

Teaching Machines to See Patterns

A Neural Networks Primer: Why We Needed Each Piece of the Puzzle

NLP Course 2025

From the 1950s mail sorting crisis to ChatGPT: How humanity taught machines to think

Where We're Going Today

Part 1: The Problem (1943-1969)

- The mail sorting crisis
- First mathematical neurons
- The perceptron revolution
- The XOR catastrophe

Intermission: Understanding the Basics

- How neurons calculate
- Why we need layers
- Following the forward pass

Part 2: The Breakthrough (1980s-1990s)

- Hidden layers save the day
- Backpropagation breakthrough
- Universal approximation proof

Part 3: The Revolution (2000s-Present)

- Deep learning explosion
- Modern architectures
- Real-world impact

Your Turn: Building Networks

- Build your first network
- Next steps

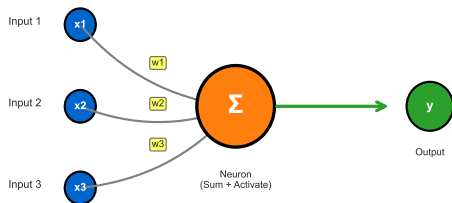
Each part builds on the previous - we'll go step by step!

What is a Neuron? The Building Block

Before we tell the story, let's understand the fundamental building block

Understanding a Single Neuron

Anatomy of a Single Neuron



What Happens Inside a Neuron

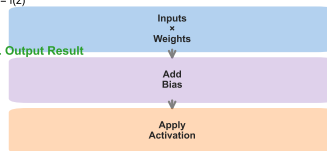
1. Weighted Sum

$$z = \sum(w_i \times x_i) + \text{bias}$$

2. Apply Activation

$$y = f(z)$$

3. Output Result

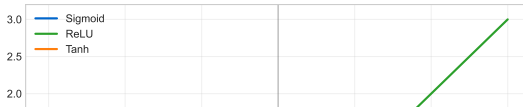


Concrete Example

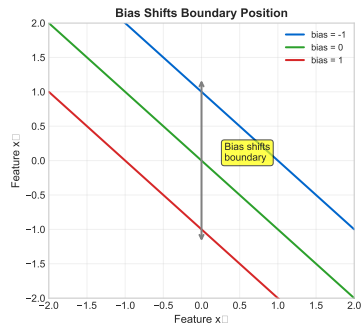
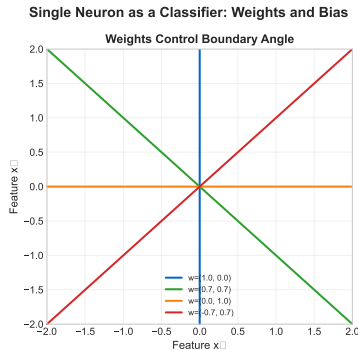
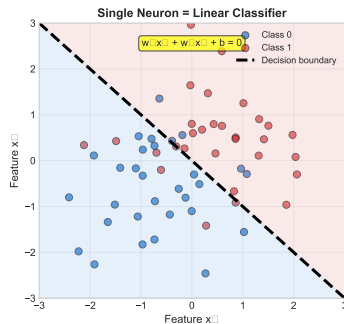
Example: Is this fruit ripe?

Inputs:
 $x_1 = 2$ (color: red=high)
 $x_2 = 3$ (softness: high)

Common Activation Functions



How Weights and Bias Control Decisions

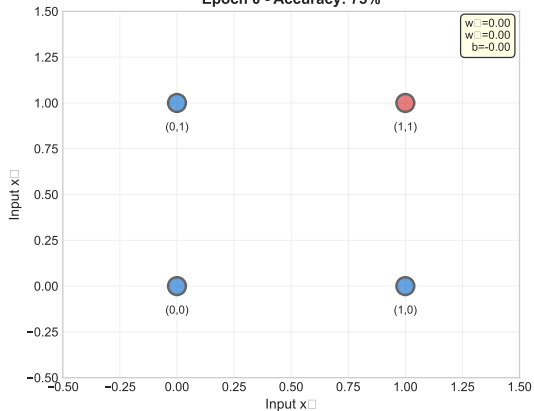


Weights control the angle, bias shifts the position - together they define the decision boundary

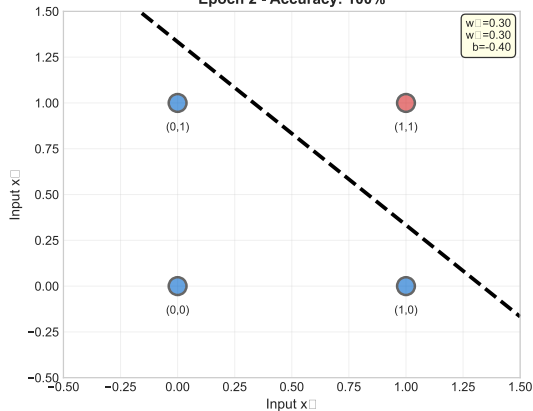
Learning = Gradually Adjusting Weights to Fit the Data

How a Neuron Learns: Training Progress on AND Logic

Epoch 0 - Accuracy: 75%



Epoch 2 - Accuracy: 100%



Epoch 5 - Accuracy: 100%



Epoch 10 - Accuracy: 100%



The Core Idea: Neural Networks are Function Approximators

What does this actually mean?

The Problem:

- We have inputs (x)
- We want outputs (y)
- But we don't know the formula!
- Examples:
 - Size \rightarrow Price
 - Image \rightarrow Label
 - Text \rightarrow Sentiment

Traditional Approach:

- Guess the formula
- Write explicit rules
- Hope it works
- **Problem:** Real world is too complex!

Example:

Price = $a \times \text{Size} + b$
(Too simple for real data!)

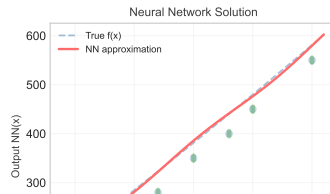
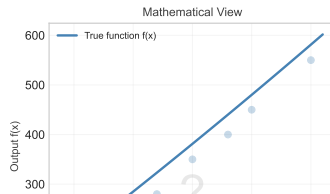
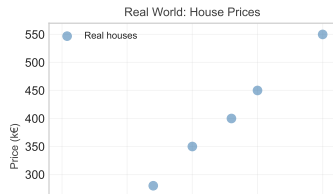
Neural Network Approach:

- Learn from examples
- Build the formula automatically
- Adjust until it fits
- **Works for ANY pattern!**

Magic:

NN learns: $f(x) \approx y$
No formula needed!

