

Artificial Intelligence

For naturally intelligent beings

CSE422 Lecture Notes

First Edition

Ipshita Bonhi Upoma



Inspiring Excellence

Department of Computer Science and Engineering
School of Data and Sciences

NOTE TO STUDENTS

Welcome to our course on Artificial Intelligence (AI)! This course is designed to explore the core strategies and foundational concepts that power intelligent systems, solving complex, real-world problems across various domains. From navigating routes to optimizing industrial processes, and from playing strategic games to making predictive models, the use of AI is pivotal in creating efficient and effective solutions.

These lecture notes have been prepared by Ipshita Bonhi Upoma (Lecturer, BRAC University) based on the book "Artificial Intelligence: A Modern Approach" by Peter Norvig and Stuart Russell. The credit for the initial typesetting of these notes goes to Arifa Alam.

The purpose of these notes is to provide additional examples using toy problems to enhance the understanding of how basic AI algorithms work and the types of problems they can solve.

This is a work in progress, and as it is the first time preparing such notes, many parts need citations and rewriting. Several graphs, algorithms and images are directly borrowed from Russell and Norvig's book. The next draft will address these citations.

As this is the first draft, please inform us of any typos, mathematical errors, or other technical issues so they can be corrected in future versions.

Lecture Plan (Central)

Our Reference book (Artificial Intelligence: A Modern Approach 3rd Ed.): <https://drive.google.com/file/d/16pK--SRAKZzkVs8NcxxYCDFfxtvE0pmZ/view?usp=sharing>

Our Reference book (Artificial Intelligence: A Modern Approach 4th Ed.):<https://drive.google.com/file/d/1SryRW2xX9IUMAUJUtmtUv6sEnJ4qdRjD/view?usp=sharing>

CSE422 OBE Outline: <https://docs.google.com/document/d/1SJFBkfkL0wHUokhqfcBRZLhNqtw6Qhjt/edit?usp=sharing&ouid=106202042710929132457&rtpof=true&sd=true>

CSE422 Tentative Lecture Plan: <https://docs.google.com/document/d/1fAsXap2Qjm-QdAgAXKb2J2bAOedit?usp=sharing&ouid=106202042710929132457&rtpof=true&sd=true>

Links to all resources and forms, deadlines for section 13, 17, 18 will be updated in this google

sheet: <https://docs.google.com/spreadsheets/d/10GerKB3PEfhvhGsyHHSCvIOJmBkGkY4uCtn0HPYJLHK/edit?usp=sharing>

Marks Distribution

- Class Task 5%
- Quiz 10% (4 questions, Best 3 will be counted)
- Assignment 5% (3 Assignments, Bonus assignments?)
- Lab 25%
- Mid 25%
- Final 30%

Class rules for Section 13, 17 and 18:

1. Classwork on the lecture will be given after the lecture and attendance will be counted based on the classworks.
2. Feel free to bring coffee/ light snacks for yourself. Make sure to not cause any distractions in the class.
3. If you are not enjoying the lecture you are free to leave the lecture, but in no way you should do anything that disturbs me or the other students.
4. If you want me to consider a leave of absence, submit an application on with valid reason. Classwork must be submitted even if leave is considered.
5. Best 3 of 4 quizzes will be counted.
6. All assignments will be counted. 30% penalty for late submission.
7. Cheating in any form will not be tolerated and will result in a 100% penalty.
8. If bonus assignments are given, the marks of bonus will be added after completion of all other assessments.
9. Bonus marks cannot be added to the Midterm and Final marks.
10. Lab marks are totally up to the Lab instructor.
11. No grace marks will be given for any grade bump. Such requests will not be considered.

