

Comprehensive AI/ML Study Notes

Modules 1-4: Complete Guide for Semester Examinations

TABLE OF CONTENTS

1. Module 1: Machine Learning Fundamentals
 2. Module 2: Deep Learning Fundamentals
 3. Module 3: Deep Learning Architectures (CNN, RNN, LSTM)
 4. Module 4: Deep Learning for NLP
 5. Module 5: Machine Learning Case Studies
-

MODULE 1: MACHINE LEARNING FUNDAMENTALS

I. What is Machine Learning?

Definition: Machine Learning (ML) is the science of programming computers so they can learn from data and optimize a performance criterion without being explicitly programmed.

Key Characteristics:

- Subdomain of Artificial Intelligence (AI)
- Enables systems to learn patterns from data
- Improves performance through experience
- Widely used in image classification, clustering, prediction tasks

When to Use Machine Learning:

- ✓ When rules cannot be written manually
- ✓ When patterns keep changing over time
- ✓ When large amounts of data are available
- ✓ When accuracy requirements are moderate

When NOT to Use Machine Learning:

- ✗ When the task is simple and rule-based
- ✗ When you don't have enough training data
- ✗ When perfect accuracy is absolutely required (safety-critical systems)

II. Machine Learning Pipeline: 7-Step Process

Step 1: Data Collection

- Gather relevant data for the problem
- Ensure sufficient quantity and quality
- Identify data sources and collection methods

Step 2: Data Preparation

- Clean data (handle missing values, outliers)
- Normalize/scale features to comparable ranges
- Encode categorical variables
- Handle imbalanced datasets

Step 3: Choosing a Model

- Select appropriate algorithm based on problem type
- Consider trade-offs (accuracy vs interpretability, training time)
- Decide between simple vs complex models

Step 4: Training

- Feed prepared data to the model
- Model learns patterns and updates parameters
- Optimize loss function using gradient descent

Step 5: Evaluation

- Test model on unseen data (test set)
- Calculate performance metrics (accuracy, precision, recall, F1)
- Compare against baseline

Step 6: Hyperparameter Tuning

- Adjust learning rate, regularization, batch size
- Use techniques: Grid Search, Random Search, Bayesian Optimization
- Improve validation performance

Step 7: Prediction & Deployment

- Use trained model for real-world predictions
- Monitor performance over time
- Retrain periodically with new data

III. Types of Machine Learning

1. Supervised Learning

Aspect	Details
Data	Labeled data with input-output pairs
Goal	Predict outputs for unseen data
Learning	Learns mapping from input → output
Examples	Classification (email spam detection), Regression (house price prediction)
When	When labels are available

Sub-types:

- **Classification:** Predicting discrete categories (Yes/No, 0/1/2/...)
 - Logistic Regression, Decision Trees, SVM, Random Forest, Neural Networks
- **Regression:** Predicting continuous values
 - Linear Regression, Polynomial Regression, Support Vector Regression

2. Unsupervised Learning

Aspect	Details
Data	Unlabeled data (no target variable)
Goal	Find hidden structures and patterns
Learning	Groups data based on similarity
Examples	Clustering (customer segmentation), Dimensionality Reduction (PCA)
When	Labels are not available or expensive to obtain

Sub-types:

- **Clustering:** Grouping similar data points
 - K-Means, Hierarchical Clustering, DBSCAN
- **Dimensionality Reduction:** Reducing feature count while preserving information
 - PCA (Principal Component Analysis), t-SNE, Autoencoders

3. Semi-Supervised Learning

Aspect	Details
Data	Few labeled samples + many unlabeled samples
Goal	Improve learning with limited labels
Learning	Combines supervised and unsupervised approaches
Example	Web page classification with partial labels
Benefit	Reduces labeling costs while improving accuracy

Techniques:

- Self-training: Model predicts on unlabeled data, learns from high-confidence predictions
- Co-training: Multiple models trained on different feature subsets
- Graph-based methods: Propagate labels through data manifold

4. Reinforcement Learning

Aspect	Details
Data	Agent interacts with environment, receives rewards/punishments
Goal	Learn policy that maximizes cumulative reward
Learning	Improves through trial and error
Example	Robotics, game playing (AlphaGo), autonomous driving
Signal	Reward or punishment feedback

Key Concepts:

- **Agent**: Entity learning from environment
- **Environment**: System where agent operates
- **State**: Current situation description
- **Action**: Decision made by agent
- **Reward**: Feedback signal (positive/negative)
- **Policy**: Mapping from state → action

IV. Bias, Variance, and the Bias-Variance Tradeoff

Bias

Definition: Bias is the error introduced by approximating a real-world problem with a simplified model.

- **High Bias** (Underfitting):
 - Model is too simple to capture underlying patterns
 - High error on both training and test data
 - Often occurs with: linear models on complex data, insufficient features
 - Symptoms: consistently poor performance everywhere
 - Solution: Use more complex models, add features, reduce regularization
- **Low Bias:**
 - Model flexible enough to capture true patterns
 - Can fit training data well

Variance

Definition: Variance is the model's sensitivity to fluctuations in training data.

- **High Variance** (Overfitting):
 - Model memorizes noise in training data instead of learning general patterns
 - Very low training error but high test error
 - Often occurs with: complex models, insufficient regularization, small dataset
 - Symptoms:
 - Training accuracy \gg Test accuracy (large gap)
 - Model learns spurious patterns
 - Poor performance on new data
 - Solutions:
 - Use regularization (L1, L2)
 - Increase training data
 - Use dropout in neural networks
 - Apply early stopping
 - Reduce model complexity
 - Use cross-validation
- **Low Variance:**
 - Model produces similar predictions across different training sets
 - More generalizable to new data

Bias-Variance Tradeoff

$$\text{Total Error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

Sweet Spot: Balance between underfitting and overfitting

- Too simple (High Bias): Underfits
- Too complex (High Variance): Overfits
- Optimal: Minimal Total Error

Visual Representation:

- **High Bias, Low Variance:** Predictions clustered far from target (consistently wrong)

- **Low Bias, High Variance:** Predictions scattered widely around target (inconsistent)
- **Optimal:** Predictions clustered near target (accurate and consistent)

V. Overfitting and Underfitting

Overfitting

Definition: Model learns noise in training data instead of generalizable patterns.

Characteristics:

- Very low training error, high test error
- Model memorizes the dataset rather than learning patterns
- High variance, low bias
- Complex model on small/noisy dataset

Indicators:

- Training accuracy significantly higher than validation accuracy
- Model performs well on training data but poorly on test data

Solutions:

1. Regularization (L1/L2):

- L1 (Lasso): Adds penalty $\lambda \sum |w|$ to loss, encourages sparsity (zero weights)
- L2 (Ridge): Adds penalty $\lambda \sum w^2$ to loss, keeps weights small
- Prevents weights from becoming too large

2. Dropout:

- Randomly deactivates neurons during training with probability p
- Prevents co-dependency between neurons
- Forces network to learn redundant representations
- Effective in deep neural networks

3. Early Stopping:

- Monitor validation loss during training
- Stop when validation loss stops improving for N epochs
- Prevents model from training too long and memorizing noise

4. Collect More Training Data:

- More data provides better generalization
- Helps model distinguish noise from signal

5. Reduce Model Complexity:

- Use simpler model architecture
- Fewer parameters = less memorization capability
- Trade-off with underfitting

Underfitting

Definition: Model is too simple to capture underlying patterns in data.

Characteristics:

- High error on both training and test data
- High bias, low variance
- Simple model on complex dataset

Indicators:

- Both training and validation accuracy are poor
- Model fails to improve with more training

Solutions:

1. Increase model complexity
2. Add more relevant features
3. Reduce regularization strength
4. Train for more epochs
5. Use more sophisticated algorithm

VI. Cross-Validation

Definition: Statistical technique to evaluate how well a model generalizes to unseen data. Instead of a single train-test split, cross-validation uses multiple splits to ensure robust performance estimation.

Why Cross-Validation is Important:

- Single train-test split is unreliable (depends on random seed)
- Helps detect overfitting
- Works well with small datasets
- Provides stable, reproducible performance estimates
- Uses entire dataset for both training and testing (different folds)

K-Fold Cross-Validation Process

Steps:

1. Shuffle the entire dataset randomly
2. Divide data into k equal-sized parts (folds)
3. For each fold i from 1 to k :
 - o Use fold i as test set
 - o Use remaining $k-1$ folds as training set
 - o Train model on training set
 - o Evaluate on test set
 - o Record performance metric
4. Calculate average performance across all folds
5. Final accuracy = Mean of all fold scores

Example (5-Fold CV):

Fold 1: Train on [2,3,4,5] → Test on [1]

Fold 2: Train on [1,3,4,5] → Test on [2]

Fold 3: Train on [1,2,4,5] → Test on [3]

Fold 4: Train on [1,2,3,5] → Test on [4]

Fold 5: Train on [1,2,3,4] → Test on [5]

$$\text{Average Accuracy} = (\text{Acc}_1 + \text{Acc}_2 + \text{Acc}_3 + \text{Acc}_4 + \text{Acc}_5) / 5$$

Cross-Validation Variants

1. K-Fold Cross-Validation

- Standard approach
- Typical k values: 5, 10
- Good balance between bias and computation

2. Stratified K-Fold

- Maintains class distribution in each fold
- Important for imbalanced datasets
- Ensures each fold has representative class proportions

3. Leave-One-Out Cross-Validation (LOOCV)

- $k = n$ (number of samples)
- Each iteration: 1 sample for testing, $n-1$ for training
- Computationally expensive but unbiased
- Better for very small datasets

4. Hold-Out Validation

- Single split: typically 80-20 or 70-30
- Fastest but unreliable
- Should be avoided for small datasets

5. Time Series Cross-Validation

- For temporal data, maintain temporal order
- Forward-chaining: Training set increases, test set moves forward
- Prevents information leakage from future to past

Advantages of Cross-Validation

- ✓ Reduces bias in performance estimation
- ✓ Reduces variance by averaging k results
- ✓ Works well with small datasets
- ✓ Detects overfitting
- ✓ Uses entire dataset for both training and evaluation

Disadvantages of Cross-Validation

- ✗ Computationally expensive (trains model k times)
- ✗ Not ideal for time series data
- ✗ Very slow for large datasets
- ✗ May not be practical with very complex models

VII. Feature Engineering

Definition: Process of transforming raw data into meaningful input features that help ML models learn better patterns and improve accuracy.