Function Approximation: Learning Patterns from Examples



How NNs Build Complex Functions from Simple Pieces

The LEGO Principle: Combine Simple Parts to Build Anything

The Building Blocks:

1. Individual Neurons:

- Each neuron = simple decision
- "Is input i threshold?"
- Outputs: on/off (or smooth version)

2. Combine Neurons:

- Add their outputs
- Weight their importance
- Create complex shapes

3. Stack Layers:

- First layer: simple features
- Next layer: combinations
- Final layer: complete function

Real-World Analogy:

Making a Cake (Complex) from Ingredients (Simple):

- Flour → Basic structure
- Sugar → Sweetness level
- Eggs → Binding agent
- Mix right amounts → Perfect cake!

In Neural Networks:

- Neuron 1 → Detects edges
- Neuron 2 → Detects curves
- Neuron 3 → Detects colors
- Combine all → Recognize faces!

Building Complex Functions from Simple Neurons



The Universal Approximation Theorem: Why This Always Works

The Most Important Theorem in Deep Learning (Cybenko, 1989)

The Theorem (Plain English):

*"A neural network with enough neurons can approximate **ANY** continuous function to **ANY** desired accuracy"*

What This Means:

- **Universal:** Works for any smooth pattern
- **Guaranteed:** Not hoping, but proving
- **Practical:** Just add more neurons!

The Catch:

- **How many neurons?** Could be millions
- **How to find weights?** That's training
- **How long to train?** That's the art

Intuitive Proof:

Think of it like pixel art:

1. With 4 pixels: Very blocky image
2. With 100 pixels: Recognizable
3. With 10,000 pixels: Photo-realistic
4. With infinite pixels: Perfect!

Same with neurons:

1. Few neurons: Rough approximation
2. More neurons: Better fit
3. Many neurons: Nearly perfect
4. Infinite neurons: Exact function!

Why This Matters:

We don't need different architectures for different problems - just one universal tool that adapts!

The Universal Approximation Theorem (Cybenko, 1989)

Universal Approximation: More Neurons = Better Fit

Network Width

1950s: The Mail Sorting Crisis

The Challenge:

- 150 million letters per day
- Hand-written addresses
- Human sorters: slow, expensive, error-prone
- Traditional programming: useless

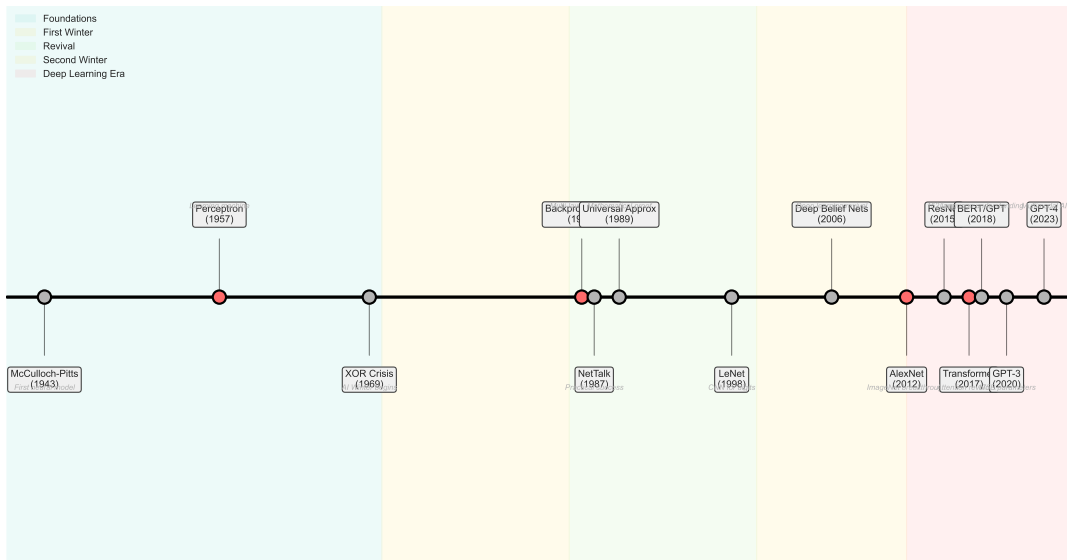
Why Traditional Code Failed:

- Can't write rules for every handwriting style
- Too many variations of each letter
- Context matters: "l" vs "I" vs "1"
- This wasn't computation—it was **pattern recognition**

This problem would take 40 years to solve properly

80 Years of Neural Networks: The Complete Journey

Neural Networks: 80 Years of Evolution



Why Can't We Just Write Rules?

Problem: Recognize the Letter "A"

Traditional Approach (Failed):

```
if (has_triangle_top AND  
    has_horizontal_bar AND  
    two_diagonal_lines) {  
  return "A"  
}
```

But what about...

- Handwritten A's?
- Different fonts?
- Rotated A's?
- Partial A's?

The Challenge: Infinite Variations of "A"

A A A A

A a A A

Just for the letter "A", we'd need thousands of rules!

The breakthrough: What if machines could learn patterns like children do?

The Birth of Computational Neuroscience

The Revolutionary Paper:

- "A Logical Calculus of Ideas Immanent in Nervous Activity"
- First mathematical model of neurons
- Proved: Networks can compute ANY logical function
- Inspired von Neumann's computer architecture

Key Insight:

- Neurons = Logic gates
- Brain = Computing machine
- Thinking = Computation

The Model:

- Binary neurons (0 or 1)
- Threshold activation
- Fixed connections
- No learning yet!

Historical Impact:

- Founded field of neural networks
- Influenced cybernetics movement
- Set stage for AI research
- "The brain is a computer" metaphor

14 years later, Rosenblatt would add the missing piece: learning

Frank Rosenblatt's Radical Idea: Neurons That Learn

Beyond McCulloch-Pitts:

- Adjustable weights (not fixed!)
- Learning from mistakes
- Physical machine built (Mark I)
- Could recognize simple patterns

The Hardware:

- 400 photocells (20×20 “retina”)
- 512 motor-driven potentiometers
- Weights adjusted by electric motors
- Took 5 minutes to learn patterns

Mathematical Model:

- Inputs: x_1, x_2, \dots, x_n
- Weights: w_1, w_2, \dots, w_n
- Sum: $z = \sum_{i=1}^n w_i x_i + b$
- Output: $y = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases}$

In plain words: Each input gets a vote (weight). We add up all votes plus a bias. If total is positive, output 1; otherwise 0.

Learning Rule: If wrong: $w_i = w_i + \eta \cdot \text{error} \cdot x_i$

The New York Times, 1958: "The Navy revealed the embryo of an electronic computer that will be able to walk, talk, see, write, reproduce itself and be conscious of its existence."

