many neurons with zero activations during training

#### Limitation

Introduces an additional hyperparameter (the negative slope) which needs tuning. The benefit over ReLU is marginal for shallow networks.

### Sigmoid and Tanh

Sigmoid maps inputs to  $[0, 1]$ , while Tanh maps to  $[-1, 1]$ . Both were historically popular but have largely been replaced by ReLU for hidden layers.

### Limitation: Vanishing Gradients

Both functions saturate at extreme values, producing very small gradients that slow or stall training in deep networks. They remain useful for output layers but are rarely used in hidden layers anymore.

#### Best Practice

Default recommendation: Use ReLU for all hidden layers. Only switch to Leaky ReLU if you observe widespread dying neurons or are training very deep networks (> 6 layers).

## 2.2 Output Layer Activation Functions

The output activation function determines how the final layer transforms its inputs into predictions. The choice depends entirely on your problem type. Using the wrong output activation can make your model unable to learn.

### Linear (No Activation)

Default choice for regression problems where you need to predict continuous values without bounds — temperature forecasting, stock prices, revenue prediction.

- Outputs: Unbounded real numbers ( $-\infty$  to  $+\infty$ )
- Use case: Any regression task without constraints

### ReLU (Output Layer)

Use for regression when outputs must be non-negative — house prices, age, counts, volumes.

- Outputs: Non-negative real numbers ( $\geq 0$ )
- Benefit: Encodes domain knowledge directly into the architecture

### Sigmoid

Standard choice for binary classification problems — customer churn, fraud detection, spam classification.

- Outputs: Probabilities in range  $[0, 1]$
- Interpretation: Probability of the positive class
- Decision threshold: Typically 0.5 (adjustable based on business costs)

### Softmax

Standard choice for multi-class classification — product categorization, sentiment analysis (positive/neutral/negative).

- Outputs: Probability distribution over all classes (sums to 1)
- Prediction: Class with highest probability

### Tanh (Output Layer)

Occasionally used for regression when outputs should be bounded between  $-1$  and  $1$  — normalized targets, certain time series problems.

### ⚠ Warning

Critical: Always match output activation to problem type:

- Regression → Linear or ReLU (if non-negative)
- Binary classification → Sigmoid
- Multi-class classification → Softmax

Mismatching output activation and loss function will prevent the model from learning.

## 3. Weight Initialization

Weight initialization defines the starting values for the network's parameters before training begins. Proper initialization is critical because initializing all weights identically causes symmetry — every neuron in a layer computes the same function and receives identical gradient updates, effectively reducing each layer to a single neuron.

Good initialization breaks symmetry and prevents gradients from becoming too large (exploding gradients) or too small (vanishing gradients) as they flow through the network.

Method	Best For	Advantages	Limitations
<b>He / Kaiming</b>	ReLU, Leaky ReLU	Prevents variance shrinkage in deep networks	Not optimal for sigmoid/tanh
<b>Xavier / Glorot</b>	Sigmoid, Tanh, Linear output	Controls variance; prevents activation saturation	Suboptimal for ReLU
<b>Uniform</b>	General / shallow networks	Simple; ensures symmetry breaking	Less principled than He/Xavier
<b>Normal</b>	Fallback / simple models	Standard initialization	Less effective than Gaussian specialized methods

### 3.1 He / Kaiming Initialization

Developed specifically for ReLU-based networks, He initialization scales weights using the number of input units in a layer. This prevents the variance of activations from shrinking or exploding as data flows through deep networks.

#### When to Use

- Hidden layers with ReLU or Leaky ReLU activations

- This is the recommended default for modern neural networks

### Why It Works

ReLU sets half of its inputs to zero, halving the variance of activations. He initialization compensates by drawing weights from a distribution with variance proportional to  $2/n_{in}$ , where  $n_{in}$  is the number of inputs to the layer.

## 3.2 Xavier / Glorot Initialization

Xavier initialization scales weights based on the average of the number of input and output units. It is effective for activation functions like Sigmoid and Tanh, which require careful control over activation variance because they saturate easily.