Importance:

- Directly impacts model accuracy and performance
- Reduces overfitting by removing irrelevant features
- Speeds up training time
- Handles missing or noisy data
- Converts non-linear relationships into linear form

Feature Engineering Steps

A. Data Cleaning

1. Handling Missing Values:

- Deletion: Remove rows/columns with missing values
- Imputation: Fill with mean, median, mode, or predictive model
- Interpolation: For time series data
- Choose based on missing data percentage and importance

2. Handling Outliers:

- Identify: Use Z-score ($|z| > 3$), IQR method, or visualization
- Treatment:
 - Remove outliers (if data entry errors)
 - Keep outliers (if they represent true anomalies)
 - Transform using log/sqrt to reduce impact
 - Use robust statistics (median instead of mean)

B. Feature Construction (Creating New Features)

1. Polynomial Features:

- Create interaction terms: x_1x_2, x_1^2, x_2^2
- Capture non-linear relationships
- Example: Price prediction might benefit from (area²)

2. Temporal Features (for time series):

- Extract: day, month, year, day-of-week, hour, minute
- Lag features: Previous values of time series
- Rolling averages: Moving window statistics

3. Domain-Specific Features:

- Product: Create features like price/quantity ratio
- Text: Word count, unique words, sentiment score
- Image: Edge detection, color histograms, texture

C. Feature Selection

- Remove irrelevant or redundant features
- Methods: Filter (correlation, variance), Wrapper (recursive elimination), Embedded (L1 regularization)
- Reduces training time and improves interpretability

Feature Engineering Challenges

- ✗ Requires deep domain knowledge
- ✗ Time-consuming and laborious
- ✗ Not always automated
- ✗ Can create spurious relationships
- ✗ Risk of information leakage

Best Practices

- ✓ Understand business domain deeply
- ✓ Analyze data distributions and relationships
- ✓ Create interpretable features
- ✓ Validate feature importance
- ✓ Document feature creation logic
- ✓ Avoid data leakage (don't use test set information in training)

VIII. Gradient Descent

Definition: Optimization algorithm used to find optimal values of model parameters (weights) that minimize the cost/loss function.

Intuition: Imagine a ball rolling down a hill. Gradient descent is the process of taking steps downhill until reaching the valley (minimum).

How Gradient Descent Works

Mathematical Foundation:

- **Cost Function:** $J(W)$ measures prediction error
- **Gradient:** $\nabla J(W) = \partial J / \partial W$ (slope/direction of steepest increase)
- **Goal:** Move in direction opposite to gradient (steepest descent)

Algorithm Steps:

1. Initialize parameters W to small random values (or zeros)
2. Compute cost function: $J(W)$
3. Compute gradient (derivative): $\partial J / \partial W$
4. Update parameters: $W_{\text{new}} = W - \alpha \cdot \nabla J(W)$
 - o α = learning rate (controls step size)
5. Repeat steps 2-4 until convergence ($\nabla J(W) \approx 0$)

Learning Rate (α): Critical Hyperparameter

Impact of Learning Rate:

Learning Rate	Effect	Result
Too Small	Takes tiny steps	Converges slowly, computationally expensive
Too Large	Takes huge jumps	May overshoot minimum, diverge, oscillate
Optimal	Balanced steps	Fast convergence to minimum

Learning Rate Scheduling (Adapt α during training):

1. **Learning Rate Annealing/Decay:**

- Reduce α as training progresses
 - Large steps initially, fine-tune later
 - Formula: $\alpha_t = \alpha_0 / (1 + \text{decay} \cdot t)$
- 2. Exponential/Cosine Decay:**
- Exponential: $\alpha_t = \alpha_0 \cdot e^{(-\text{decay} \cdot t)}$
 - Cosine: $\alpha_t = \alpha_0 (1 + \cos(\pi t/T))/2$
- 3. Reduce on Plateau:**
- Reduce α when validation loss plateaus
 - Helps escape local minima
- 4. Warm-up:**
- Start with small α , gradually increase
 - Prevents training instability early on

Advantages of Gradient Descent

- ✓ Simple and easy to implement
- ✓ Works for large-scale problems
- ✓ Can optimize non-linear, multi-parameter models
- ✓ Computationally efficient

Limitations of Gradient Descent

- ✗ Sensitive to learning rate selection
- ✗ Can get stuck in local minima (non-convex functions)
- ✗ Slow convergence for large datasets
- ✗ Requires computing derivatives

IX. Loss Functions

Definition: Mathematical function that measures how wrong the model's predictions are compared to actual target values. Used to guide model training by updating weights.

Common Loss Functions

1. Mean Squared Error (MSE)

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

Characteristics:

- Penalizes large errors heavily (quadratic penalty)
- Not robust to outliers
- Single global minimum (convex)
- Used for regression tasks
- Smooth gradient for optimization

2. Mean Absolute Error (MAE)

$$L(y, \hat{y}) = (1/n) \cdot \sum |y - \hat{y}|$$

Characteristics:

- Measures absolute distance between prediction and actual
- More robust to outliers than MSE
- Linear penalty for errors

- Sharper gradients near zero
- Used for regression

3. Huber Loss

$$L(y, \hat{y}) = \begin{cases} 0.5 \cdot (y - \hat{y})^2, & \text{if } |y - \hat{y}| \leq \delta \text{ (Quadratic)} \\ \delta \cdot (|y - \hat{y}| - 0.5 \cdot \delta), & \text{if } |y - \hat{y}| > \delta \text{ (Linear)} \end{cases}$$

Characteristics:

- Combines MSE (small errors) and MAE (large errors)
- More robust than MSE, smoother than MAE
- Useful when some outliers are expected
- Hyperparameter δ controls transition

4. Binary Cross-Entropy (Log Loss)

$$L(y, p) = -[y \cdot \log(p) + (1-y) \cdot \log(1-p)]$$

Where:

- y = true label (0 or 1)
- p = predicted probability

Characteristics:

- Measures uncertainty using entropy concept
- Loss decreases as predicted probability becomes more accurate
- Used for binary classification
- Outputs meaningful probabilities
- Prevents sharp predictions

5. Categorical Cross-Entropy (Multi-class)

$$L(Y, P) = -\sum_i y_i \cdot \log(p_i)$$

Characteristics:

- Generalizes binary cross-entropy to multi-class
- Used with softmax activation
- Each true class has probability 1, others 0
- Effective for multi-class classification

6. Hinge Loss (SVM Loss)

$$L(y, f(x)) = \max(0, 1 - y \cdot f(x))$$

Characteristics:

- Encourages correct prediction with margin
- Not only correct prediction, but confident prediction
- Penalizes even correct predictions if too close to decision boundary
- Used in SVM and margin-based learning
- Creates well-separated decision boundaries

Choosing Loss Function

Problem Type	Loss Function	Reason
Regression (continuous)	MSE, MAE, Huber	Predict continuous values
Binary Classification	Binary Cross-Entropy	Probabilistic output
Multi-class Classification	Categorical Cross-Entropy	Multiple mutually exclusive classes
Imbalanced Classification	Weighted Cross-Entropy, Focal Loss	Handle class imbalance
Ranking/Margin-based	Hinge Loss, Triplet Loss	Maximize margin between classes

MODULE 2: DEEP LEARNING FUNDAMENTALS

I. Machine Learning vs Deep Learning

Comparison Table

Aspect	Machine Learning	Deep Learning
Features	Hand-crafted by humans	Learns automatically from data
Model	Mostly shallow models	Deep networks with many hidden layers
Data	Works with small to medium	Needs large datasets (100K+)
Computation	CPU sufficient	GPU/TPU required for efficiency
Training	Easier, faster	Complex, slower
Interpretability	More interpretable	"Black box", harder to interpret
Feature Engineering	Manual, domain-dependent	Automatic feature learning
Performance	Good for simple tasks	Excellent for complex patterns
Example	Decision Trees, SVM, Random Forest	CNN, RNN, LSTM, Transformers

II. Neural Networks: Foundation of Deep Learning

Definition: Computational model inspired by biological neurons in the human brain. Consists of interconnected units (neurons) arranged in layers that process information through non-linear transformations.

Key Insight: Neural networks use non-linear mapping to compute complex decision boundaries, enabling them to learn intricate patterns that linear models cannot.

Neural Network Architecture Components

A. Neurons (Artificial Neuron)

Definition: Basic computational unit that receives inputs, multiplies by weights, adds bias, and applies activation function.