Let's Understand How This Actually Works

We've Seen the History...

- McCulloch-Pitts invented the neuron
- Rosenblatt made it learn
- The perceptron was born

Now Let's See the Science:

- How does a neuron calculate?
- What does learning mean?
- Why was XOR so hard?

Next 5 slides: Hands-on calculations and exercises
Get your pencil ready - we're going to work through real examples!

Don't worry - we'll return to the story once you understand the basics

Let's Make Sure We're Together

Quick Questions:

1. Why couldn't traditional programming solve mail sorting?
2. What does a weight represent in simple terms?
3. Why do we need the bias term?
4. What was revolutionary about Rosenblatt's perceptron?

Think About It:

- A weight is like the importance/trust we give to each input
- Bias shifts our decision threshold
- Learning = adjusting these weights
- The perceptron was the first machine that could learn!

Try It Yourself: Draw a simple perceptron with 2 inputs. Label the weights, bias, and output. What would the weights be to compute AND logic?

If any of these are unclear, revisit the previous slides before continuing

Problem: Learn OR function (output 1 if ANY input is 1)

Training Data:

x_1	x_2	Output
0	0	0
0	1	1
1	0	1
1	1	1

The Perceptron:

$$z = w_1 \cdot x_1 + w_2 \cdot x_2 + b$$

$$\text{output} = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases}$$

In plain words: Multiply first input by first weight, second input by second weight, add bias, check if positive

Learning Process:

1. Start with random weights
2. For each example:
 - Calculate output
 - If wrong: adjust weights
 - If correct: keep weights
3. Repeat until all correct

Final Solution: $w_1 = 1$, $w_2 = 1$, $b = -0.5$

Success! But this was just the beginning...

Let's Calculate Together: Is This Email Spam?

A Real Perceptron Calculation You Can Follow

The Email:

"FREE money! Click here NOW for amazing offer!!!"

Our Features (Inputs):

- $x_1 = \text{Has "FREE"}? = 1$
- $x_2 = \text{Has "money"}? = 1$
- $x_3 = \text{Many "!"?} = 1$
- $x_4 = \text{From friend?} = 0$

Learned Weights:

- $w_1 = +3$ (FREE is very spammy)
- $w_2 = +2$ (money is suspicious)
- $w_3 = +2$ (!!! is aggressive)
- $w_4 = -5$ (friends are trusted)
- $b = -2$ (threshold)

Let's Calculate:

$$\begin{aligned} z &= w_1 \cdot x_1 + w_2 \cdot x_2 + w_3 \cdot x_3 + w_4 \cdot x_4 + b \\ &= 3 \cdot 1 + 2 \cdot 1 + 2 \cdot 1 + (-5) \cdot 0 + (-2) \\ &= 3 + 2 + 2 + 0 - 2 \\ &= 5 \end{aligned}$$

Decision:

- $z = 5 > 0$
- Output = 1 = SPAM!

Try It Yourself: What if this email WAS from a friend ($x_4 = 1$)? Recalculate! Would it still be spam?

Answer: $z = 5 - 5 = 0$, borderline!

This is exactly how early spam filters worked - and why they failed on clever spam

Breaking Down the Math Symbols

Inputs and Weights:

- x_i = input value (what we see)
- w_i = weight (importance/strength)
- b = bias (threshold adjuster)

The Computation:

$$z = \sum_{i=1}^n w_i x_i + b$$

This means:

- Multiply each input by its weight
- Add them all up
- Add the bias

This simple math would evolve into deep learning

Real Example:

Should I go outside?

Factor	Value	Weight
Sunny?	1	+2
Raining?	0	-3
Weekend?	1	+1

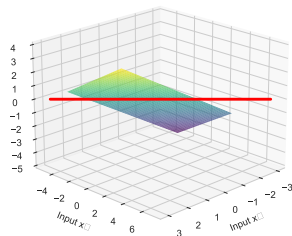
$$z = (1 \times 2) + (0 \times -3) + (1 \times 1) = 3$$

Decision: $z > 0$, so YES!

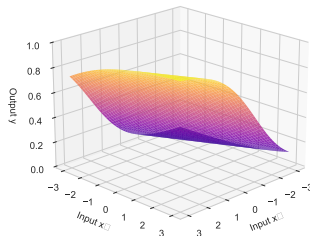
Visualizing How Activation Functions Transform the Output Space

Single Neuron Visualization: Effect of Activation Functions

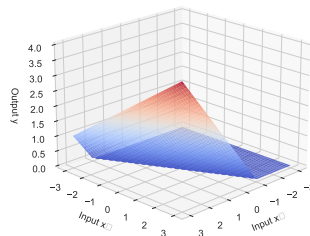
Linear (No Activation)
 $z = w_1x_1 + w_2x_2 + b$



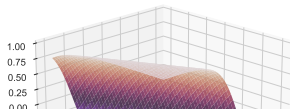
Sigmoid Activation
 $y = 1/(1 + e^{-z})$