### When to Use

- Hidden layers with Sigmoid or Tanh activations
- Output layers with Linear activation (regression)

### Why It Works

Sigmoid and Tanh are sensitive to input magnitude — very large or small inputs push them into saturation regions where gradients vanish. Xavier initialization keeps activations in the responsive range.

## 3.3 Uniform and Normal Initialization

These are fallback methods that draw weights from uniform or normal distributions. They remain effective for shallow or moderately deep networks and ensure symmetry breaking.

### When to Use

- Simple models where specialized initialization is unnecessary
- As a baseline comparison during experimentation

### Best Practice

Recommended defaults:

- ReLU hidden layers → He initialization
- Sigmoid/Tanh hidden layers → Xavier initialization
- Linear output (regression) → Xavier initialization

If unsure, use He initialization — it works well for most cases and is the safest general-purpose choice.

## 4. Batch Normalization

Batch normalization normalizes the inputs of each layer by adjusting and scaling activations using batch statistics (mean and variance). This stabilizes the distribution of activations during training.

### Benefits

- Accelerates convergence by allowing higher learning rates
- Reduces sensitivity to weight initialization
- Acts as a regularizer — can reduce need for dropout
- Stabilizes training of deep networks

### When to Use

- Deep networks ( $\geq 4$  layers) where training instability is a concern
- When you want to train with higher learning rates for faster convergence

### Limitations

#### **Warning**

Batch normalization is incompatible with:

- Very small batch sizes ( $< 8$ ) — insufficient samples to compute stable statistics
- Stochastic gradient descent with batch size = 1 — no batch statistics available

For small batches, consider Layer Normalization or Group Normalization instead (not covered here).

#### **Note**

For shallow networks (2–3 layers) on tabular data, batch normalization is often unnecessary. Reserve it for deeper architectures or when you observe training instability.

## 5. Loss Functions

The loss function quantifies how far the model's predictions are from the actual target values. It directly determines what the optimizer minimizes during training. Choosing the right loss function is critical.

### 5.1 Regression Loss Functions

#### Mean Squared Error (MSE)

MSE is the default loss function for regression. It computes the mean of squared differences between predicted and actual values.

#### When to Use

- Default choice for most regression tasks

- When large errors are particularly undesirable (financial forecasting, demand estimation)

### Characteristics

- Penalizes large errors more heavily due to squaring
- Sensitive to outliers — extreme values disproportionately inflate the loss
- Smooth gradient makes optimization straightforward

## Mean Absolute Error (MAE)

MAE computes the mean of absolute differences between predictions and targets.

### When to Use

- Datasets with outliers or natural irregularities
- When robustness matters more than penalizing large mistakes

### Characteristics

- Treats all errors equally without amplifying extreme values
- More interpretable — same units as target variable
- Gradient is constant, which can slow convergence near the optimum

#### Best Practice

Start with MSE for regression. If residual analysis reveals significant outliers that are legitimate (not errors), switch to MAE to reduce their influence on training.

## 5.2 Classification Loss Functions

### Log Loss (Binary Cross-Entropy)

Log loss measures the accuracy of probability predictions for binary classification. It heavily penalizes confident incorrect predictions.

### When to Use

- All binary classification problems (paired with Sigmoid output)
- When probability calibration matters (fraud detection, medical diagnosis)

### Characteristics

- Encourages the model to output well-calibrated probabilities
- Exponentially penalizes predictions far from the true label
- Works on the principle of maximum likelihood estimation

### Categorical Cross-Entropy

Categorical cross-entropy extends log loss to multi-class problems. It compares the predicted probability distribution against the true one-hot encoded label.

### When to Use

- All multi-class classification problems (paired with Softmax output)
- Requires one-hot encoded target labels

#### **Warning**

Critical: Match loss function to problem type and output activation:

- Regression + Linear/ReLU output → MSE or MAE
- Binary classification + Sigmoid output → Log Loss (Binary Cross-Entropy)
- Multi-class + Softmax output → Categorical Cross-Entropy

Mismatching these will prevent the model from learning correctly.

## 6. Regularization

Regularization prevents neural networks from overfitting by adding constraints that encourage the model to generalize better to unseen data. This is critical for production reliability — a model that overfits training data may perform well in testing but fail catastrophically when deployed on new customers, market conditions, or seasonal patterns.