Mathematical Operation:

$$z = \sum(w_i \cdot x_i) + b$$

$$a = f(z) \text{ [activation function]}$$

Where:

- x_i = input signals
- w_i = weights (connection strengths)
- b = bias (shifts activation function)
- z = weighted sum
- $f(\cdot)$ = activation function
- a = output

Purpose:

- Weights learn importance of each input
- Bias allows network to learn even when all inputs are zero
- Non-linearity from activation function enables complex learning

B. Layers

1. Input Layer

- Receives raw data from environment
- No computation performed
- Only passes information forward
- Number of neurons = number of input features

2. Hidden Layers

- Perform most computation
- Apply weights, biases, activation functions
- Extract hierarchical features
- Multiple layers enable learning complex patterns
- Number of hidden layers determines network depth

3. Output Layer

- Produces final prediction
- Activation depends on task:
 - Binary classification: 1 neuron with sigmoid
 - Multi-class: n neurons with softmax
 - Regression: 1 neuron with linear activation

C. Weights

- **Definition:** Parameters representing connection strength between neurons
- **Range:** Typically continuous values in $[-1, 1]$ or wider
- **Learning:** Updated during training via backpropagation
- **Initialization:** Important for training stability
 - Random initialization breaks symmetry
 - Too large/small values cause problems
 - Common: Xavier, He initialization

D. Biases

- **Definition:** Additional parameters added to weighted sum
- **Purpose:** Allows network to learn translations (shifts in data)
- **Equation:** $z = Wx + b$
- **Importance:** Enables network to learn patterns even when $x=0$

Activation Functions

Purpose: Introduce non-linearity, enabling networks to learn complex patterns. Without activation, stacking layers gives same result as single layer (linear composition).

1. Sigmoid Function

$$\sigma(z) = 1 / (1 + e^{-z})$$

Range: $(0, 1)$

Properties:

- S-shaped curve
- Smooth, differentiable
- Squashes values to probability-like range $[0,1]$
- Historical popularity (logistic regression)

Disadvantages:

- Saturates (gradient ≈ 0 at extremes) → vanishing gradient
- Not zero-centered (outputs mostly positive)
- Slow convergence

2. Tanh (Hyperbolic Tangent)

$$\tanh(z) = (e^z - e^{-z}) / (e^z + e^{-z}) = 2\cdot\sigma(2z) - 1$$

Range: $(-1, 1)$

Properties:

- S-shaped, steeper than sigmoid
- Zero-centered output
- Stronger gradients than sigmoid

Disadvantages:

- Still suffers from saturation and vanishing gradient
- Better than sigmoid but not perfect

3. ReLU (Rectified Linear Unit) * Most Popular

$$f(z) = \max(0, z) = \{$$

z , if $z > 0$

0 , if $z \leq 0$

}

Properties:

- Simple, computationally efficient
- Non-zero gradient for positive inputs (prevents saturation)
- Encourages sparse representations

Disadvantages:

- Not differentiable at $z=0$ (though practically negligible)
- "Dying ReLU" problem: neurons can become inactive (output 0 for all inputs)

4. Leaky ReLU

$$f(z) = \max(\alpha \cdot z, z) = \begin{cases} z, & \text{if } z > 0 \\ \alpha \cdot z, & \text{if } z \leq 0 \end{cases}$$

Where α is small constant (e.g., 0.01)

Properties:

- Fixes "dying ReLU" problem
- Small negative slope allows gradient flow
- More stable training than ReLU

Advantage over ReLU:

- Neurons can still update even when inactive

5. ELU (Exponential Linear Unit)

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

Properties:

- Smoother than Leaky ReLU
- Negative saturation (approaches $-\alpha$)
- Helps center activations

6. Linear Activation

$$f(z) = z$$

Properties:

- No transformation
- Used in output layer for regression
- Preserves value ranges

When to Use Which Activation

Layer Type	Activation	Reason
Hidden	ReLU (default)	Fast, prevents saturation
Hidden	Leaky ReLU, ELU	When ReLU dies
Output (Binary)	Sigmoid	Probability output [0,1]
Output (Multi-class)	Softmax	Probability distribution
Output (Regression)	Linear	Unrestricted output

III. Neural Network Training: Forward and Backward Propagation

1. Forward Pass (Propagation)

- Data flows from input layer through hidden layers to output layer
- Each neuron computes: $z = Wx + b$, $a = f(z)$
- Final layer produces prediction \hat{y}

Mathematical Notation:

- Layer l : $z^l = W^{l-1}a^{l-1} + b^l$
- Activation: $a^l = f^l(z^l)$

2. Loss Function Computation

- Compare prediction with ground truth
- Example: $L = (\hat{y} - y)^2$ for MSE

3. Backward Pass (Backpropagation)

Definition: Algorithm to compute gradients of loss w.r.t. all parameters, enabling weight updates.

Core Idea: Chain rule of calculus

- $\partial L / \partial W = \partial L / \partial a^l \cdot \partial a^l / \partial z^l \cdot \partial z^l / \partial W$
- Errors propagate backward through network

Steps:

1. Compute $\partial L / \partial a^l$ (error at output layer)
2. For each layer l from L to 1:
 - Compute $\partial L / \partial z^l = \partial L / \partial a^l \cdot f'(z^l)$
 - Compute $\partial L / \partial W^l$, $\partial L / \partial b^l$
 - Pass gradient to previous layer: $\partial L / \partial a^{l-1} = (W^l)^T \cdot \partial L / \partial z^l$

4. Weight Updates

- $W_{\text{new}} = W - \alpha \cdot \partial L / \partial W$
- $b_{\text{new}} = b - \alpha \cdot \partial L / \partial b$
- Repeat forward-backward cycle (epochs)

IV. Training Parameters and Configuration

Training Components

1. **Weights (W)**: Updated during training
2. **Biases (b)**: Updated during training
3. **Learning Rate (α)**: Hyperparameter, controls step size
4. **Batch Size**: Number of samples per update
5. **Epochs**: Number of complete passes through dataset
6. **Regularization**: Prevent overfitting

V. Learning Rate: Critical Hyperparameter

Definition: Controls how big a step optimizer takes in direction opposite to gradient.

Formula: $W_{\text{new}} = W - \alpha \cdot \nabla L(W)$

Effects of Learning Rate

Learning Rate	Effect	Outcome
Too Small	Takes tiny steps downhill	Converges very slowly, computationally expensive
Too Large	Takes huge jumps	Overshoots minimum, diverges, loss increases
Optimal	Balanced steps	Fast convergence to good minimum

Visual Intuition

- Imagine walking down a mountain blindfolded
- Too small steps: Takes forever to reach valley
- Too large steps: Jump over valley, go up other side
- Optimal: Steady descent reaching valley quickly

Learning Rate Scheduling

Adapt learning rate during training:

1. Step Decay

- Reduce α by factor every N epochs
- Example: $\alpha = \alpha_0 / (1 + \text{decay} \cdot \text{epoch})$

2. Exponential Decay

- $\alpha = \alpha_0 \cdot e^{(-\text{decay} \cdot \text{epoch})}$

3. Cosine Decay

- Smooth decrease following cosine curve
- $\alpha = \alpha_0 \cdot (1 + \cos(\pi \cdot \text{epoch}/\text{max_epoch}))/2$

4. Reduce on Plateau

- Monitor validation loss
- If loss doesn't improve, reduce α

5. Warm-up

- Start with small α , gradually increase over first epochs
- Stabilizes training in early phases

VI. Regularization: Preventing Overfitting

Goal: Prevent model from memorizing training data and improve generalization.

1. L2 Regularization (Ridge)

$$L_{\text{reg}} = L(X, y) + \lambda/2 \cdot \sum(W^2)$$

Effect:

- Adds penalty proportional to weight magnitudes
- Forces weights toward zero
- Large weights heavily penalized
- Encourages distributed feature usage

When to use:

- When you suspect most features are relevant
- Want to reduce all weights evenly

Gradient Update:

- $W_{\text{new}} = W - \alpha \cdot (\partial L / \partial W + \lambda \cdot W)$
- Weights pulled toward zero each update

2. L1 Regularization (Lasso)

$$L_{\text{reg}} = L(X, y) + \lambda \cdot \sum |W|$$

Effect:

- Penalty proportional to absolute weight values
- Encourages sparsity (many weights $\rightarrow 0$)
- Selects relevant features

When to use:

- When you want feature selection
- Suspect only some features are important
- Want interpretable sparse models

Difference from L2:

- L1: Sharp corner at zero (encourages zeros)
- L2: Smooth curve (encourages small values)

3. Dropout

During training:

- Randomly deactivate (drop) neurons with probability p
- Share of neurons removed: p (often 0.2-0.5)
- Forward pass: Set dropped neuron outputs to 0
- Backward pass: Don't update dropped neurons

Effect:

- Prevents co-adaptation of neurons
- Reduces complex interdependencies
- Each neuron learns robust features independently
- At test time: Use all neurons, scale activations by $(1-p)$

Intuition:

- Like training ensemble of smaller networks
- Forces learning of redundant representations
- Each neuron must be useful independently

How It Works:

1. Create random binary mask $M \sim \text{Bernoulli}(1-p)$
2. Forward: $a_{\text{dropout}} = a \odot M / (1-p)$
3. Backward: Gradients only for active neurons

4. Early Stopping

Monitoring:

1. Split data: Train / Validation
2. Train model, track validation loss
3. If validation loss increases for N consecutive epochs:
 - Stop training immediately
 - Revert to best checkpoint

Effect:

- Prevents model from training too long
- Escapes point where training fits noise
- Simple but effective

5. Data Augmentation

- Artificially increase training set size
- Apply transformations: rotation, flip, zoom, noise
- Helps model learn invariances

Comparison of Regularization Methods

Method	How It Works	When to Use	Strength
L2	Penalize weight magnitude	General purpose	Distributed weight reduction
L1	Penalize absolute weights	Feature selection needed	Sparse solutions
Dropout	Random neuron removal	Deep networks	Prevents co-adaptation
Early Stop	Monitor validation loss	Any model	Simple, effective

VII. Optimizers: Beyond Basic Gradient Descent

Problem with Basic GD: Uniform learning rate for all parameters often suboptimal. Different parameters may need different learning rates.