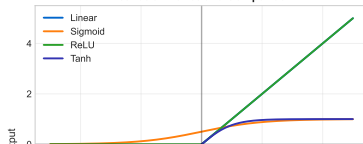
ReLU Activation
 $y = \max(0, z)$



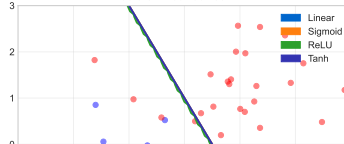
Tanh Activation
 $y = \tanh(z)$



Activation Functions Comparison



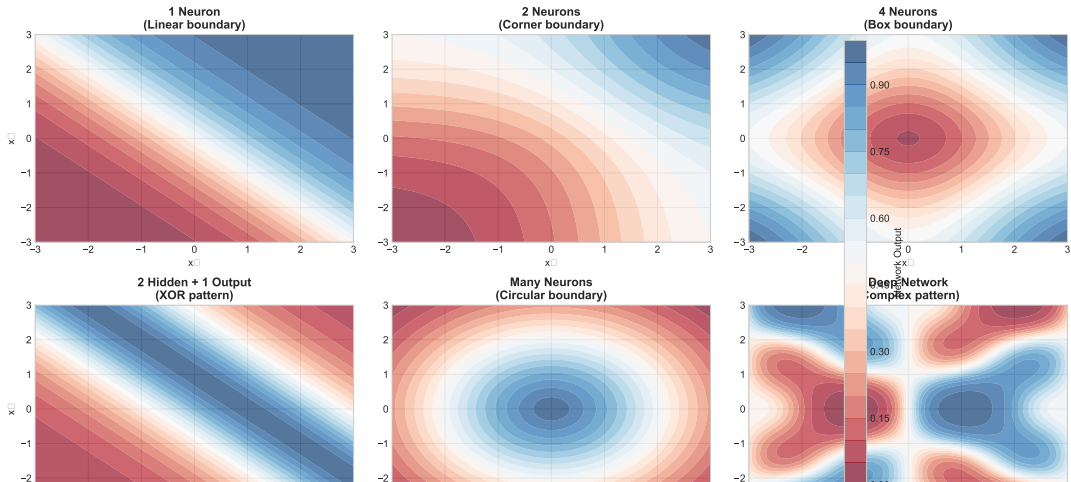
Decision Boundaries in 2D



From Simple to Complex: Network Depth Creates Complexity

How More Neurons Enable More Complex Decision Boundaries

How Network Complexity Grows with Neurons and Layers



1969: The Crisis - XOR Problem

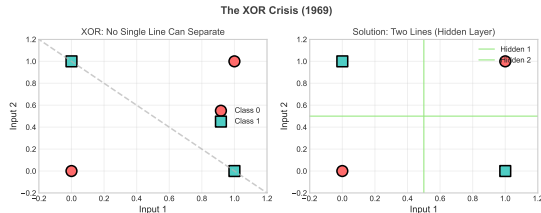
XOR (Exclusive OR):

x_1	x_2	Output
0	0	0
0	1	1
1	0	1
1	1	0

The Problem:

- Can't draw a single line to separate
- Perceptron only learns linear boundaries
- Real-world problems are non-linear!

The field would be dormant for over a decade...

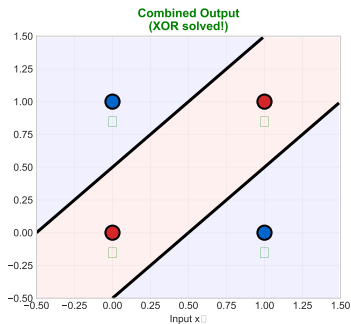
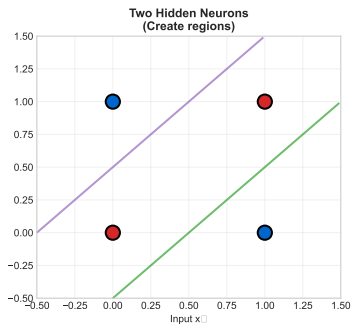
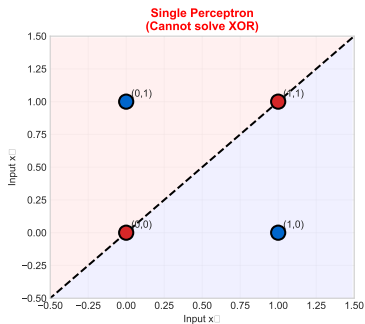


Impact:

- Funding dried up
- "AI Winter" begins
- Neural networks abandoned

Why We Need Hidden Layers: The XOR Solution

Solving XOR: Why We Need Hidden Layers



Two hidden neurons working together can solve what one neuron cannot

When One Line Isn't Enough: Real Problems Need More

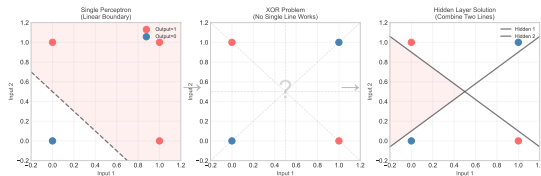
Let's See Why We Need Hidden Layers

Problem 1: Spam Detection (Easy)

- Has many spam words? → SPAM
- Has few spam words? → NOT SPAM
- One line (threshold) works!

Problem 2: Cat or Dog Photo (Hard)

- Small + fluffy? Could be either!
- Large + smooth? Could be either!
- Pointy ears + whiskers? → Cat
- Floppy ears + wet nose? → Dog
- Need multiple feature detectors!



The Solution:

1. First layer: Multiple detectors
 - Detector 1: "Has cat features?"
 - Detector 2: "Has dog features?"
2. Second layer: Combine detections
 - If cat features \wedge dog features → Cat

This is why deep learning works: each layer builds more complex detectors from simpler ones

1980s: The Hidden Layer Revolution

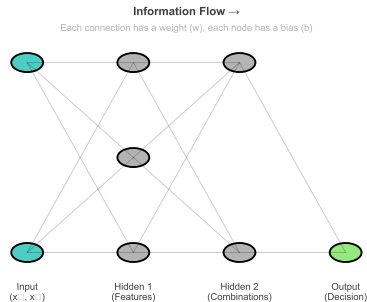
The Insight:

- Stack multiple layers!
- First layer: detect simple features
- Hidden layer: combine features
- Output layer: final decision

Solving XOR:

- Hidden neuron 1: Is it (0,1)?
- Hidden neuron 2: Is it (1,0)?
- Output: OR of hidden neurons

Multi-Layer Network: Solving Complex Problems



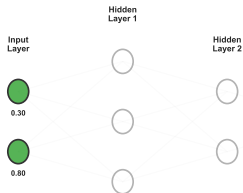
New Architecture:

- Input layer: raw data
- Hidden layer(s): feature extraction

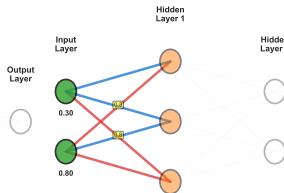
Forward Pass: Signal Propagation Step-by-Step

Following Data as it Flows Through the Network

Frame 1: Input Data

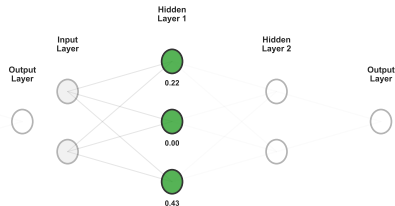


Forward Pass: Step-by-Step Signal Propagation
Frame 2: First Layer Computation

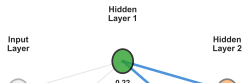


$$\text{Computing: } h_1 = f(W_1 \times h_0)$$

Frame 3: Hidden Layer 1 Activated



Frame 4: Second Layer Computation



Frame 5: Hidden Layer 2 Activated



Frame 6: Final Output



Rumelhart, Hinton, Williams: The Learning Algorithm

The Problem:

- Perceptron learning only works for 1 layer
- How to adjust hidden layer weights?
- No direct error signal for hidden neurons

The Solution:

- Propagate error backwards!
- Each layer gets blame for output error
- Use calculus (chain rule) to distribute blame

The Algorithm:

1. Forward: Calculate output
2. Compare: Find error
3. Backward: Distribute blame
4. Update: Adjust all weights

In plain words: Like a teacher marking an essay: finds the final error, then traces back to see which paragraphs, sentences, and words caused it

Impact:

- Finally could train deep networks!
- Neural networks reborn
- Foundation of all modern AI

This algorithm runs billions of times to train ChatGPT

Cybenko, Hornik: The Ultimate Proof

The Theorem:

A feedforward network with:

- One hidden layer
- Enough neurons
- Non-linear activation

Can approximate ANY continuous function to arbitrary accuracy!