### 6.1 Norm-Based Regularization (L1 / L2)

#### L1 Regularization (Lasso)

L1 adds the sum of absolute values of weights to the loss function. This encourages sparsity — many weights are driven to exactly zero.

#### When to Use

- When you suspect many input features are irrelevant
- When you need a compressed model for faster inference

#### Lambda Range

Typical values: 0.001 to 1.0

- Small lambda ( $\approx 0.001$ ): Minimal regularization
- Medium lambda ( $\approx 0.01 - 0.1$ ): Balanced sparsity
- Large lambda ( $\approx 1.0$ ): Aggressive sparsity; many weights eliminated

#### L2 Regularization (Ridge / Weight Decay)

L2 adds the sum of squared weights to the loss function. This penalizes large weights without forcing any to zero.

#### When to Use

- Most common regularization choice for neural networks
- When you want to prevent any single neuron from dominating

### Lambda Range

Typical values: 0.0001 to 0.1

- Small lambda ( $\approx 0.0001$ ): Light regularization
- Medium lambda ( $\approx 0.001 - 0.01$ ): Standard choice
- Large lambda ( $\approx 0.1$ ): Strong weight shrinkage

#### Best Practice

For neural networks, L2 regularization (weight decay) is the default. Use lambda  $\approx 0.001 - 0.01$  as a starting point. L1 is less common but useful when you want automatic feature selection through sparsity.

## 6.2 Dropout Regularization

Dropout reduces overfitting by randomly disabling a fraction of neurons during each training step. This prevents neurons from co-adapting too strongly and compels the network to learn more robust, redundant representations.

### How It Works

During training, each neuron has a probability  $p$  of being temporarily "dropped" (set to zero). The remaining neurons must compensate, forcing the network to learn features that work in multiple contexts.

### Dropout Probability (Dropout Rate)

Typical values: 0.2 to 0.5 (20% to 50% of neurons dropped)

- Low dropout (0.1 – 0.2): Gentle regularization; suitable for moderate overfitting
- Medium dropout (0.3 – 0.4): Standard choice for most networks
- High dropout (0.5+): Aggressive regularization; use when severe overfitting is observed

### Where to Apply

- Typically applied to hidden layers, not the output layer
- Can be applied to different layers with different rates
- Often paired with L2 regularization for stronger effect

#### Warning

Do NOT use dropout during inference (testing/production). Dropout is only applied during training. During inference, all neurons are active but their outputs are scaled by  $(1 - \text{dropout\_rate})$  to compensate.

## 7. Optimization Techniques

Optimization methods govern how model parameters (weights) are updated to minimize the loss function. The choice of optimizer significantly affects convergence speed, stability, and final model quality.

Optimizer	Hyperparameters	Strengths	When to Use
<b>Gradient Descent</b>	Alpha: 0.001–0.1	Simple; predictable	Sensitive to learning rate; can be slow
<b>Momentum</b>	Alpha: 0.001–0.1 Beta: 0.8–0.99	Faster convergence; smooths oscillations	One extra hyperparameter
<b>Adam</b>	Alpha: 0.001–0.1 Beta1: 0.9 Beta2: 0.999	Adaptive rates; robust default for deep learning	Slight computational overhead; can overshoot

### 7.1 Gradient Descent

Gradient descent is the fundamental optimization method. At each iteration, it computes the gradient of the loss with respect to the parameters and moves in the direction of steepest descent.

#### Learning Rate (Alpha)

Typical range: 0.001 to 0.1

- Controls the step size for each parameter update
- Small alpha ( $\approx 0.001$ ): Slow, stable convergence; suitable for noisy gradients
- Large alpha ( $\approx 0.1$ ): Fast initial progress; risk of overshooting the minimum

#### Characteristics

- Simple and predictable
- Sensitive to learning rate — requires tuning
- Can be slow on complex loss surfaces with plateaus or valleys

### 7.2 Gradient Descent with Momentum

Momentum improves gradient descent by incorporating information from previous update steps. This accelerates convergence in consistent directions and reduces oscillations.