Variants of Gradient Descent

1. Batch Gradient Descent (BGD)

- Computes gradient using entire training dataset
- One weight update per epoch
- Very stable but slow for large datasets
- Guaranteed to converge (for convex functions)

2. Stochastic Gradient Descent (SGD)

- Updates weights after each single sample
- Very fast, noisy updates
- Escapes local minima (noise helps)
- High variance, training unstable
- May oscillate around optimum

3. Mini-batch Gradient Descent

- Updates using small batch of samples (e.g., 32, 64, 128)
- Best of both worlds: speed + stability

- Parallelizable, GPU-friendly
- Most commonly used in practice

Advanced Optimizers

1. Momentum

Equation:

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

Effect:

- Accumulates gradient over time (momentum term)
- Accelerates movement in consistent directions
- Dampens oscillations perpendicular to gradient
- Smooths out noisy updates

Intuition:

- Like a ball rolling downhill gaining speed
- Maintains velocity from previous steps
- Helps escape shallow local minima

Typical β values: 0.9 (90% of previous velocity)

2. Nesterov Accelerated Momentum

Equation:

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

Improvement over Momentum:

- "Look ahead" before updating
- Evaluates gradient at future approximated position
- Faster convergence than standard momentum
- Less oscillation

3. AdaGrad (Adaptive Gradient)

Equation:

- $G_t = G_{t-1} + (\nabla L)^2$ [accumulate squared gradients]
- $W = W - (\alpha / \sqrt{G_t + \epsilon}) \cdot \nabla L(W)$

Effect:

- Adaptive learning rate per parameter
- Large gradients \rightarrow smaller step (divide by large G)
- Small gradients \rightarrow larger step (divide by small G)
- Parameters with infrequent updates get larger steps

Advantage:

- No manual learning rate tuning

- Sparse data friendly (rare features get big updates)

Disadvantage:

- Learning rate decreases monotonically
- Eventually becomes too small (stops learning)
- Not ideal for non-convex problems

4. RMSProp (Root Mean Square Propagation)

Equation:

- $E[g^2]_t = \beta \cdot E[g^2](t-1) + (1-\beta) \cdot (\nabla L)^2$
- $W = W - (\alpha / \sqrt{E[g^2]_t + \epsilon}) \cdot \nabla L(W)$

Improvement over AdaGrad:

- Uses exponential moving average instead of sum
- Learning rate doesn't shrink forever
- Fixes AdaGrad's diminishing LR problem
- Divides by RMS of gradients

Effect:

- Adapts learning rate per parameter
- Combines benefits of momentum and AdaGrad
- Better convergence than AdaGrad

5. Adam (Adaptive Moment Estimation) ★ Most Popular

Equations:

- $m_t = \beta_1 \cdot m_{(t-1)} + (1-\beta_1) \cdot \nabla L(W)$ [1st moment - mean]
- $v_t = \beta_2 \cdot v_{(t-1)} + (1-\beta_2) \cdot (\nabla L)^2$ [2nd moment - variance]
- $\hat{m}_t = m_t / (1 - \beta_1^t)$ [bias correction]
- $\hat{v}_t = v_t / (1 - \beta_2^t)$ [bias correction]
- $W = W - \alpha \cdot (\hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon))$

Combines:

- Momentum (exponential moving average of gradients)
- RMSProp (exponential moving average of squared gradients)
- Bias correction (accounts for zero initialization)

Default values:

- $\beta_1 = 0.9$ (momentum coefficient)
- $\beta_2 = 0.999$ (RMSProp coefficient)
- $\epsilon = 10^{-8}$
- $\alpha = 0.001$ (typical)

Advantages:

- Works well across most DL tasks
- Computationally efficient
- Little hyperparameter tuning needed
- Handles sparse data well

- Fast convergence

Disadvantages:

- More complex (4 equations vs 1 for SGD)
- May generalize worse than SGD in some cases

Optimizer Comparison Table

Optimizer	Speed	Stability	Sparse Data	Comments
SGD	Fast	Noisy	Poor	Basic but solid
Momentum	Faster	Better	Poor	Smooths noise
Nesterov	Faster	Better	Poor	Momentum + lookahead
AdaGrad	Fast	Good	Excellent	Frequency-based
RMSProp	Fast	Good	Good	Fixes AdaGrad
Adam	Fast	Excellent	Good	Best general choice

VIII. Hyperparameter Tuning

Definition: Process of selecting optimal values for model hyperparameters (not learned from data).

Key Hyperparameters in Neural Networks

1. Learning Rate (α)

- Range: 0.0001 to 0.1
- Start with 0.001, adjust if training unstable

2. Batch Size

- Range: 16, 32, 64, 128, 256
- Larger batch → more stable, less frequent updates
- Smaller batch → noisy, escapes local minima faster

3. Number of Epochs

- Use early stopping to determine
- Monitor validation loss

4. Number of Layers (Depth)

- More layers → more capacity, risk of overfitting
- 2-5 hidden layers typical for many tasks

5. Neurons per Layer (Width)

- More neurons → more parameters, more computation
- Typical: 64-512 per layer

6. Regularization (λ , dropout p)

- λ range: 1e-5 to 0.1
- p range: 0.1 to 0.5

7. Optimizer

- Default: Adam
- Alternatives: SGD with momentum, RMSProp

Hyperparameter Tuning Methods

1. Grid Search

- Exhaustive search over predefined values
- Example: $\text{batch_size} \in \{32, 64, 128\}$, $\text{lr} \in \{0.001, 0.01, 0.1\}$
- Try all combinations
- Guarantees finding best in grid
- Computationally expensive

2. Random Search

- Sample hyperparameters randomly from distributions
- More efficient than grid search
- Often finds better values than grid

3. Bayesian Optimization

- Model performance as probability function
- Intelligently select next hyperparameters to try
- More efficient than random or grid
- Requires more setup

4. Manual Tuning

- Adjust based on training curves
- Intuition and experience matter
- Still widely used by practitioners

IX. Shallow vs Deep Networks

Comparison

Aspect	Shallow Network	Deep Network
Hidden Layers	1-2 layers	3+ layers
Features	Hand-engineered	Learned automatically
Complexity	Simple patterns	Complex hierarchical patterns
Representational Power	Limited	High
Training Ease	Easier	Harder (vanishing gradients)
Computation	Less computation	More computation, needs GPU
Data Requirement	Works with small datasets	Needs large datasets
Performance	Simple tasks	Complex tasks (images, text, etc.)
Interpretability	More interpretable	Less interpretable ("black box")
Underfitting Risk	Higher	Lower
Overfitting Risk	Lower	Higher (more capacity)

When to Use

Shallow Networks:

- ✓ Small datasets
- ✓ Simple relationships
- ✓ Need interpretability
- ✓ Limited computation

Deep Networks:

- ✓ Large datasets (100K+)
- ✓ Complex patterns (images, language)
- ✓ State-of-the-art performance
- ✓ GPU available

MODULE 3: DEEP LEARNING ARCHITECTURES

PART A: CONVOLUTIONAL NEURAL NETWORKS (CNN)

I. What are Convolutional Neural Networks?

Definition: Specialized feed-forward neural network designed to process data with grid-like topology (images). Takes advantage of 2D spatial structure by using local receptive fields, drastically reducing parameters compared to fully connected networks.

Key Innovation: Instead of connecting every neuron to every input, each neuron only connects to small local regions (windows) of input. This exploits **local connectivity** and **parameter sharing**.

II. How CNNs Work

Core Idea:

1. Use learnable filters (kernels) that slide across the image
2. Each filter performs **convolution operation** (dot product with local regions)
3. Detect local features (edges, corners, textures)
4. Stack multiple filters to build hierarchical features

Feature Detection Hierarchy:

- **Layer 1:** Low-level features (edges, colors, simple patterns)
- **Layer 2:** Mid-level features (textures, simple shapes)
- **Layer 3:** High-level features (object parts, complex shapes)
- **Output:** Object classification

III. CNN Architecture Components

1. Convolution Layer

Operation: Applies multiple learnable filters across input image.

Mathematical Formula:

$$\text{Output}[i,j,k] = \sum_m \sum_n \text{Input}[i+m, j+n] \times \text{Filter}_k[m,n] + \text{bias}_k$$

Parameters:

- **Filter Size:** Typically 3×3 , 5×5 , or 7×7
 - 3×3 most common (balance speed and receptive field)
- **Stride:** How many pixels filter moves each step
 - Stride=1: Dense feature extraction
 - Stride=2: Reduces spatial dimensions (faster)
- **Padding:** Adding zeros around input
 - 'Same' padding: Output size = input size
 - 'Valid' padding: Output size smaller than input
- **Number of Filters:** How many feature detectors

- Earlier layers: 32-64 filters
- Deeper layers: 128-512 filters

Output Size Calculation:

$$\text{Output_size} = (\text{Input_size} - \text{Filter_size} + 2 \times \text{Padding}) / \text{Stride} + 1$$

Example:

- Input: 28×28, Filter: 3×3, Padding: 1, Stride: 1
- Output: $(28 - 3 + 2) / 1 + 1 = 28 \times 28$ (same padding)

Advantages:

- ✓ Local receptive fields reduce parameters
- ✓ Parameter sharing (same filter applied everywhere)
- ✓ Detects features regardless of location (invariance)
- ✓ Hierarchical feature extraction

2. Activation Function (ReLU Layer)

Operation: Apply non-linearity after convolution.