What This Means:

- Neural networks are universal computers
- No function is too complex
- Just need enough neurons and data

Real-World Analogy:

LEGO blocks can build anything:

- Few blocks = rough shape
- Many blocks = detailed model
- Infinite blocks = perfect replica

Same with neurons:

- Few neurons = rough approximation
- Many neurons = good function
- Infinite neurons = exact function

Try It Yourself: Think of any pattern or function. This theorem guarantees a neural network can learn it!

This gave theoretical backing to the neural network revolution

The Need for Non-Linearity

Problem with Linear:

- Stack of linear layers = still linear!
- $f(g(x)) = (wx + b_1)w' + b_2 = w'wx + \dots$
- Can't learn complex patterns

Solution: Activation Functions

- Add non-linearity after each layer
- Allows learning complex boundaries
- Different functions for different needs

Common Activation Functions:

- **Sigmoid:** $\sigma(x) = \frac{1}{1+e^{-x}}$
 - Smooth, outputs 0-1
 - Good for probabilities

In plain words: Squashes any input to range 0-1. Large positive becomes 1, large negative becomes 0

- **ReLU:** $f(x) = \max(0, x)$
 - Simple, fast
 - Solves vanishing gradient
- **Tanh:** $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
 - Outputs -1 to 1
 - Zero-centered

ReLU's simplicity revolutionized deep learning in 2011

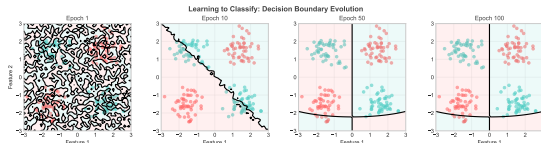
Teaching a Network to Separate Red from Blue Points

The Setup:

- Input: (x, y) coordinates
- Output: Red or Blue class
- Network: $2 \rightarrow 4 \rightarrow 2$ neurons

Training Process:

1. Epoch 1: Random boundary
2. Epoch 10: Rough separation
3. Epoch 50: Good boundary
4. Epoch 100: Perfect fit



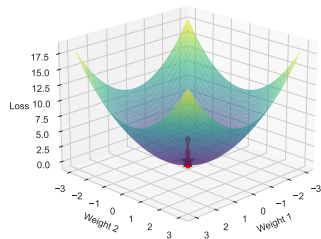
What Each Layer Learns:

- Layer 1: Simple boundaries
- Hidden: Combine boundaries
- Output: Final decision

This same principle scales to millions of parameters

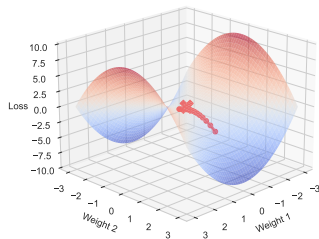
Gradient Descent: Finding the Valley in 3D Space

Ideal Case: Convex Loss
(Easy optimization)

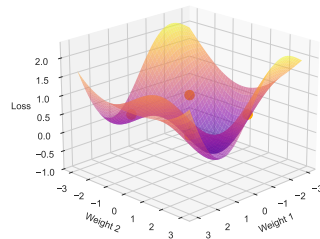


Gradient Descent: Navigating the Loss Landscape

Saddle Point Problem
(Gradient = 0, not minimum)



Multiple Local Minima
(Can get stuck)



Convex Loss (Top View)



Effect of Learning Rate



Loss Decrease Over Iterations



1998-2012: From Digits to ImageNet

1998 - LeNet: First Success

- Yann LeCun's CNN for digits
- 32×32 pixels \rightarrow 10 classes
- 60,000 parameters
- Banks adopt for check reading

Key Innovation: Convolutions

- Share weights across image
- Detect features anywhere
- Build complexity layer by layer

2012 - AlexNet: The Revolution

- 1000 ImageNet classes
- 60 million parameters
- GPUs enable training
- Error rate: 26% \rightarrow 16%

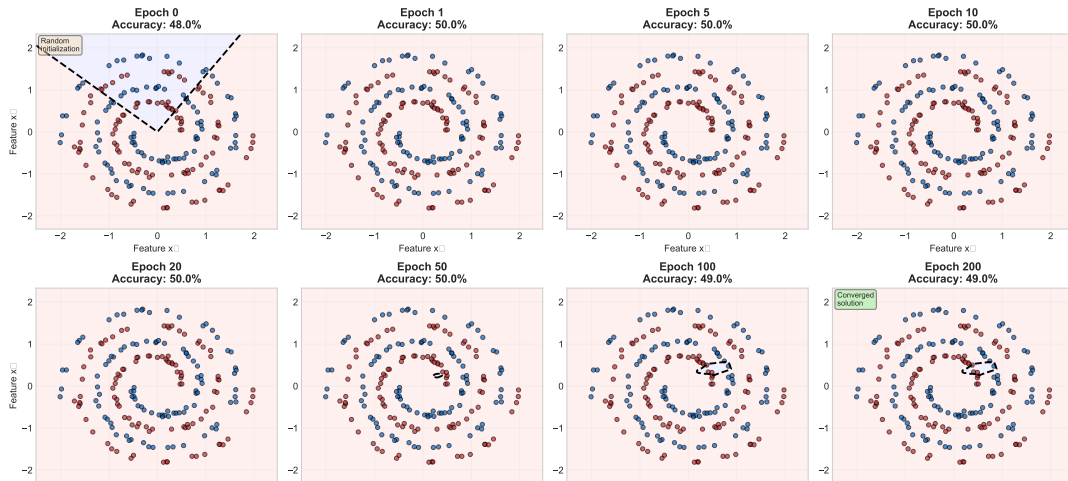
What Changed:

- Big Data (millions of images)
- GPU computing (100x faster)
- ReLU activation
- Dropout regularization

This victory ended the second AI winter permanently

Watching Decision Boundaries Evolve During Training

Neural Network Learning: Decision Boundary Evolution



2014-Present: Networks That Changed the World

The Depth Revolution:

- 2014 - VGGNet: 19 layers
- 2015 - ResNet: 152 layers
- 2017 - Transformers: Attention
- 2020 - GPT-3: 175B parameters

Why Depth Matters:

- Each layer = abstraction level
- Deep = complex reasoning
- Hierarchical feature learning

Real-World Impact:

- **Vision:** Self-driving cars
- **Language:** Google Translate
- **Speech:** Siri, Alexa
- **Medicine:** Disease diagnosis
- **Science:** Protein folding

The Scale:

- Billions of parameters
- Trained on internet-scale data
- Months of GPU time
- Emergent abilities appear

We went from recognizing digits to passing the bar exam in 25 years

Problem: Networks Couldn't Get Deeper

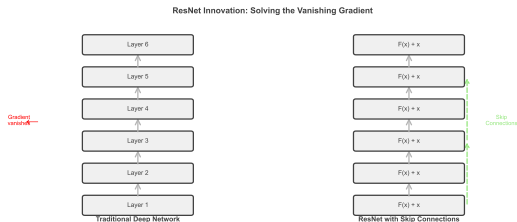
The Vanishing Gradient:

- Gradients multiply through layers
- Become exponentially small
- Deep layers stop learning
- 20 layers was the limit

The Breakthrough: Skip Connections

- Add input directly to output
- $F(x) + x$ instead of just $F(x)$
- Gradients flow directly backward
- Can train 1000+ layers!

This simple trick enabled the deep learning revolution



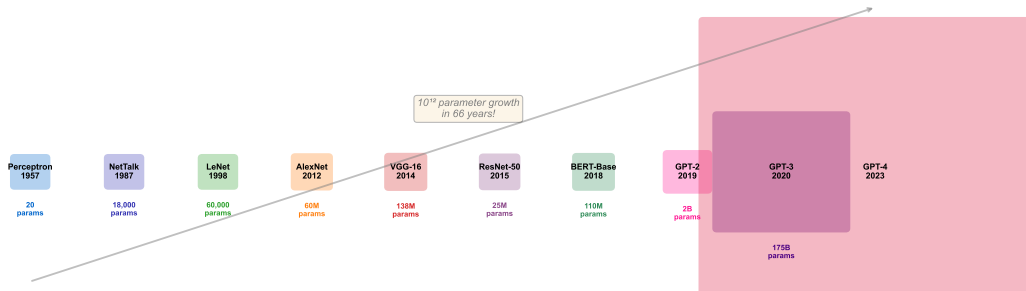
Why It Works:

- Learn residual (difference) only
- Identity mapping is easy default
- Gradients have direct path
- Each layer refines previous result

From 20 Parameters to 1.8 Trillion: The Growth of Neural Networks

Neural Network Evolution: From Perceptron to GPT-4

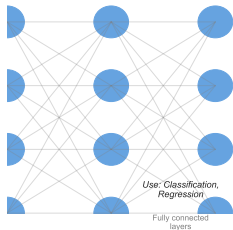
Box size represents relative number of parameters



Different Architectures for Different Problems

Neural Network Architecture Types

**Feedforward
(MLP)**



Transformer

**Convolutional
(CNN)**



Use: Images, Video

Spatial feature detection

**Graph Neural
(GNN)**

**Recurrent
(RNN/LSTM)**



Use: Text, Time-series

Sequential processing

**Generative
(GAN/VAE)**

Feedforward Networks:

- Information flows forward only
- Fixed-size input and output
- Good for: Classification, regression

Convolutional (CNN):

- Spatial feature detection
- Translation invariance
- Good for: Images, video

Recurrent (RNN):

- Process sequences
- Maintain memory/state
- Good for: Text, time-series

Transformer:

- Attention mechanism
- Parallel processing
- Good for: Language, everything else

Each architecture encodes different assumptions about the data

Transfer Learning:

- Start with pre-trained network
- Fine-tune on your task
- 100x less data needed
- Days → Hours training

Data Augmentation:

- Create variations of training data
- Rotations, crops, color shifts
- Prevents overfitting
- Free performance boost

Simple Optimizers to Start:

- **SGD:** Basic gradient descent
- **Adam:** Adaptive learning rates (use this!)

Mixed Precision:

- Use 16-bit floats where possible
- Keep 32-bit for critical ops
- 2-3x speedup
- Same accuracy

These techniques make deep learning practical for everyone

Misconceptions That Will Confuse You

WRONG: "Neurons are like brain neurons"

- **Brain neurons:** Complex, chemical, adaptive
- **Artificial neurons:** Simple math functions
- Just multiply and add!
- No biology involved

WRONG: "Networks understand concepts"

- **What you think:** "It knows what a cat is"
- **Reality:** It found statistical patterns
- No understanding, just correlation
- Can be fooled by tiny changes

WRONG: "More layers = always better"

- **Too deep:** Vanishing gradients
- **Too deep:** Overfitting
- **Right depth:** Depends on problem complexity
- Simple problems need shallow networks

WRONG: "It learns like humans"

- **Humans:** Learn from few examples
- **Humans:** Transfer knowledge easily
- **Networks:** Need thousands of examples
- **Networks:** Struggle with new situations

Remember: Neural networks are just fancy pattern matchers.
They don't think, understand, or reason - they find correlations in data.

Understanding these limits helps you use neural networks effectively

1. Data Explosion:

- Internet = infinite training data
- ImageNet: 14M labeled images
- Common Crawl: 300TB of text
- YouTube: 500 hours/minute

2. Hardware Revolution:

- GPUs: 100x faster than CPUs
- TPUs: Built for neural nets
- Cloud computing: Rent supercomputers
- Mobile chips with NPUs

3. Algorithm Breakthroughs:

- ReLU activation (2011)
- Batch normalization (2015)
- Skip connections (2015)
- Attention mechanism (2017)

4. Open Source Culture:

- TensorFlow, PyTorch free
- Pre-trained models shared
- Papers with code
- Collaborative research

The same ideas from 1980s finally had the resources to work

Building a Digit Classifier in 10 Lines

PyTorch Implementation:

```
import torch
import torch.nn as nn

class SimpleNet(nn.Module):
    def __init__(self):
        super().__init__()
        self.fc1 = nn.Linear(784, 128)
        self.fc2 = nn.Linear(128, 10)

    def forward(self, x):
        x = torch.relu(self.fc1(x))
        return self.fc2(x)

# Train
model = SimpleNet()
optimizer = torch.optim.Adam(model.parameters())
criterion = nn.CrossEntropyLoss()
```

This simple network achieves 97% accuracy on MNIST

What This Does:

- Input: 28×28 pixel image
- Hidden: 128 neurons
- Output: 10 digit classes
- Activation: ReLU
- Training: Adam optimizer

Training Loop:

- Forward pass
- Calculate loss
- Backward pass
- Update weights
- Repeat

Systematic Debugging Saves Hours

Try It Yourself: Save this checklist - you'll need it for every project!