CONTENTS

Note To Students	3
Contents	5
I Classical AI	11
1 Introduction: Past, present, future	13
1.1 Some of the earliest problems that were solved using Artificial Intelligence.	13
1.2 Problems we are trying to solve using these days:	15
2 Solving Problems with AI	17
2.1 Solving problems with artificial intelligence	17
2.1.1 Searching Strategies	17
2.1.2 Constraint Satisfaction Problems (CSP)	18
2.1.3 Machine Learning Techniques	18
2.2 Some keywords you will hear every now and then	19
2.3 Properties of Task Environment	19
2.4 Types of Agents	21
2.4.1 Learning Agent	22
2.5 Problem Formulation and choice of strategies	23
2.6 Steps of Problem Formulation	23
2.7 Examples of problem formulation	25
3 Informed Search Algorithms	29
3.1 Heruistic Function	30
3.1.1 Key Characteristics of Heuristics in Search Algorithms	30
3.1.2 Why do we use heuristics?	31
3.2 Heuristic functions to solve different problems	31
3.3 Greedy Best First Search- Finally an algorithm	35
3.3.1 Algorithm: Greedy Best-First Search	37
3.3.2 Key Points	37
3.4 A* Search algorithm	38
3.4.1 Algorithm: A* Search	38
3.5 Condition on heuristics for A* to be optimal:	42
3.6 How to choose a better Heuristic:	48
3.6.1 Why is a dominant heuristic for efficient?	49
4 Local Search	51

4.1	Local Search	51
4.1.1	The Need for Local Search	52
4.1.2	Examples of problems that can be solved by local search:	52
4.1.3	Some Constraint Satisfaction Problems can also be solved using local search strategies:	53
4.2	Local Search Algorithms:	54
4.2.1	State-Space and Objective Function:	54
4.3	Hill-Climbing Search:	56
4.3.1	Examples:	57
4.3.2	Key Characteristic:	58
4.3.3	Drawbacks:	59
4.3.4	Remedies to problems of Hill-Climbing Search Algorithm:	60
4.3.5	Variants of Hill Climbing:	60
4.4	Introduction to Simulated Annealing	61
4.5	How does simulated annealing navigate the solution space:	63
4.5.1	Probability of Accepting a New State	63
4.5.2	Cooling Schedule	63
4.5.3	Random Selection of Neighbors	64
4.5.4	Mathematical Convergence	64
4.6	Example Problem:	64
4.6.1	Traveling Salesman Problem (TSP)	64
4.7	Genetic algorithm	65
4.7.1	Algorithm	66
4.7.2	Explanation of the Pseudocode	66
4.7.3	Evaluate Fitness:	66
4.8	Mutation	68
4.8.1	Purpose of Mutation	68
4.8.2	How Mutation Works	68
4.8.3	Common Types of Mutation	69
4.8.4	Mutation Rate	70
4.8.5	Diversity in Genetic Algorithm	70
4.8.6	Advantages and Applications	70
4.9	Examples	71
4.9.1	8-Queen Problem	71
4.9.2	Traveling Salesman Problem (TSP)	75
4.9.3	0/1 Knapsack Problem	79
4.9.4	Graph Coloring Problem:	81
4.9.5	Max-cut Problem	83
4.10	Application of Genetic Algorithm in Machine Learning	85
4.10.1	Problem Setup: Feature Selection for Predictive Modeling	85
4.10.2	Genetic Algorithm Process for Feature Selection	85
4.10.3	Problem Setup: Hyperparameter Optimization for a Neural Network	86
4.10.4	Genetic Algorithm in AI Games:	88
4.10.5	Genetic algorithms in Finance	89

5 Lecture 6: Adversarial Search/Games	91
5.1 Introduction	91
5.1.1 Two player zero-sum games	91
5.1.2 Maximizer and Minimizer	92
5.2 Making optimal decisions using Minimax Search	95
5.3 Minimax algorithm	95
5.4 Alpha—Beta Pruning	96
5.4.1 Worst Case Scenario for Alpha-Beta Pruning	100
II Modern AI	103
1 Machine Learning Basics	105
1.1 Introduction	105
1.2 Why do we need machine learning	105
1.3 Paradigms of Machine Learning	106
1.3.1 Major Paradigms of Machine Learning	107
1.4 Supervised Learning	107
1.5 Steps of Supervised Learning	108
1.6 Types of Supervised Learning	109
1.7 Model Selection	109
1.8 Hypothesis in Supervised Learning	110
1.9 Hypothesis Selection	111
2 Probability theory	115
2.1 Probability Theory in AI	115
2.2 Basic Concepts	115
2.3 Basic Probability Rules	117
2.4 Bayes' Rule: Derivation and Examples	119
2.4.1 Derivation of Bayes' Rule	119
2.4.2 Significance of Bayes' Rule	120
2.5 Discrete Random Variables	121
2.5.1 Notation for Discrete Random Variables	121
2.5.2 Properties of Discrete Random Variables	122
2.5.3 Basic Vector Notation	122
2.5.4 Discrete Random Variable with Categorical Outcomes	123
2.6 Data to Probability	124
2.6.1 Marginal Probability	124
2.6.2 Marginal Probability Distribution	125
2.6.3 Joint Probability	125
2.6.4 Joint Probability Distribution	125
2.7 Understanding Conditional Probability using Probability Distribution Table	125
2.7.1 Absolute independence	126
2.7.2 Conditional Independence	126
2.7.3 Example: Checking Conditional Independence	128
2.8 Law of Total Probability	128