#### Hyperparameters

- **Alpha (Learning Rate):** 0.001 to 0.1
- **Beta (Momentum Coefficient):** 0.8 to 0.99
  - Higher beta ( $\approx 0.99$ ): Smooth, long-term trajectory; slower reaction to changes
  - Lower beta ( $\approx 0.8$ ): Faster response to new gradient directions

#### Characteristics

- Faster convergence than vanilla gradient descent on most problems
- Helps escape shallow local minima and plateaus
- Adds one additional hyperparameter to tune

## 7.3 Adam Optimizer

Adam (Adaptive Moment Estimation) is the most popular optimizer in deep learning. It computes adaptive learning rates for each parameter using moving averages of both gradients and squared gradients.

### Hyperparameters

- Alpha (Learning Rate): 0.001 to 0.1 — typically start with 0.001
- Beta1 (First Moment): 0.9 — controls momentum
- Beta2 (Second Moment): 0.999 — stabilizes updates using gradient variance

### Characteristics

- Adaptive per-parameter learning rates — minimal manual tuning needed
- Robust default choice across a wide range of architectures and datasets
- Combines benefits of momentum and adaptive learning rates
- Slight computational overhead compared to basic gradient descent

#### Best Practice

Recommended default: Start with Adam optimizer (alpha = 0.001, beta1 = 0.9, beta2 = 0.999). This works well for most neural networks without extensive tuning. Only switch to simpler optimizers if you have specific reasons (e.g., Adam is overshooting, or you want more control).

## 8. Learning Rate Decay

Learning rate decay gradually reduces the learning rate during training, helping the model transition from coarse adjustments early on to fine-tuning in later epochs.

### When Learning Rate Decay Is Unnecessary

Learning rate decay is often not needed for:

- Small datasets (< 10K samples)
- Transfer learning tasks
- Simple models (2–3 layers)
- When using Adam optimizer (it adapts learning rates automatically)

### When to Use Learning Rate Decay

- Large datasets with long training runs (100+ epochs)
- When validation loss plateaus but training loss continues to decrease
- When fine-tuning near convergence is important

## 8.1 Exponential Decay

Exponential decay multiplies the learning rate by a constant factor less than 1 after each epoch.

### Characteristics

- Smooth, continuous reduction
- Useful for long training runs requiring gradual refinement
- Formula:  $\alpha_{\text{new}} = \alpha_{\text{initial}} \times \text{decay\_rate}^{\text{epoch}}$

## 8.2 Step-wise Decay

Step-wise decay reduces the learning rate by a fixed factor at predetermined intervals (e.g., every 30 epochs, reduce by 50%).

### Characteristics

- Simple and predictable — easy to configure
- Common schedule: drop by 10× at 50%, 75% of total epochs
- Works well when training duration is known in advance

## 8.3 Loss-Deviation Decay (Adaptive)

Loss-deviation decay monitors validation loss and reduces the learning rate only when progress stalls.

### Characteristics

- Data-driven approach — adapts to actual training dynamics
- Reduces learning rate by a fixed factor when validation loss stops improving for N consecutive epochs
- Requires validation set

### Best Practice

Most neural networks on tabular data don't need learning rate decay, especially when using Adam. If validation loss plateaus, try loss-deviation decay (reduce-on-plateau) as it automatically responds to training dynamics without requiring manual scheduling.

## 9. Convergence Criteria

Convergence criteria determine when to stop training. Proper stopping conditions prevent underfitting (stopping too early), overfitting (training too long), and wasted computational resources.

## 9.1 Loss Threshold

Training stops when the training loss reaches or drops below a predefined target value.

### When to Use

- When the application requires a specific accuracy level
- When you have a clear target performance metric

### Limitations

- Requires knowing an appropriate threshold in advance
- May stop prematurely if threshold is set too loosely
- Does not account for overfitting — training loss can be low while validation loss increases

## 9.2 Iteration Count (Fixed Epochs)

Training runs for a fixed number of epochs (passes through the dataset).