Step 1: Sanity Checks

- ☐ Can you overfit a single batch?
- ☐ Are inputs normalized?
- ☐ Is output layer correct?
- ☐ Loss function matches task?

Step 2: Data Checks

- ☐ Plot sample inputs
- ☐ Check label distribution
- ☐ Verify train/val split
- ☐ Look for data leakage

Step 3: Training Checks

- ☐ Plot loss curves
- ☐ Check gradient norms
- ☐ Monitor weight updates
- ☐ Try different learning rates

Step 4: Architecture

- ☐ Start with known working model
- ☐ Add complexity gradually
- ☐ Check activation distributions
- ☐ Verify dimensions match

Common Confusion: 90% of bugs are in data preprocessing, not the model!

Print this slide and keep it handy

The Journey So Far

Core Concepts:

1. **Neurons:** $y = f(\sum w_i x_i + b)$
2. **Learning:** Adjust weights to minimize error
3. **Depth:** Each layer adds abstraction
4. **Backpropagation:** Distribute error backwards
5. **Non-linearity:** Enables complex functions

Historical Lessons:

1. Every limitation spawned innovation
2. Simple ideas + scale = revolution
3. Biology inspires but doesn't limit
4. Persistence through AI winters
5. Theory + engineering = breakthroughs

You now understand the fundamentals that power all modern AI

Let's Build Something Real!

Complete MNIST Classifier:

```
import torch
import torch.nn as nn
from torchvision import datasets, transforms
from torch.utils.data import DataLoader
```

1. Define Network

```
class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        self.fc1 = nn.Linear(784, 128)
        self.fc2 = nn.Linear(128, 64)
        self.fc3 = nn.Linear(64, 10)

    def forward(self, x):
        x = x.view(-1, 784) # Flatten
        x = torch.relu(self.fc1(x))
        x = torch.relu(self.fc2(x))
        return self.fc3(x)
```

2. Load Data

```
transform = transforms.ToTensor()
train_data = datasets.MNIST('.', train=True,
                             download=True,
                             transform=transform)
train_loader = DataLoader(train_data,
                           batch_size=64,
                           shuffle=True)
```

3. Setup Training

4. Training Loop

```
for epoch in range(3):
    for batch_idx, (data, target) in
        enumerate(train_loader):
        # Forward pass
        output = model(data)
        loss = criterion(output, target)

        # Backward pass
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()

    # Print progress
    if batch_idx % 100 == 0:
        print(f'Epoch-{epoch}:-{loss:.4f}')
```

5. Test One Example

```
model.eval()
test_image = train_data[0][0]
prediction = model(test_image.unsqueeze(0))
print(f" Predicted :-{prediction.argmax()}")
```

What You Just Built:

- 3-layer neural network
- 60K training images
- 97% accuracy in 3 epochs

Continue Your Neural Network Journey

Next Topics to Learn:

1. **CNNs:** Computer vision
2. **RNNs/LSTMs:** Sequences
3. **Transformers:** Modern NLP
4. **GANs:** Generation
5. **RL:** Decision making

Practical Projects:

- Image classifier for your photos
- Sentiment analysis of tweets
- Chatbot for customer service
- Style transfer for art
- Game-playing agent

Resources:

- **Fast.ai:** Practical deep learning
- **PyTorch Tutorials:** Official guides
- **Papers with Code:** Latest research
- **Kaggle:** Competitions and datasets
- **3Blue1Brown:** Visual explanations

Remember:

- Start simple, build up
- Theory + practice together
- Join communities
- Build projects you care about
- Share what you learn

**You've learned how humanity taught machines to think.
Now it's your turn to push the boundaries!**

Deep Dives for the Curious

This section contains advanced material that goes beyond the core BSc curriculum.

Topics covered:

- The Lottery Ticket Hypothesis
- Inductive Biases in Neural Architectures
- Scaling Laws and Performance Prediction
- Deep vs Wide Network Architectures
- Emergent Abilities at Scale
- Advanced Optimization Algorithms

These topics are valuable for understanding state-of-the-art research but not essential for getting started with neural networks.

Most Network Weights Don't Matter!

The Discovery:

- Networks contain "winning tickets"
- Subnetworks that train well alone
- 90-95% of weights can be removed
- Performance stays the same!

The Hypothesis: "Dense networks succeed because they contain sparse subnetworks that are capable of training effectively"

Implications:

- We massively overparameterize
- Training finds the needle in haystack
- Future: Train small from start?
- Mobile deployment possible

Why It Matters:

- Explains why big networks train better
- Pruning after training works
- Efficiency revolution starting
- Changes how we think about learning

A 1 billion parameter model might only need 50 million

The Right Architecture for the Right Problem

What Are Inductive Biases?

- Assumptions built into architecture
- Guide learning toward solutions
- Trade flexibility for efficiency
- "Priors" about the problem

Examples:

- **CNN:** Spatial locality matters
- **RNN:** Order/time matters
- **GNN:** Graph structure matters
- **Transformer:** All positions can interact

Why They Matter:

- Reduce search space
- Faster convergence
- Better generalization
- Less data needed

The Tradeoff:

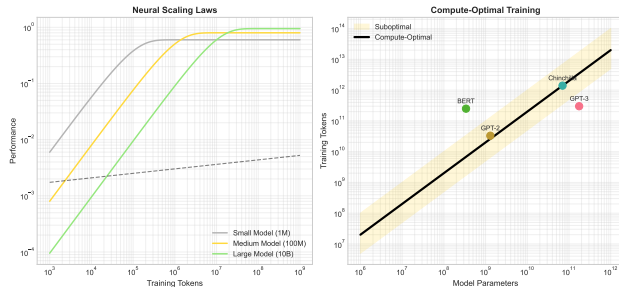
- Right bias = 10x better
- Wrong bias = 10x worse
- General architectures = safe but slow
- Specialized = fast but limited

Choosing the right inductive bias is still an art

Scaling Laws: How Performance Grows with Data

The Predictable Relationship Between Data, Model Size, and Performance

Data Scaling Laws (Chinchilla, 2022)



The Chinchilla Law (2022):

- Optimal ratio: 20 tokens per parameter
- 10B model needs 200B tokens
- Most models are undertrained
- Data quality matters more than quantity

Power Law Scaling:

Practical Implications:

- 10x data \rightarrow 2x performance
- 10x parameters \rightarrow 1.7x performance
- 10x compute \rightarrow 3x performance
- Diminishing returns always

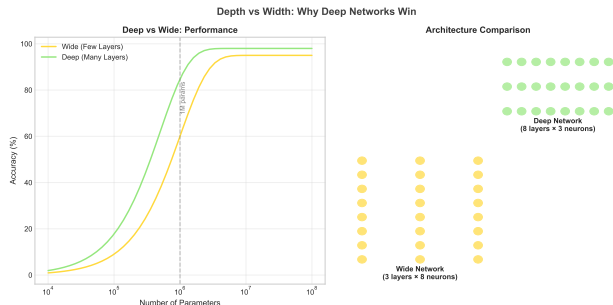
Data Efficiency Tricks:

- Data augmentation
- Synthetic data generation
- Active learning
- Curriculum learning
- Multi-task training

Why it matters: These laws predict costs before training

Current Limits

The Fundamental Tradeoff in Neural Architecture



Deep Networks (Many Layers):

- Complex hierarchical features
- Exponential expressiveness growth
- Harder to train (vanishing gradients)
- Better for vision, NLP

Wide Networks (Many Neurons):

The Sweet Spot:

- Vision: Deep (100+ layers)
- Language: Very deep (24-96 layers)
- Tabular: Wide and shallow (2-4 layers)
- Time series: Moderate (5-10 layers)

Modern Insights:

- Depth beats width for same parameters
- Skip connections enable extreme depth
- Width helps with memorization
- Depth helps with generalization

Scaling Laws:

- Performance \propto depth^{0.8}

Capabilities That Appear Suddenly with Scale

The Phenomenon:

- Small models: Can't do task at all
- Medium models: Still can't
- Large models: Suddenly can!
- No gradual improvement

Examples:

- 3-digit arithmetic (≈ 10 B params)
- Chain-of-thought reasoning (≈ 50 B)
- Code generation (≈ 20 B)
- Multilingual translation (≈ 100 B)

Why It Happens:

- Complex patterns need capacity
- Phase transitions in learning
- Composition of simpler abilities
- "Grokking" - sudden understanding

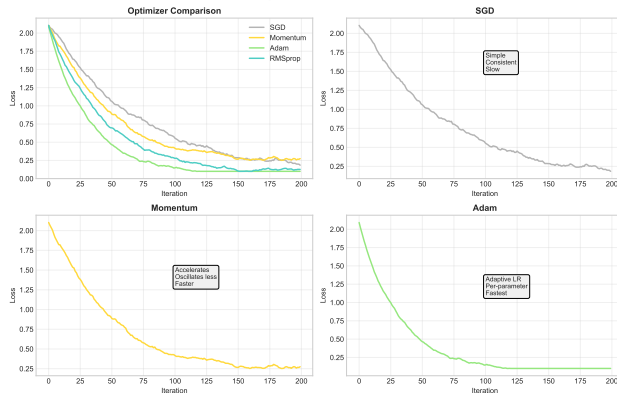
Implications:

- We can't predict what's next
- Scaling might unlock AGI
- Or hit fundamental limits
- Active area of research

GPT-3 showed abilities nobody expected or programmed

The Evolution of Gradient Descent

Modern Optimizers: From SGD to Adam



SGD (1951):

- Basic gradient descent
- Learning rate: Fixed

Adam (2014):

- Adaptive learning rates per parameter
- Combines momentum + RMSprop
- De facto standard
- Works out-of-the-box

Modern Variants:

- **AdamW:** Decoupled weight decay
- **RAadam:** Rectified Adam
- **LAMB:** Large batch training
- **Sophia:** 2nd-order approximation

Choosing an Optimizer:

- Start with Adam ($\text{lr}=3\text{e-}4$)

The Full Story: Historical Deep Dives

This section contains fascinating historical details that enrich the narrative but aren't essential for understanding the technical concepts.

Topics covered:

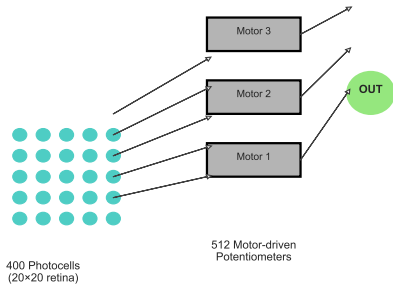
- The Mark I Perceptron: Physical Hardware
- NetTalk: Networks Learn to Speak (1987)
- Batch Normalization: Keeping Networks Stable

These stories show how each breakthrough built on previous work and overcame specific limitations.

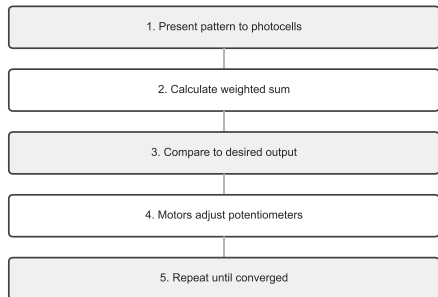
The Mark I Perceptron: A Physical Learning Machine

The Mark I Perceptron (1957): A Physical Learning Machine

Mark I Perceptron Architecture



Physical Learning Process



The first neural network wasn't software—it was a room-sized machine with motors and photocells

Sejnowski & Rosenberg: The First Viral NN Demo

The Challenge:

- Convert written text to speech
- English is irregular (tough, though, through)
- Rule-based systems had 1000s of exceptions

The Network:

- 7×29 input (7-letter window)
- 80 hidden neurons
- 26 output phonemes
- Trained overnight on DEC workstation

The Magic:

- Started: Random babbling
- Hour 1: Vowel-consonant patterns
- Hour 5: Recognizable words
- Hour 10: 95% accuracy!

Hidden Neurons Learned:

- Vowel detectors
- Consonant clusters
- Word boundaries
- Nobody programmed these!

Common Confusion: The network discovered linguistic concepts on its own - features linguists took centuries to identify!

Media sensation: "Computer teaches itself to read aloud overnight"

The Internal Covariate Shift Problem

BatchNorm Algorithm:

$$\mu_B = \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$y_i = \gamma \hat{x}_i + \beta$$

The Issue:

- Each layer's input distribution changes
- As previous layers update
- Makes learning unstable
- Requires tiny learning rates

The Solution:

- Normalize inputs to each layer
- Mean = 0, Variance = 1
- Learn scale and shift parameters
- Apply during training and testing

In plain words: 1) Find average, 2) Find spread, 3) Normalize to standard range, 4) Scale and shift as needed

Benefits:

- 10x faster training
- Higher learning rates OK
- Less sensitive to initialization