3 Naive Bayes Classification	131
3.1 Key Idea	131
3.2 Example: Why Independence Helps?	131
3.3 Steps in Naive Bayes Classification	132
3.4 Example: Probability of HIV	132
3.5 Example: Naive Bayes for Spam Detection	133
3.6 Example: Predicting Tennis Play from dataset	134
3.7 Example: Bayesian Diagnosis with another HIV Example	138
3.8 Types of Naive Bayes Classifiers	138
3.9 Advantages and Limitations	138
3.9.1 Advantages	138
3.9.2 Limitations	139
4 Decision Tree	141
4.1 Introduction	141
4.2 Entropy: A Measure of Uncertainty	142
4.3 Conditional Entropy and Information Gain	144
4.4 Example: Toy Dataset	144
4.5 The ID3 Algorithm	145
4.6 Example: Full ID3 Walkthrough	146
4.7 Overfitting and Pruning	152
4.8 Handling Special Cases	154
4.8.1 Missing Data	154
4.8.2 Continuous Attributes	154
4.8.3 Regression Trees: Handling Numerical Target Variables	155
5 Linear Regression and Gradient Descent	157
5.1 Introduction	157
5.2 Mathematical Model	157
5.3 Learning Objective	158
5.4 Examples of Linear Regression	158
5.5 Univariate Linear Regression	159
5.5.1 Multivariable Linear Regression	160
5.6 Gradient Descent Algorithm	161
5.7 Advantages of Linear Regression	167
6 Linear Classifier and Logistic Regression	169
6.1 Linear Functions with Thresholds for Classification	169
6.1.1 Perceptron Learning	170
6.1.2 Linear Classifier with Logistic Regression	172
6.1.3 Gradient Descent for Logistic Regression	173
6.1.4 Computing Gradient	173
6.1.5 Logistic regression: Computing Gradient for Mean Squared Error	174
6.1.6 Logistic Regression: Computing gradient of Binary Cross-Entropy	176
6.2 Example: Logistic Regression with Binary Cross-Entropy	177
7 Neural Networks	179
7.1 Introduction	179

7.2	Perceptron: The Basic Unit	179
7.3	Network Layers	180
7.3.1	Input Layer	180
7.3.2	Hidden Layers	181
7.3.3	Activation Function	182
7.3.4	Output Layer	184
7.4	Understanding the Weight Matrix and Activation Function	184
7.5	Feedforward Neural Networks	186
7.6	Activation and Loss Functions	189
7.6.1	Sigmoid (Binary Classification)	190
7.6.2	Softmax (Multiclass Classification)	190
7.6.3	Linear (Regression)	190
7.7	Training: Gradient Descent and Backpropagation	190
7.7.1	Backpropagation in Feedforward Neural Networks	190
7.7.2	Computing the Gradient at the Output Layer	191
7.7.3	Understanding the Derivative of Activation with Respect to Pre-Activation	195
7.7.4	Understanding the Derivative of Loss with Respect to Output Activation	196
7.7.5	Defining Error, $\frac{d\mathcal{L}}{dZ^{(l)}}$	197
7.7.6	Update Rule at the output layer	197
7.7.7	Computing gradient at the Hidden Layer ($l = 1$)	197
7.7.8	Simplified Example: Backpropagation in 2-layer Feed-forward Neural Network	201
7.7.9	Example: Backpropagation on a 3-Layer Neural Network	202

Part I

Classical AI

CHAPTER 1

INTRODUCTION: PAST, PRESENT, FUTURE

1.1 Some of the earliest problems that were solved using Artificial Intelligence.

In the 1950s, the field of Artificial Intelligence (AI) began to take shape as researchers pondered the question: Can machines think? Alan Turing, one of the pioneers in computing, proposed the Turing Test as a way to measure if a machine could think like a human. This sparked an interest in developing computers that could simulate human reasoning, laying the groundwork for AI. Researchers focused on translating human problem-solving abilities into mathematical rules and algorithms, forming the very foundation of the field.

