

Module 1: Supervised vs Unsupervised Learning

1.1 What is Machine Learning?

Definition

Arthur Samuel's Definition: "The field of study that gives computers the ability to learn without being explicitly programmed"

Key Example

- Samuel's checkers program (1950s)
- Computer learned by playing thousands of games against itself
- Identified winning/losing positions through experience
- Became better than Samuel himself

Core Principle

More training opportunities → Better performance

Two Main Types

1. **Supervised Learning** - Most widely used in real-world applications
2. **Unsupervised Learning** - Second most common type

Course Structure

- Courses 1-2: Supervised Learning
 - Course 3: Unsupervised Learning, Recommender Systems, Reinforcement Learning
-

1.2 Supervised Learning

Definition

Learning algorithm that maps **inputs** (\mathbf{x}) to **outputs** (\mathbf{y})

Key Characteristic: Training with labeled examples (correct input-output pairs)

Applications

- **Email Spam Filtering:** Email → Spam/Not Spam
- **Speech Recognition:** Audio → Text transcript
- **Machine Translation:** English → Other languages
- **Online Advertising:** User/Ad info → Click prediction
- **Self-Driving Cars:** Image/Sensor data → Position of objects
- **Manufacturing:** Product image → Defect detection

Type 1: Regression

Predicts continuous numerical values from infinitely many possibilities

Example: Housing Price Prediction

- **Input (x):** House size (sq ft)
- **Output (y):** Price (\$)
- Can fit straight line or curve to data
- Predicts any number (e.g., \$150,000, \$183,000, \$200,000)

Mathematical representation: $y = f(x)$ where $y \in \mathbb{R}$

Type 2: Classification

Predicts discrete categories from a small, finite set of possibilities

Example: Breast Cancer Detection

- **Input (x):** Tumor size
- **Output (y):** Benign (0) or Malignant (1)
- Can have 2+ categories (e.g., Type 0, Type 1, Type 2 cancer)

Key Difference from Regression:

- Classification: Small finite set of categories
- Regression: Infinitely many possible numbers

Multiple Input Features Can use multiple inputs simultaneously:

- Tumor size AND patient age
- Algorithm finds decision boundary separating categories
- Real applications use many features (thickness, uniformity, cell size, etc.)

Mathematical representation: $y \in \{0, 1\}$ or $y \in \{0, 1, 2, \dots, n\}$

1.3 Unsupervised Learning

Definition

Given data with **only inputs (x)**, no output labels (y)

Goal: Find structure, patterns, or interesting relationships in data

Type 1: Clustering

Groups similar data points together automatically

Applications 1. Google News

- Groups related articles together
- Finds common words (e.g., “panda,” “twin,” “zoo”)
- Adapts to new topics daily without human supervision

2. DNA/Genetic Data Analysis

- Groups individuals by genetic similarities
- Identifies “types” of people based on gene expression
- No predefined categories

3. Market Segmentation

- Groups customers by behavior/motivation
 - Example: Knowledge seekers, Career developers, Industry updaters
 - Enables targeted marketing strategies
-

Type 2: Anomaly Detection

Detects unusual events or outliers

Applications

- Fraud detection in financial systems
 - Identifying unusual transactions
 - Security monitoring
-

Type 3: Dimensionality Reduction

Compresses large datasets while preserving information

- Reduces data size
 - Maintains essential patterns
 - Improves computational efficiency
-

Key Comparisons

Aspect	Supervised Learning	Unsupervised Learning
Data	Inputs (x) + Labels (y)	Only inputs (x)

Aspect	Supervised Learning	Unsupervised Learning
Goal	Predict y for new x	Find patterns/structure
Examples	Regression, Classification	Clustering, Anomaly detection
Training	Learn from “right answers”	Discover hidden structure

Summary

Supervised Learning

- **Regression:** Predict continuous numbers (housing prices, temperature)
- **Classification:** Predict categories (spam/not spam, disease diagnosis)

Unsupervised Learning

- **Clustering:** Group similar data (news articles, customer segments)
- **Anomaly Detection:** Find unusual patterns (fraud detection)
- **Dimensionality Reduction:** Compress data efficiently

Best Practice

Success requires both:

1. Understanding the algorithms (tools)
2. Knowing how to apply them effectively (practical skills)