### Characteristics

- Simplest and most widely used criterion
- Guarantees predictable training time
- Does not guarantee convergence on its own
- Often combined with other criteria (especially early stopping)

### Typical Values

- Small datasets: 50–200 epochs
- Medium datasets: 100–500 epochs
- Large datasets: 10–100 epochs (each epoch covers more data)

## 9.3 Loss Gradient Threshold

Training stops when the improvement in loss between successive epochs falls below a threshold.

### When to Use

- When you want to stop automatically once learning plateaus
- Prevents unnecessary training when further learning yields negligible improvement

### Typical Threshold

Epsilon values: 1e-6 to 1e-4

- Stop if  $|\text{loss\_epoch\_t} - \text{loss\_epoch\_}{t-1}| < \text{epsilon}$

## 9.4 Early Stopping (Validation-Based)

Training halts when validation performance stops improving for a predetermined number of consecutive epochs (patience). This is the most important convergence criterion for neural networks.

### How It Works

- Monitor validation loss during training
- If validation loss does not improve for N consecutive epochs, stop training
- Typical patience values: 10–50 epochs

### Benefits

- Prevents overfitting by stopping before the model memorizes training noise
- Saves computational resources by avoiding unnecessary epochs
- Automatically balances underfitting and overfitting

### Requirements

- Requires a held-out validation set (typically 10–20% of training data)
- Adds slight overhead to compute validation loss each epoch

#### Best Practice

Always use early stopping for neural networks. Set a high maximum epoch count (e.g., 500–1000) and let early stopping find the optimal point. Use patience = 20–50 to avoid stopping prematurely due to temporary validation loss fluctuations.

## 10. Status Notification & Monitoring

During training, displaying loss values at regular intervals provides visibility into the optimization process. This helps diagnose issues like divergence, slow convergence, or overfitting.

### 10.1 What to Monitor

#### Training Loss Only

Displays the loss computed on the training data at each logging interval.

#### When to Use

- When no validation set is available
- When computational efficiency is a priority (validation adds overhead)

#### Limitation

Provides no indication of generalization — training loss may decrease while validation loss increases (overfitting).

## Training + Validation Loss

Displays both training loss and validation loss at each interval.

### When to Use

- Always, when a validation set is available (strongly recommended)
- Essential for early stopping
- Helps diagnose overfitting, underfitting, or stagnation

### Benefits

- Shows both learning progress (training) and generalization (validation)
- Clear visual signal when overfitting begins (training loss decreases, validation increases)

#### Best Practice

Always monitor both training and validation loss. The computational cost is minimal, and the insight into model behavior is invaluable for diagnosing problems early.

## 10.2 Logging Frequency

The frequency parameter controls how often loss values are displayed.

### Trade-offs

- **High frequency (e.g., every 10 iterations):** More granular updates; useful for debugging convergence issues; slight training slowdown
- **Low frequency (e.g., every 100 iterations):** Minimal overhead; suitable for long training runs; less visibility into iteration-by-iteration progress

### Recommendation

For production pipelines, log at a moderate frequency (e.g., once per epoch or every 50–100 iterations) and save the full training history to disk for post-hoc analysis.

# 11. Troubleshooting Guide for Neural Networks

Neural networks have many hyperparameters and settings, which provides flexibility but can also create challenges. This section provides a diagnostic framework for common training problems.

Symptom		Solutions	Root Cause
Loss decreasing	not	Check feature normalization; increase learning rate; add capacity (layers/neurons)	Model not learning

Symptom	Solutions	Root Cause
<b>Train Validation accuracy &gt;&gt;</b>	Add dropout/L2 regularization; reduce model complexity; use early stopping	Overfitting
<b>Both accuracies low</b>	Increase model capacity; reduce regularization; increase learning rate	Underfitting
<b>Loss oscillating</b>	Reduce learning rate; use gradient clipping; increase batch size	Training instability
<b>Very slow training</b>	Increase learning rate; reduce regularization; increase batch size	Slow convergence

## 11.1 Model Not Learning (Loss Not Decreasing)

If the loss fails to decrease over multiple epochs, the model is not learning. This is the most fundamental problem.