Function: $f(z) = \max(0, z)$

Effect:

- Removes negative values
- Keeps computation efficient
- Enables learning of non-linear patterns

3. Pooling Layer

Purpose: Downsample feature maps to reduce spatial dimensions and computation.

Operation: Apply sliding window, compute statistic.

Types:

A. Max Pooling

$\text{Output}[i,j] = \max(\text{Input}[\text{window centered at } i,j])$

Effect:

- Keeps maximum activation (most important feature)
- Provides invariance to small spatial shifts
- Reduces parameters and computation

B. Average Pooling

$\text{Output}[i,j] = \text{average}(\text{Input}[\text{window centered at } i,j])$

Effect:

- Smoother, less aggressive than max pooling
- Preserves more information
- Less common than max pooling

Common Pool Size: 2×2 with stride 2

- Reduces spatial dimensions by half
- Reduces computation by 4x

Advantages:

- ✓ Reduces computational load
- ✓ Controls overfitting
- ✓ Provides spatial invariance
- ✓ Makes features more interpretable

4. Fully Connected (Dense) Layer

Purpose: Transform spatial features into class predictions.

Operation:

- Flatten pooled feature maps into 1D vector
- Connect to dense layer(s)
- Acts like standard neural network

Characteristics:

- Processes features globally (not locally like conv)
- Often includes dropout for regularization
- Final layer uses softmax for classification

IV. Typical CNN Architecture Sequence

Input ($28 \times 28 \times 3$ image)

↓

Convolution: 32 filters 3×3 , ReLU

↓

Pooling: 2×2 Max Pool

↓

Convolution: 64 filters 3×3 , ReLU

↓

Pooling: 2×2 Max Pool

↓

Flatten: Convert to 1D vector

↓

Dense: 128 neurons, ReLU

↓

Dropout: $p=0.5$

↓

Dense: 10 neurons, Softmax (output)

↓

Output: 10-class probabilities

V. Famous CNN Architectures

Architecture	Year	Key Innovation	Parameters
LeNet-5	1998	First successful CNN	60K
AlexNet	2012	Deep CNN, GPU training, ReLU	60M
VGG-16	2014	Stacked small filters, depth importance	138M
ResNet-50	2015	Skip connections, very deep	25M
MobileNet	2017	Lightweight, mobile-friendly	4M
EfficientNet	2019	Balanced efficiency	Variable

PART B: RECURRENT NEURAL NETWORKS (RNN)

I. What are Recurrent Neural Networks?

Definition: Neural network architecture designed for sequential/temporal data where previous outputs influence current prediction.

Key Characteristic: Contains feedback loops that allow information to persist across time steps, enabling the network to maintain memory of past inputs.

Difference from Feed-Forward NN:

- **Feed-Forward:** One-way information flow, no loops
- **Recurrent:** Loops allow information to cycle back

II. Why RNNs for Sequential Data?

Sequential Data Characteristics:

- **Time Series:** Stock prices, temperature, sensor data
- **Text:** Words depend on previous words for meaning
- **Speech:** Phonemes in context
- **Video:** Frames related to previous frames

Problem with Feed-Forward:

- Each input processed independently

- Cannot understand context or dependencies
- No memory of what came before

RNN Solution:

- Maintains hidden state across time steps
- Hidden state captures information from past
- Uses previous hidden state + current input for prediction

III. RNN Mathematics

Recurrence Relation:

$$h_t = f(h_{t-1}, x_t) \text{ [general form]}$$

$$h_t = \tanh(W_{hh} \cdot h_{t-1} + W_{xh} \cdot x_t + b_h) \text{ [specific form]}$$

$$y_t = W_{hy} \cdot h_t + b_y$$

Where:

- h_t = hidden state at time t (captures information)
- x_t = input at time t
- y_t = output at time t
- W_{hh} = weight matrix for hidden-to-hidden (recurrence)
- W_{xh} = weight matrix for input-to-hidden
- W_{hy} = weight matrix for hidden-to-output

Intuition: Hidden state h_t summarizes all information from past sequence.

IV. Unrolling RNNs Through Time

Concept: RNNs are unrolled across time steps to visualize computation.

At each time step:

- Use same W matrices (parameter sharing)
- Apply same recurrence relation
- Inputs: previous hidden state + current input
- Output: prediction + updated hidden state

Example (3 time steps):

$$x_1 \rightarrow \text{RNN} \rightarrow h_1 \rightarrow y_1$$

↓

$$h_1$$

↓

$$x_2 \rightarrow \text{RNN} \rightarrow h_2 \rightarrow y_2$$

↓

$$h_2$$

↓

$$x_3 \rightarrow \text{RNN} \rightarrow h_3 \rightarrow y_3$$

V. RNN Types Based on Input-Output Architecture

1. One-to-One: Fixed Input → Fixed Output

$x \rightarrow \text{RNN} \rightarrow y$

Use Case: Single input to single output (not truly sequential)

Example: Image classification

2. One-to-Many: Single Input → Sequence Output

$x \rightarrow \text{RNN} \rightarrow y_1 \rightarrow y_2 \rightarrow y_3$

Use Case: Generate sequence from single input

Example: Image captioning (image → words)

3. Many-to-One: Sequence Input → Single Output

$x_1 \rightarrow \text{RNN} \rightarrow (\text{aggregate}) \rightarrow y$

$x_2 \rightarrow \text{RNN} \nearrow$

$x_3 \rightarrow \text{RNN} \nearrow$

Use Case: Classification of entire sequence

Example: Sentiment analysis (words → sentiment)

4. Many-to-Many: Sequence → Sequence

$x_1 \rightarrow \text{RNN} \rightarrow y_1$

$x_2 \rightarrow \text{RNN} \rightarrow y_2$

$x_3 \rightarrow \text{RNN} \rightarrow y_3$

Use Case: Process and transform entire sequence

Example: Machine translation (English → French)

VI. Feed-Forward NN vs Recurrent NN

Aspect	Feed-Forward	Recurrent
Data Flow	One direction only	Contains feedback loops
Memory	No memory	Has memory (hidden state)
Input Processing	Each input independent	Context from previous inputs
Suitability	Non-sequential data	Sequential/temporal data
Parameter Sharing	None	Over time steps (same W used)
Output Dependency	Current only	Current + history
Example	Tabular data classification	Stock price prediction

VII. Training RNNs: Backpropagation Through Time (BPTT)

Challenge: Gradients flow backward through time steps, causing:

- **Vanishing Gradients:** Gradients become very small, learning stops
- **Exploding Gradients:** Gradients become very large, unstable training

Solution Techniques:

1. **Gradient Clipping:** Cap gradient magnitude to prevent explosions
2. **LSTM/GRU:** Special architectures with gating mechanisms
3. **Careful Initialization:** Initialize weights properly
4. **Learning Rate:** Use appropriate learning rate

VIII. Problems with Standard RNNs

1. Vanishing Gradients

Problem:

- Gradient multiplied by recurrent weight matrix at each step
- If weight < 1 , gradients shrink exponentially
- After many steps: gradient ≈ 0
- Network forgets long-term dependencies

Effect:

- Cannot learn long-term patterns
- Model mostly relies on recent inputs
- Struggles with long sequences

2. Exploding Gradients

Problem:

- If recurrent weight > 1, gradients grow exponentially
- Causes numerical instability
- Weight updates become massive

Effect:

- Training becomes unstable (loss diverges)
- Model outputs become NaN

PART C: LONG SHORT-TERM MEMORY (LSTM)

I. What is LSTM?

Definition: Enhanced version of RNN that can capture long-term dependencies in sequential data, making it ideal for tasks requiring understanding of context over long sequences.

Key Innovation: Introduces "gates" that control information flow, allowing gradients to flow through time without vanishing or exploding.

Motivation: Solves vanishing/exploding gradient problems of standard RNNs.

