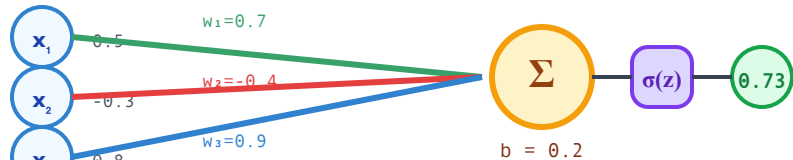


# Neural Network Basics Cheat Sheet

Traditional Feedforward Architectures - Prof. Laio Oriel Seman



## Single Neuron (Building Block)



### Step-by-Step Calculation:

#### 1. Linear Combination:

$$z = 0.7 \times 0.5 + (-0.4) \times (-0.3) + 0.9 \times 0.8 + 0.2$$
$$z = 0.35 + 0.12 + 0.72 + 0.2 = 1.39$$

#### 2. Activation Function:

$$\text{output} = \sigma(1.39) = 1/(1+e^{-1.39}) \approx 0.73$$

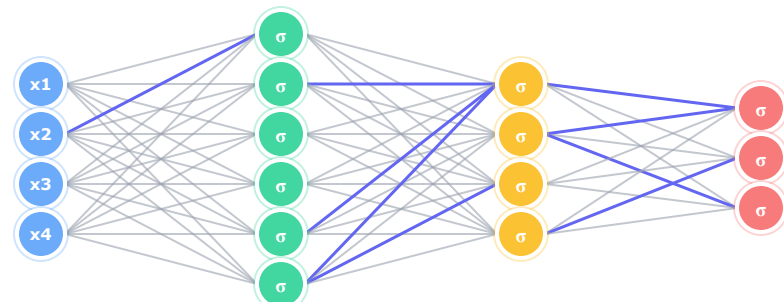
💡 Each neuron is a weighted sum followed by a non-linear transformation



## Feedforward Architecture

**INPUT** (4) **HIDDEN 1** (6) **HIDDEN 2** (4) **OUTPUT** (3)

24 connections 24 connections 12 connections



**Fully Connected: Every neuron connects to every neuron in the next layer**

Raw features Pattern detection Feature combination Final prediction



## Loss Functions

### Regression

#### Mean Squared Error

$$MSE = (1/n) \sum (y - \hat{y})^2$$

Penalizes large errors more heavily. Use when: Regression problems, outliers are important

#### Mean Absolute Error

$$MAE = (1/n) \sum |y - \hat{y}|$$

Robust to outliers. Use when: Regression with outliers

### Classification

#### Cross Entropy

$$CE = -\sum y_i \log(\hat{y}_i)$$

Measures difference between probability distributions. Use when: Multi-class classification

#### Binary Cross Entropy

$$BCE = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$$

Specialized for binary classification. Use when: Binary classification problems



## Optimization

### Batch Processing

#### SGD

Batch: 1 (sample)  
✓ Simple and fast  
✗ Noisy updates

#### Mini-Batch SGD

Batch: 32-128  
✓ Good balance  
✗ Requires tuning

#### Full Batch

Batch: All data  
✓ Stable gradients  
✗ Very slow on large datasets

### Popular Optimizers

#### Adam

$m^- = \beta_1 m + (1-\beta_1)g$ ,  $v^- = \beta_2 v + (1-\beta_2)g^2$   
✓ Adaptive learning rate  
✓ Combines momentum & RMS  
✓ Works well out of the box

Default lr: 0.001

Default choice

#### RMSprop

$v = \gamma v + (1-\gamma)g^2$ ,  $w \leftarrow w - \frac{g}{\sqrt{v}}$   
✓ Good for RNNs  
✓ Handles non-stationary loss

Default lr: 0.001

RNNs, streaming tasks

#### AdaGrad

$w \leftarrow w * g / (\sqrt{G} + \epsilon)$   
✓ Good for sparse data  
✓ Automatically adjusts lr

Default lr: 0.01

NLP, sparse features



## Regularization

#### L1 Regularization

$$\text{Loss} + \lambda ||w||_1$$

Encourages sparsity (some weights become 0)  
💡 Feature selection  
sparse models

#### L2 Regularization

$$\text{Loss} + \lambda ||w||_2^2$$

Reduces weight magnitudes evenly  
💡 Prevent overfitting  
smooth weights

#### Early Stopping

Stop when validation loss increases  
Prevents over-training  
💡 Automatic regularization



## Weight Initialization

#### Xavier/Glorot Initialization

$\text{Var}(W) = 1/n_{\text{in}}$   
Best for: Tanh, Sigmoid activations  
Maintains variance across layers

#### He Initialization

$\text{Var}(W) = 2/n_{\text{in}}$   
Best for: ReLU family activations  
Accounts for ReLU killing half the neurons



## Learning Rate Schedulers

#### Step Decay

$\text{lr} = \text{lr}_0 \times \gamma^{\lfloor \text{epoch}/s \rfloor}$   
Drops lr at fixed steps  
💡 Plateaued loss  
Scheduled decay

#### Cosine Annealing

$\text{lr} = \eta_{\text{min}} + 0.5(\eta_{\text{max}} - \eta_{\text{min}})(1 - \cos(\pi \cdot \text{epoch}/\text{steps}))$   
Smooth cyclical decay  
💡 Training from scratch  
Vision tasks