### Common Causes and Solutions

- **Features not normalized:** Neural networks are highly sensitive to feature magnitudes. Use standard scaling (zero mean, unit variance) or min-max normalization for all input features.
- **Learning rate too low:** Updates become negligible, leading to flat learning curves. Try increasing alpha from 0.001 to 0.01 or 0.1.
- **Learning rate too high:** Gradients overshoot the optimal trajectory entirely, causing loss to stay high or diverge. Try reducing alpha.
- **Model too simple:** The architecture lacks capacity to capture underlying relationships. Add more layers or more neurons per layer.
- **Insufficient data:** Neural networks require substantial data to learn patterns. If you have < 10K samples, consider simpler models (linear, tree-based).
- **Wrong output activation or loss function:** Verify that output activation matches problem type (Sigmoid for binary, Softmax for multi-class, Linear for regression) and that loss function is correctly paired.

## 11.2 Overfitting (Train Accuracy >> Validation Accuracy)

Overfitting occurs when training accuracy remains high but validation accuracy drops. The model is memorizing training patterns instead of generalizing.

### Common Causes and Solutions

- **Insufficient regularization:** Add dropout (0.2–0.5) or increase L2 weight decay (0.001–0.01).
- **Model too complex:** Reduce the number of layers or neurons per layer.

- **Not enough training data:** Collect more data or apply data augmentation techniques if applicable.
- **Missing early stopping:** Enable early stopping with patience = 20–50 to halt training when validation loss stops improving.
- **Biased validation set:** Ensure the validation split is representative of the problem space. Use stratified sampling for classification.

### 11.3 Underfitting (Both Train and Validation Accuracy Low)

Underfitting occurs when the model performs poorly on both training and validation sets. The architecture lacks capacity to model underlying patterns.

#### Common Causes and Solutions

- **Model too simple:** Add more layers or more neurons per layer to increase capacity.
- **Excessive regularization:** Reduce dropout rate or decrease L2 weight decay.
- **Learning rate too low:** Model cannot explore parameter space effectively. Increase learning rate.
- **Training stopped prematurely:** Increase maximum epoch count or adjust early stopping patience.

### 11.4 Training Instability (Loss Oscillating)

Training instability manifests as erratic, oscillating loss curves that prevent smooth convergence.

#### Common Causes and Solutions

- **Learning rate too high:** Large updates prevent smooth convergence. Reduce learning rate by 10× (e.g., 0.1 → 0.01).
- **Exploding gradients:** Gradient magnitudes grow uncontrollably in deep networks. Apply gradient clipping (clip norm to 1.0–5.0).
- **Batch size too small:** Very small batches (< 8) create noisy gradients. Increase batch size to 32, 64, or 128.
- **Poor weight initialization:** Ensure you're using He initialization for ReLU or Xavier for Sigmoid/Tanh.

### 11.5 Slow Convergence

Slow convergence means the model is learning, but progress is inefficiently slow.

#### Common Causes and Solutions

- **Learning rate too low:** Optimizer moves through loss landscape incrementally. Increase learning rate.

- **Heavy regularization:** High dropout ( $> 0.5$ ) or strong L2 ( $> 0.1$ ) slows progress. Reduce regularization strength.
- **Batch size too small:** Increase batch size to improve computational efficiency, especially on GPUs.
- **Suboptimal optimizer:** Switch from basic gradient descent to Adam for faster, more adaptive convergence.

## 12. Practical Recommendations

The following guidelines synthesize best practices from the sections above into a concise decision framework.

### Architecture

- Start with 2–3 hidden layers of 64–128 neurons each for tabular data
- Only increase depth ( $\geq 4$  layers) if you have  $> 100K$  samples and high-dimensional features
- Use wide and shallow over deep and narrow for structured data

### Activation Functions

- Hidden layers: ReLU (default)
- Output layer: Linear (regression), Sigmoid (binary classification), Softmax (multi-class)
- Only use Leaky ReLU if you observe widespread dying neurons in deep networks

### Weight Initialization

- ReLU networks: He initialization (default)
- Sigmoid/Tanh networks: Xavier initialization
- When in doubt: He initialization works for most cases

### Regularization

- L2 weight decay:  $\lambda = 0.001 - 0.01$  (default)
- Dropout: 0.2 – 0.5 for hidden layers (apply to deeper networks)
- Always use early stopping with patience = 20–50