II. LSTM Architecture

Core Components:

A. Cell State (Memory)

C_t : Long-term memory that carries information through time

- Can be kept unchanged or modified at each step
- Enables information flow without degradation
- "Highway" for gradient flow

B. Hidden State (Short-term)

h_t : Short-term output, used for current prediction

- Filtered version of cell state
- Contains immediate information

C. Three Gates

All gates use sigmoid activation $\sigma(z) = 1/(1+e^{-z})$, output in [0,1]

1. Forget Gate: f_t

$$f_t = \sigma(W_f[h_{t-1}, x_t] + b_f)$$

Purpose: Decides what to forget from memory

- $f_t=0$: Forget everything (erase memory)
- $f_t=1$: Remember everything (keep memory)

- $f_t \in (0,1)$: Partially forget

Operation: $C_t = f_t \odot C_{(t-1)}$ [element-wise multiply]

2. Input Gate: i_t

$$i_t = \sigma(W_i \cdot [h_{(t-1)}, x_t] + b_i)$$

Purpose: Decides what new information to store

- $i_t=0$: Ignore new input
- $i_t=1$: Accept all new input

Also: Compute candidate values

$$\tilde{C}_t = \tanh(W_c \cdot [h_{(t-1)}, x_t] + b_c)$$

Operation: Update memory with new info

$$C_t = C_{(t-1)} + i_t \odot \tilde{C}_t$$

3. Output Gate: o_t

$$o_t = \sigma(W_o \cdot [h_{(t-1)}, x_t] + b_o)$$

Purpose: Decides what to output

- $o_t=0$: Output nothing
- $o_t=1$: Output full memory

Operation: $h_t = o_t \odot \tanh(C_t)$

III. LSTM Step-by-Step Operations

At time t:

1. Forget Gate:

$$f_t = \sigma(W_f \cdot [h_{(t-1)}, x_t] + b_f)$$

2. Input Gate:

$$i_t = \sigma(W_i \cdot [h_{(t-1)}, x_t] + b_i)$$

3. Candidate State:

$$\tilde{C}_t = \tanh(W_c \cdot [h_{(t-1)}, x_t] + b_c)$$

4. Update Cell State (Memory):

$$C_t = f_t \odot C_{(t-1)} + i_t \odot \tilde{C}_t$$

↑ ↑ ↑

forget previous add new info

5. Output Gate:

$$o_t = \sigma(W_o \cdot [h_{(t-1)}, x_t] + b_o)$$

6. Update Hidden State:

$$h_t = o_t \odot \tanh(C_t)$$

7. Prediction (if output at this step):

$$y_t = \text{softmax}(W_y \cdot h_t + b_y)$$

IV. Why LSTM Fixes Vanishing Gradients

Key Insight: Cell state C_t has additive update:

$$C_t = f_t \odot C_{(t-1)} + i_t \odot \tilde{C}_t$$

↑ ↑

forget gate add new info

Gradient flow:

$$\partial C_t / \partial C_{(t-1)} = f_t$$

If $f_t \approx 1$ (forget gate learns to keep memory):

- Gradient $\partial C_t / \partial C_{(t-1)} \approx 1$
- Gradients don't vanish (don't multiply by small values)
- Information flows through without degradation

Standard RNN Problem:

$$C_t = \tanh(W \cdot C_{(t-1)} + \dots)$$

$$\text{Gradient: } \partial C_t / \partial C_{(t-1)} = W \cdot (1 - \tanh^2(\dots))$$

- Often < 1 , so gradients shrink each step

V. LSTM Advantages

- ✓ **Captures long-term dependencies:** Information persists through cell state
- ✓ **Prevents vanishing gradients:** Additive path with multiplicative gates
- ✓ **Selective memory:** Gates control what to remember/forget
- ✓ **Flexible:** Can be stacked in multiple layers
- ✓ **Works well for:** Language modeling, machine translation, time series

VI. LSTM Training: Backpropagation Through Time (BPTT)

Backward Pass Through LSTM:

1. Compute gradient for forget gate weights
2. Compute gradient for input gate weights
3. Compute gradient for candidate state weights
4. Compute gradient for output gate weights
5. Update weights using gradient descent

Key Point: Gradients have cleaner path through cell state, reducing vanishing gradient problem.

VII. Variants and Extensions

Gated Recurrent Unit (GRU):

- Simplified LSTM with 2 gates instead of 3
- Combines forget and input gates
- Slightly faster, comparable performance
- Fewer parameters

Bidirectional LSTM (BiLSTM):

- Process sequence forward and backward
 - Captures context from both directions
 - Better for tasks like named entity recognition
 - Slightly increased computation
-

PART D: SEQUENCE-TO-SEQUENCE (SEQ2SEQ) MODELS

I. What is Seq2Seq?

Definition: Neural architecture that converts one sequence into another using Encoder-Decoder structure with RNN/LSTM/GRU.

Examples:

- Machine translation: English → French
- Summarization: Long document → Short summary
- Image captioning: Image features → Caption words
- Speech recognition: Audio → Text

II. Seq2Seq Architecture

A. Encoder

Process entire input sequence

↓

Extract information

↓

Produce context vector (final hidden state)

$h_1 \rightarrow h_2 \rightarrow h_3 \rightarrow h_4 \rightarrow$ Context Vector

$x_1 x_2 x_3 x_4$

Function:

- Takes variable-length input sequence
- Updates hidden state at each step
- Final hidden state = context vector
- Context vector summarizes entire input

Mathematical:

$$h_1 = \text{LSTM}(x_1, h_0)$$

$$h_2 = \text{LSTM}(x_2, h_1)$$

$$h_3 = \text{LSTM}(x_3, h_2)$$

...

Context = h_n (final hidden state)

B. Decoder

Start with context vector

↓

Generate output one token at time

↓

Use previous output + context for next prediction

Context → y_1
→ y_2
→ y_3
→ y_4

Function:

- Takes context vector from encoder
- Generates output sequence one token at a time
- Uses previous prediction to inform next
- Generates until EOS (end-of-sequence) token

Mathematical:

$h'_1 = \text{LSTM}(y_0, \text{Context})$ [y_0 = START token]

$y_1 = \text{softmax}(W h'_1)$

$h'_2 = \text{LSTM}(y_1, h'_1)$

$y_2 = \text{softmax}(W h'_2)$

... continue until EOS token

III. Teacher Forcing

Concept: During training, use ground truth previous token instead of model's prediction.

Benefits:

- ✓ Faster convergence (feeds correct token)
- ✓ Reduces error propagation (doesn't learn to handle its own mistakes)
- ✓ More stable training

Drawback:

- ✗ Mismatch between training (ground truth) and inference (predictions)
- ✗ Model can be brittle when fed its own predictions

Solution: Gradually reduce teacher forcing during training (curriculum learning).

IV. Challenges and Solutions

1. Information Bottleneck

Problem: Context vector must contain all info about input

- For long sequences, hard to compress into single vector
- Important info may be lost

Solution: Attention Mechanism

- Focus on relevant parts of input for each output
- Weight different input positions differently

2. Seq2Seq vs Attention

Without Attention:

- Encoder: $x_1, x_2, x_3, x_4 \rightarrow \text{Context}$

- Decoder: Context \rightarrow y_1, y_2, y_3, y_4
- Same context used for all outputs

With Attention:

- Decoder focuses on different inputs for each output
- Weights computed dynamically
- More flexible information flow

V. Applications of Seq2Seq

Application	Input	Output
Machine Translation	English sentence	French sentence
Summarization	Long document	Short summary
Image Captioning	Image features	English caption
Speech Recognition	Audio features	Text transcription
Chatbot	User message	Response
Time Series	Historical values	Future prediction

MODULE 4: DEEP LEARNING FOR NATURAL LANGUAGE PROCESSING

I. What is Natural Language Processing?

Definition: Branch of AI that enables computers to understand, interpret, analyze, and generate human language in meaningful ways.

Scope:

- Text processing and understanding
- Speech recognition and generation
- Meaning extraction and context understanding
- Language generation and translation

Key Challenge: Language is complex, ambiguous, context-dependent, evolving.

II. NLP Applications

Application	Task	Example
Machine Translation	Convert one language to another	Google Translate
Speech to Text	Convert audio to transcription	Voice assistants
Text Classification	Assign category to text	Email spam detection
Sentiment Analysis	Determine emotional tone	Product review polarity
Named Entity Recognition	Identify people, locations, organizations	Extract names from text
Question Answering	Answer natural language questions	Reading comprehension
Chatbots	Conversational AI	Virtual assistants
Text Summarization	Condense long text	Automatic article summary

III. NLP Approaches

1. Rule-Based Approach

Characteristic:

- Relies on predefined linguistic rules and patterns
- Experts/programmers manually create rules
- Encode language knowledge

Example:

If email contains ["viagra", "buy now"] → SPAM

If email from known sender AND contains response keywords → NOT SPAM

Advantages:

- Good for controlled domains
- Transparent and interpretable
- Quick to implement for simple rules

Disadvantages:

- Doesn't scale to complex tasks
- Requires expert knowledge
- Hard to capture all exceptions

- Maintenance burden

2. Statistical Approach

Characteristic:

- Learn patterns from labeled datasets
- Use probabilistic techniques
- Estimate probabilities from training data

Example:

- $P(\text{spam} \mid \text{words})$ computed from training emails
- Naive Bayes classifier learns word-to-class probabilities

Advantages:

- More flexible than rules
- Learns from data, reduces manual effort
- Handles some variations

Disadvantages:

- Requires sufficient labeled data
- Features still hand-crafted
- May not capture complex semantic relationships

3. Neural Network Approach

Characteristic:

- Uses RNN, LSTM, CNN for NLP
- Learns deep representations
- End-to-end learning from raw text

Example:

- Seq2Seq for machine translation
- LSTM for language modeling
- BERT for understanding

Advantages:

- Captures semantic and contextual meaning
- Handles long-term dependencies effectively
- State-of-the-art performance
- Automatic feature learning

Disadvantages:

- Requires large labeled datasets
- Computationally expensive
- Less interpretable

IV. NLP Pipeline: Key Components

Two Main Parts

A. NLU (Natural Language Understanding)

Focus: Enabling machines to understand language

Key Functions:

1. **Text Interpretation:** Extract meaning from text
2. **Syntax Analysis:** Analyze grammatical structure
3. **Semantic Analysis:** Understand meaning of words/phrases
4. **Entity Recognition:** Identify important entities (names, dates, locations)
5. **Context Understanding:** Grasp meaning from surrounding context

B. NLG (Natural Language Generation)

Focus: Enabling machines to produce human-like language

Key Functions:

1. **Content Planning:** What to say
2. **Sentence Planning:** How to structure sentences
3. **Lexical Choice:** Which words to use
4. **Surface Realization:** Generate actual text
5. **Response Generation:** Create output text

V. NLP Processing Phases

Phase 1: Lexical & Morphological Analysis

Purpose: Break text into meaningful units and assign linguistic properties.