#### Reduce LR on Plateau

if val\_loss > lr \* factor  
Adaptive on stagnation  
💡 Validation-driven decay  
Auto-tuning



## Layer Types

### Dense Layer (Fully Connected)

$$z = Wx + b$$

Every neuron connects to all neurons in next layer • Feature transformation • Classification • Regression

### Dropout (Regularization)

$$y = x * \text{mask (training)}, y = x * p \text{ (inference)}$$

Randomly zeros neurons during training • Prevent overfitting • Regularization • Improve generalization

### Batch Normalization (Normalization)

$$y = \gamma((x-\mu)/\sigma) + \beta$$

Normalizes inputs to stabilize training • Faster training • Higher learning rates • Reduce internal covariate shift

### Layer Normalization (Normalization)

$$y = \gamma((x-\mu)/\sigma) + \beta$$

Normalizes across features instead of batch • RNNs, Transformers • Stabilize training in sequential models



## Training Pipeline

#### Forward Pass

$$x \rightarrow \hat{y}$$

Compute predictions through network

#### Loss Calculation

$$L(y, \hat{y})$$

Compare predictions with true labels

#### Backward Pass

$$\nabla L / \nabla w$$

Compute gradients via backpropagation

#### Parameter Update

$$w \leftarrow w - \eta \nabla w$$

Update weights using optimizer

Repeat for epochs

### Mathematical Foundation:

$$z^{(l+1)} = W^{(l+1)}a^{(l)} + b^{(l+1)}$$

$$a^{(l)} = \sigma(z^{(l)})$$

$$\partial L / \partial W^{(l+1)} = \partial L / \partial z^{(l+1)} \times (a^{(l+1)})^T$$

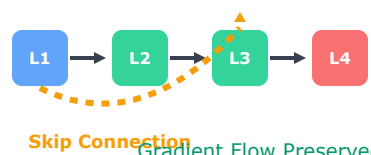


## Skip Connections

### Without Skip Connections



### With Skip Connections



$$y = F(x) + x$$

### Key Benefits:

- ✓ Better gradient flow
- ✓ Enables very deep networks
- ✓ Easy to learn identity mapping

**Skip connections solve the vanishing gradient problem by providing shortcut paths for gradient flow during backpropagation.**



## Evaluation Metrics Guide

### Confusion Matrix

		Predicted	
		Pos	Neg
Actual	Pos	TP	FP
	Neg	FN	TN

### Choose Your Metric:

Medical Diagnosis → High Recall (don't miss diseases) | Spam Detection → High Precision (avoid false positives) | Balanced Dataset → Accuracy

### Classification Metrics

#### Key Metrics:

**Accuracy:**  $(TP + TN) / (TP + TN + FP + FN)$   
"Overall correctness"

**Precision:**  $TP / (TP + FP)$   
"When I predict positive, how often am I right?"

**Recall:**  $TP / (TP + FN)$   
"How many actual positives did I catch?"

**F1-Score:**  $2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$   
"Harmonic mean of precision and recall"

**Roc-Auc:** Area under ROC curve (TPR vs FPR)  
"Measures separability across all thresholds"

### Regression Metrics

#### Mean Absolute Error

$$MAE = \frac{1}{n} \sum |y_i - \hat{y}_i|$$

✓ Easy to interpret (same units as target)

✓ Robust to outliers

#### Mean Squared Error

$$MSE = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

✓ Penalizes large errors more heavily

✓ Smooth gradient

#### Root Mean Squared Error

$$RMSE = \sqrt{\frac{1}{n} \sum (y_i - \hat{y}_i)^2}$$

✓ Same units as target

✓ Penalizes large errors

#### R<sup>2</sup> Score (Coefficient of Determination)

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

✓ Scale-independent

✓ Easy interpretation



## Best Practices Quick Reference

### Architecture Design

- Start simple: 1-2 hidden layers first
- Layer sizes: Powers of 2 (64, 128, 256, 512)
- Activations: ReLU for hidden, task-specific for output
- Add skip connections for networks deeper than 5 layers

### Training Tips

- Use Adam optimizer with lr=0.001 as default
- Batch size: 32-128 for most problems
- Monitor both training and validation metrics
- Use early stopping to prevent overfitting

### Debugging Guide

- Check learning curves for overfitting/underfitting
- Monitor gradient magnitudes
- Visualize weight distributions
- Start with smaller learning rates if loss explodes



## Debugging & Troubleshooting

### ⚠️ Vanishing Gradients

**Symptoms:** Deep networks don't train • Early layers learn slowly

**Solutions:** Use ReLU activation • Add skip connections • Proper weight initialization

### ⚠️ Exploding Gradients

**Symptoms:** Loss becomes NaN • Weights grow very large

**Solutions:** Gradient clipping • Lower learning rate • Better weight initialization

### ⚠️ Overfitting

**Symptoms:** High training accuracy, low validation accuracy • Large gap between training and validation loss

**Solutions:** Add dropout • L2 regularization • More training data



**Remember: Start simple, validate early, iterate often**

Keep experimenting and learning! 🚀