### Optimization

- Optimizer: Adam ( $\alpha = 0.001$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ )
- Batch size: 32, 64, or 128 depending on dataset size and memory
- Learning rate decay: Usually unnecessary with Adam; use reduce-on-plateau if needed

### Loss Functions

- Regression: MSE (default) or MAE (if outliers present)
- Binary classification: Log Loss (Binary Cross-Entropy)
- Multi-class classification: Categorical Cross-Entropy

## Convergence & Monitoring

- Always monitor both training and validation loss
- Use early stopping — set high max epochs (500–1000) and let early stopping find the optimal point
- Log loss every epoch or every 50–100 iterations
- Save training history to disk for post-training analysis

## Preprocessing

- ALWAYS normalize/standardize input features (zero mean, unit variance)
- Check for missing values and handle appropriately
- For classification, use stratified sampling to preserve class ratios in train/val/test splits

### Best Practice

Neural networks are powerful but require careful configuration. For structured/tabular data, always establish baselines with simpler models (linear, tree-based) first. Neural networks should be chosen when:

- Data is large (> 100K samples)
- Features are high-dimensional or contain complex interactions
- Simpler models have plateaued in performance

When neural networks are justified, follow the guidelines in this document to maximize your chances of training a reliable, production-ready model.

## 13. Underfitting – Overfitting Trade-off

Achieving strong model performance is a critical stage in any machine learning pipeline, where the goal is not just to learn patterns from training data, but to ensure those patterns hold on unseen data. While high training accuracy may appear promising, real-world success depends on balancing underfitting and overfitting.

Neural Networks are universal function approximators whose capacity scales with depth and width. This power makes them highly susceptible to both failure modes: a network that is too small or trained too briefly will underfit complex patterns, while a large network trained too long on limited data will memorise training examples rather than generalise. Careful capacity planning and regularization are therefore central to training neural networks well.

Aspect	▽ UNDERFITTING		△ OVERFITTING	
	Causes	Actions to Fix	Causes	Actions to Fix
<b>Network Capacity (Architecture)</b>	Network is too shallow or too narrow; it lacks the representational power to capture complex feature interactions	Add more layers to increase depth and/or more neurons per layer to increase width	Network has far more parameters than the training set can constrain, allowing it to memorise examples	Reduce the number of layers or neurons; use architecture search to find the smallest sufficient model
<b>Regularization</b>	Dropout or weight decay is set too aggressively, suppressing the network's ability to learn genuine patterns	Reduce dropout rate and weight decay; validate that training loss itself is decreasing adequately	No regularization is applied; weights grow large and the network fits training noise perfectly	Apply dropout between layers; use weight decay; add batch normalisation to stabilise activations
<b>Training Duration</b>	Training is stopped too early before the network has converged to a good solution	Train for more epochs; monitor both training and validation loss to confirm convergence	Training continues long after the validation loss has stopped improving, memorising the training set	Use early stopping by monitoring validation loss and halting when it no longer improves
<b>Learning Rate</b>	Learning rate is too low, causing extremely slow convergence or getting trapped in poor local minima	Increase the learning rate; use adaptive optimisers that adjust the rate per parameter automatically	Learning rate is too high, causing unstable training and the model jumping past good solutions	Reduce the learning rate; use a learning rate schedule that decays the rate as training progresses
<b>Training Data Volume</b>	The dataset is too small relative to the network's capacity, providing insufficient signal	Collect more data; leverage transfer learning from a pre-trained model to compensate for limited data	Small dataset relative to model size — the network memorises every example rather than learning patterns	Apply data augmentation to artificially expand the training set; use a pre-trained backbone and fine-tune
<b>Batch Size</b>	Very large batches produce smooth gradient estimates that converge to sharp minima	Use smaller batches to introduce beneficial stochasticity and encourage convergence to flatter minima	Very small batches produce noisy gradient estimates that destabilise training and	Tune batch size to a moderate level; use gradient clipping to prevent instability from noisy batches

	with poor generalisation		slow convergence	
--	--------------------------	--	------------------	--

Have a question or suggestion?  
Reach us at [support@tvaritam.ai](mailto:support@tvaritam.ai)