Operations:

1. **Tokenization:** Split text into tokens (words, punctuation)
 - o "I'm learning NLP." → ["I", "m", "learning", "NLP", "."]
2. **POS Tagging:** Assign part-of-speech tags
 - o "The cat sat" → [Article, Noun, Verb]
 - o Identifies: nouns, verbs, adjectives, adverbs, prepositions, etc.
3. **Lemmatization:** Convert to dictionary base form
 - o "running", "runs", "ran" → "run"
 - o "better" → "good"
 - o Uses vocabulary and morphological analysis

Morphological Relations:

- **Inflectional:** Change form without changing part-of-speech
 - o cat/cats, run/running, happy/happier
- **Derivational:** Create new words with different meaning
 - o run → runner, joy → joyful

Phase 2: Syntactic Analysis

Purpose: Understand sentence structure and grammar.

Operations:

1. **Syntax Analysis:** Parse sentence structure
 - o Verify grammatically correct sentences
 - o Identify main clause, dependent clauses, phrases
2. **Parsing:** Create parse tree showing grammatical relationships
Sentence
/
Noun Verb Phrase
(I) /
Verb Object
(am) (learning)

Phase 3: Semantic Analysis

Purpose: Extract meaning from text.

Operations:

1. **Word Sense Disambiguation:** Resolve ambiguous words
 - o "bank" → financial institution vs. river bank
 - o Needs context
2. **Semantic Role Labeling:** Identify relationships
 - o Who did what to whom
 - o Extract predicate-argument structure

Phase 4: Discourse Integration

Purpose: Understand relationships between sentences.

Operations:

1. **Coreference Resolution:** Resolve pronouns and references
 - o "John gave Mary the book. He was happy."
 - o He → John (resolve reference)
2. **Discourse Relations:** Understand sentence connections
 - o Coherence and logical flow

Phase 5: Pragmatic Analysis

Purpose: Understand intended meaning beyond literal words.

Operations:

1. **Use World Knowledge:** Apply common sense
 - o "John is cold." → John needs warmth (implied)
2. **Context Understanding:** Consider situation and speaker intent
 - o Same words different meaning in different contexts
3. **Speech Act Recognition:** What is the speaker trying to do
 - o Statement, question, command, request

VI. NLP Preprocessing

Goal: Clean and prepare text for ML models.

Preprocessing Steps

1. Lowercasing

- Convert all to lowercase
- "Hello" and "hello" treated identically

2. Tokenization

- Split text into tokens
- "Hello, world!" → ["Hello", ",", "world", "!"]

3. Stopword Removal

- Remove common words (the, a, is, etc.)
- Reduces noise, speeds up processing
- Example: "The cat is on the mat" → ["cat", "mat"]

4. Stemming

- Reduce word to root form using simple rules
- "running", "runs" → "run"
- Fast but crude (can produce non-words)
- Algorithm: Snowball stemmer

5. Lemmatization

- Convert to dictionary base form using vocabulary
- "running" → "run", "better" → "good"
- More accurate but slower than stemming

6. Punctuation & Special Characters Removal

- Remove or normalize: !, ?, @, #, \$, %
- Can be important in some tasks (keep if needed)

7. Normalization

- Convert to standard format
- "Normalize": Fix typos, handle contractions
- "don't" → "do not"

8. Sentence Segmentation

- Split into sentences
- Challenge: Handle abbreviations (Dr, etc.)

Generic NLP Pipeline

Raw Text

↓

Data Acquisition

↓

Preprocessing:

- Tokenization
- Stopword removal
- Lemmatization/Stemming
- Normalization

↓

Feature Engineering:

- Text representation (embeddings)
 - ↓
 - Modeling:
- Train ML/DL model
 - ↓
 - Training:
- Learn from data
 - ↓
 - Evaluation:
- Test on holdout set
 - ↓
 - Improvement:
- Feature engineering
- Model tuning
- Ensemble methods
 - ↓
 - Deployment:
- Production system
 - ↓
 - Model Update:
- Retrain with new data

VII. Text Representation Techniques

Goal: Convert text into numerical form that ML/DL models can process.

Core Principle: Machines need numbers, so convert words → vectors.

1. One-Hot Encoding

Process: Each unique word represented as binary vector with 1 at its index, 0 elsewhere.

Example (vocabulary = {cat, dog, mat}):

cat: [1, 0, 0]

dog: [0, 1, 0]

mat: [0, 0, 1]

Characteristics:

- ✓ Simple, easy to implement
- ✓ All vectors same length
- ✗ No information about word meaning
- ✗ Vectors become very high-dimensional for large vocabulary
- ✗ Sparse (mostly zeros)
- ✗ Cannot capture similarity between words

Dimensionality: Vocabulary size (e.g., 10,000 dimensions for 10K words)

2. Bag of Words (BoW)

Process: Count word frequencies in document, create vector of counts.

Example (vocabulary = {cat, dog, mat}):

Document: "cat sat on mat"

Counts: cat=1, dog=0, mat=1, sat=1

Representation: [1, 0, 1, 1]

Characteristics:

- ✓ Simple, effective for many tasks
- ✓ Works for text classification
- ✓ Good baseline model
- ✗ Loses word order (ignores grammar)
- ✗ Creates sparse vectors
- ✗ Large vector for large vocabulary
- ✗ Treats all words equally (doesn't distinguish common vs. rare)

Use Cases:

- Text classification (spam detection, sentiment)
- Document similarity
- Information retrieval

3. TF-IDF (Term Frequency - Inverse Document Frequency)

Motivation: BoW treats all words equally, but rare words are more informative than common ones.

Formula:

$$\text{TF-IDF}(\text{term}, \text{doc}) = \text{TF}(\text{term}, \text{doc}) \times \text{IDF}(\text{term})$$

TF (Term Frequency):

- How often term appears in document
- TF = count of term / total words in document
- Range: [0, 1]

IDF (Inverse Document Frequency):

- How rare the term is across all documents
- IDF = $\log(\text{total docs} / \text{docs containing term})$
- Rare terms get higher weight
- Common terms (the, a) get low weight

Example:

Document: "machine learning is important"

- "machine": rare in corpus → high IDF weight
- "is": common in corpus → low IDF weight
- TF-IDF gives high weight to "machine", low to "is"

Characteristics:

- ✓ Improves over BoW by weighting importance
- ✓ Good for search engines
- ✓ Works better for small datasets
- ✗ Still loses word order
- ✗ Produces sparse vectors
- ✗ Doesn't capture semantic meaning
- ✗ Word context independent

Use Cases:

- Search engines (ranking relevant documents)
- Keyword extraction
- Document similarity

VIII. Word Embeddings

Definition: Dense numerical representations of words in lower-dimensional space, capturing semantic meaning and relationships.

Key Idea: Similar words have similar vectors.

Dimensions: Typically 100-300 (vs. 10,000+ for one-hot).

Why Word Embeddings?

1. **Dimensionality Reduction**
 - From 10,000+ dimensions → 100-300 dimensions
 - Reduces computation, memory, and overfitting
2. **Preserve Meaning and Relationships**
 - Similar words have similar vectors
 - "king" - "man" + "woman" ≈ "queen" (vector arithmetic!)
3. **Transfer Learning**
 - Pretrained embeddings work across tasks
 - Save training time and data
4. **Semantic Information**
 - Captures: synonymy, analogy, context

Frequency-Based Embeddings

Vector Space Model (VSM)

Represents words/documents as vectors in high-dimensional space

Distance between vectors → semantic similarity

Uses techniques like TF-IDF, Co-occurrence

Characteristics:

- Sparse, high-dimensional
- Doesn't capture deep meaning
- Simple and fast

Co-occurrence Matrix

Builds matrix counting word co-occurrences in fixed context window

Example (window size=2):

Document: "The cat sat on the mat"

Co-occurrence pairs:

- "cat" appears with "the", "sat"
- "sat" appears with "cat", "on"

Results:

- Words appearing together have relationship
- Can extract semantic relationships

Disadvantages:

- Huge matrix ($V \times V$, where V =vocabulary size)
- Very sparse (mostly zeros)
- High memory usage

Prediction-Based Embeddings (Neural Embeddings)

Word2Vec (Mikolov, 2013)

Idea: Use neural network to learn embeddings by predicting words.

Two Approaches:

A. CBOW (Continuous Bag of Words)

Predict center word from context words

Example (window=2):

Input: [The, sat, on, mat] (context)

Predict: [cat] (center word)

Architecture:

[Context words] → Average embedding → Hidden layer → Softmax → [Predicted word]

Advantages:

- Faster training
- Works well for frequent words
- Smooths out noise

Uses:

- When context is most informative

B. Skip-Gram

Opposite of CBOW: Predict context from center word

Example:

Input: [cat] (center word)

Predict: [The, sat, on, mat] (context)

Architecture:

[Word] → Hidden layer → Softmax → [Context words]

Advantages:

- Better for rare words
- Generates more training pairs

Uses:

- When center word is informative
- Learning from limited data

Training Tricks:

- **Hierarchical Softmax**: Faster prediction for large vocabulary
- **Negative Sampling**: Approximate softmax, much faster

Word2Vec Characteristics:

- ✓ Captures semantic relationships ("king-man+woman≈queen")
- ✓ Fast training
- ✓ Dense, low-dimensional vectors
- ✓ Supports vector arithmetic
- ✓ Trained on massive corpora
- ✗ Single embedding per word (ignores polysemy)

GloVe (Global Vectors for Word Representation)

Innovation: Combines global co-occurrence statistics with local context windows.

Approach:

- Learn word vectors by minimizing weighted least squares
- Uses co-occurrence probabilities
- Captures global statistics better than Word2Vec

Advantages:

- Captures both global and local information
- Better performance on analogy tasks
- Good for downstream NLP tasks

FastText (Facebook AI)

Innovation: Represents words as bag of character n-grams.

Example (n-grams for "where", n=3):

"where" → ["<wh", "whe", "her", "ere", "re>"]

Advantages:

- ✓ Better for morphologically rich languages
- ✓ Captures morphological information (prefixes, suffixes)

- ✓ Handles out-of-vocabulary (OOV) words
 - "learning" seen during training, "learner" not seen
 - Can generate representation for "learner" from its n-grams
- ✓ Better for languages with complex word formation

When to Use:

- Morphologically complex languages
- Limited training data
- Need to handle OOV words

Limitation of Static Embeddings

Problem: Single embedding per word, ignores context.

Example:

- "bank" (financial): [0.5, 0.2, 0.1, ...]
- "bank" (river): Same embedding! (Wrong!)
- Context "money bank" vs "river bank" requires different embeddings

Solution: Contextual Embeddings

- ELMO: Bidirectional LSTMs generate context-dependent embeddings
 - BERT: Bidirectional transformer pretraining
 - GPT: Autoregressive transformer
 - Each occurrence of word gets unique embedding based on context
-

IX. Text Generation

Definition: Using AI/NLP to generate new human-like written content.

How It Works

Training Phase:

1. Model reads millions of text samples
2. Learns: grammar, vocabulary, sentence flow, writing style, patterns
3. Learns associations: which words follow which words

Generation Phase:

1. **Tokenization:** Convert input text to tokens
2. **Embedding:** Convert tokens to embeddings
3. **Processing:** LSTM/Transformer processes input
4. **Prediction:** Model predicts most likely next token
5. **Sampling:** Choose next token (greedy, beam search, temperature)
6. **Repeat:** Use predicted token as input for next prediction
7. **Continue** until EOS (end-of-sequence) token or max length

Generation Strategies

1. Greedy Decoding

At each step: Choose token with highest probability

Advantage: Fast

Disadvantage: May get stuck in local optima, often suboptimal overall

2. Beam Search

Maintain k best hypotheses at each step

Explore multiple paths, keep most promising

Process:

1. Start with k top tokens
2. For each token, generate k best continuations
3. Keep overall k best combinations
4. Repeat until EOS

k=3 (beam size=3): Keep 3 best paths at each step

3. Temperature Sampling

Control randomness of predictions

Temperature τ :

- $\tau < 1$: More confident, less random
- $\tau = 1$: Default softmax probabilities
- $\tau > 1$: More random, more diverse

Higher temperature \rightarrow more creative but less coherent

Lower temperature \rightarrow more predictable but safer

Text Generation Approaches

1. Autoregressive Models (Like GPT)

Generate next token using previous tokens

$$P(y_t | y_1, y_2, \dots, y_{(t-1)})$$

One direction: left to right (past \rightarrow future)

Advantage: Natural sequential generation

Used in: Language models, chatbots, code generation

2. Seq2Seq Models (Encoder-Decoder)

Encode input sequence \rightarrow Generate output sequence

For generation task:

- Encoder: Summarize input context
- Decoder: Generate output token by token

Used in: Machine translation, summarization, image captioning

3. Fine-Tuned Models

Start with pretrained model

Adapt to specific task/domain through fine-tuning

Examples:

- Legal document generation (pretrain on general, finetune on legal)
- Medical report generation
- Financial summarization

Applications of Text Generation

Application	Use Case
Content Creation	Blog posts, product descriptions, social media
Chatbots	Customer service, virtual assistants
Code Generation	GitHub Copilot (autocomplete code)
Translation	Machine translation systems
Summarization	Automatic document/article summaries
Dialog	Conversational systems, Q&A

Limitations of Text Generation

- **Coherence:** Long sequences may lose coherence
- **Factuality:** Can generate plausible but incorrect information (hallucinations)
- **Bias:** Inherits biases from training data
- **Repetition:** May repeat same phrases
- **Context Window:** Limited memory of context

MODULE 5: MACHINE LEARNING CASE STUDIES

Overview

This module covers five important ML applications demonstrating real-world problem-solving.

Case Study 1: Churn Analysis and Prediction

Problem: Customer Churn

- **Definition:** Customers leaving/stopping service usage
- **Business Impact:** High customer acquisition cost, loss of revenue
- **Goal:** Identify at-risk customers early for retention efforts

Solution Approach:

1. **Data Collection:** Customer data (usage, payments, complaints)
2. **Feature Engineering:** Create features indicating churn signals
3. **Classification Model:** Predict churn probability
4. **Retention Strategy:** Target high-churn-probability customers with offers

Case Study 2: Sentiment Analysis / Topic Mining from News

Problem: Understanding News Content

- **Challenge:** Large volume of unstructured text
- **Goal:** Extract sentiment and main topics automatically

Solution:

1. **Text Preprocessing:** Clean text data
2. **Feature Extraction:** TF-IDF, word embeddings
3. **Sentiment Analysis:** Classify positive/negative/neutral
4. **Topic Modeling:** Extract main themes from documents
5. **Visualization:** Topic trends over time

Case Study 3: Customer Segmentation and Value

Problem: Understanding Customer Groups

- **Goal:** Identify different customer segments for targeted strategies
- **Application:** Marketing, pricing, service customization

Solution:

1. **Data Preparation:** Customer features (spending, frequency, loyalty)
2. **Feature Normalization:** Scale features to comparable ranges
3. **Clustering:** K-Means to group customers
4. **Analysis:** Characterize each segment
5. **Strategy:** Different approaches per segment

Case Study 4: Netflix Movie Recommendation System

Problem: Recommending Relevant Movies

- **Challenge:** Millions of movies, users can't browse all
- **Goal:** Suggest movies user will enjoy

Solution Approaches:

A. Content-Based:

- Extract movie features (genre, director, actors)
- User profile: movies they liked + features
- Recommend similar movies

B. Collaborative Filtering:

- Find users with similar preferences
- Recommend movies liked by similar users
- Effective with large user base

C. Hybrid:

- Combine content and collaborative approaches
- Use matrix factorization for scalability
- Modern systems use deep learning

Key Metrics:

- Precision: Recommendation accuracy
- Recall: Coverage of relevant items
- RMSE: Prediction error

Case Study 5: Image & Text Classification using TensorFlow/PyTorch

A. Image Classification

Problem: Classify images into categories

Solution:

1. **Data:** Collect labeled image dataset
2. **Preprocessing:** Resize, normalize, augment
3. **Model:** CNN (transfer learning often used)
4. **Training:** Backpropagation
5. **Evaluation:** Accuracy, confusion matrix
6. **Deployment:** Serve model predictions

Transfer Learning:

- Use pretrained CNN (ImageNet)
- Fine-tune last layers on your task
- Fast, requires less data

B. Text Classification

Problem: Assign category to text (spam, sentiment, topic)

Solution:

1. **Data:** Collect labeled text samples
2. **Preprocessing:** Tokenization, lemmatization
3. **Representation:** Word embeddings or TF-IDF
4. **Model:** LSTM, CNN for text, or Transformer
5. **Training:** Cross-entropy loss, backpropagation
6. **Evaluation:** Accuracy, precision, recall, F1

Modern Approaches:

- **BERT:** Pretrained bidirectional model, finetune for your task
 - **Transformers:** State-of-the-art for NLP
 - **GPT:** Generate text in addition to classification
-

Conclusion and Key Takeaways

Machine Learning:

- Foundational algorithms and techniques
- Works well for tabular/structured data
- Requires feature engineering

Deep Learning:

- Automatic feature learning
- Requires large datasets
- Excellent for unstructured data (images, text, speech)

Specialized Architectures:

- **CNNs:** For spatial data (images)
- **RNNs/LSTMs:** For sequential data (time series, text)
- **Seq2Seq:** For sequence transformation (translation)
- **Transformers:** State-of-the-art for NLP and vision

Best Practices:

1. Understand problem deeply
 2. Prepare and clean data thoroughly
 3. Start simple, increase complexity gradually
 4. Use appropriate evaluation metrics
 5. Validate with cross-validation
 6. Use regularization to prevent overfitting
 7. Document and reproduce results
 8. Consider computational resources
 9. Monitor performance in production
 10. Keep learning - field evolves rapidly
-

Appendix: Important Formulas and Notation

Key Equations

Bias-Variance Tradeoff:

Total Error = Bias² + Variance + Irreducible Error

Gradient Descent Update:

$$W_{\text{new}} = W - \alpha \cdot \nabla L(W)$$

Loss Functions:

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

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

Cross-Entropy: $L = -\sum y \log(\hat{y})$

Neural Network Forward Pass:

$$z = Wx + b$$

$a = f(z)$ [activation function]

Backpropagation Chain Rule:

$$\partial L / \partial W = \partial L / \partial a \cdot \partial a / \partial z \cdot \partial z / \partial W$$

LSTM Update:

$$C_t = f_t \odot C_{(t-1)} + i_t \odot \tilde{C}_t$$

$$h_t = o_t \odot \tanh(C_t)$$

TF-IDF:

$$\text{TF-IDF} = \text{TF(term, doc)} \times \text{IDF(term)}$$

$$\text{IDF} = \log(\text{total docs} / \text{docs containing term})$$

This comprehensive document covers all four modules with detailed explanations, examples, and practical insights. Use this as a complete study guide for your semester examinations!