AI research quickly pushed forward, especially in Bayesian networks and Markov Decision Processes, which allowed machines to reason under uncertainty. Games also became a testing ground for AI. In 1951, Turing designed a theoretical chess program, followed soon by Arthur Samuel's checkers program, one of the earliest examples of a machine that could learn from experience—a major step towards machine learning.

AI wasn't just about games, though. In 1956, The Logic Theorist program, created by Allen Newell, Herbert Simon, and Cliff Shaw, was designed to mimic human reasoning, even proving mathematical theorems from Principia Mathematica. This breakthrough demonstrated AI's ability to solve complex logical tasks. Meanwhile, machine translation experiments, like the Georgetown project in 1954, successfully converted Russian sentences into English, showcasing AI's potential for language processing.

Even early speech recognition made strides in this decade, with Bell Labs developing Audrey, a system that recognized spoken numbers. These advancements in the 1950s laid a strong foundation for AI's growth.

In the 1960s, AI research advanced further, focusing on problem-solving algorithms and early neural networks like Perceptrons. During this time, control theory also became part of

robotics, helping machines perform tasks efficiently and adaptively.

The 1970s saw AI face skepticism. Reports like the Lighthill Report highlighted AI's limitations, but this decade also saw graph theory grow in importance, allowing for knowledge representation in early expert systems.

The 1980s brought significant advances with expert systems that made decisions using rule-based logic, as well as the backpropagation algorithm, which greatly improved neural networks.

In the 1990s, AI reached new heights. IBM's Deep Blue made history by defeating a world chess champion, and statistical learning models like Support Vector Machines gained popularity, leading to a data-driven approach in AI.

The 2000s introduced game theory for multi-agent systems, and deep learning brought renewed interest in neural networks for complex pattern recognition.

In the 2010s, AI achieved remarkable feats with IBM's Watson and Google DeepMind's AlphaGo, showcasing AI's ability to understand language and excel in strategic games, demonstrating the depth of AI's problem-solving abilities.

Large Language Models (LLMs) like GPT and BERT are significant developments in natural language processing. Starting with early neural network foundations in the 2000s, advancements such as Word2Vec laid the groundwork for sophisticated word embeddings. The 2017 introduction of the Transformer architecture revolutionized NLP by efficiently handling long-range dependencies in text.

In 2018, Google's BERT and OpenAI's GPT used the Transformer model to understand and generate human-like text, with BERT improving context understanding through bidirectional training, and GPT enhancing generative capabilities. Recent iterations like GPT-3 and GPT-4 have scaled up in size and performance, expanding their application range from content generation to conversational AI.

Today, in the 2020s, AI is focusing on issues like fairness and bias and exploring the potential of quantum computing to revolutionize the field even further.

1.2 Problems we are trying to solve using these days:

Healthcare—

Disease Diagnosis: AI algorithms analyze medical imaging data to detect and diagnose diseases early, such as cancer or neurological disorders.

Personalized Medicine: AI helps tailor treatment plans to individual patients based on their genetic makeup and specific health profiles.

Transportation—

Autonomous Vehicles: AI powers self-driving cars, aiming to reduce human error in driving and increase road safety.

Traffic Management: AI optimizes traffic flow, reduces congestion, and enhances public transport systems through predictive analytics and real-time data processing.

Finance—

Fraud Detection: AI systems analyze transaction patterns to identify and prevent fraudulent activities in real time.

Algorithmic Trading: AI uses complex mathematical formulas to make high-speed trading decisions to maximize investment returns.

Retail—

Customer Personalization: AI enhances customer experience by providing personalized recommendations based on past purchases and browsing behaviors.

Inventory Management: AI predicts future product demand, optimizing stock levels and reducing waste.

Education—

Adaptive Learning Platforms: AI tailors educational content to the learning styles and pace of individual students, improving engagement and outcomes.

Automated Grading: AI systems grade student essays and exams, reducing workload for educators and providing timely feedback.

Environment—] **Climate Change Modeling:** AI analyzes environmental data to predict changes in climate patterns, helping in planning and mitigation strategies.

Wildlife Conservation: AI assists in monitoring and protecting wildlife through pattern recognition in animal migration and population count.

Manufacturing—] **Predictive Maintenance:** AI predicts when equipment will require maintenance, preventing unexpected breakdowns and saving costs.

Quality Control: AI automatically inspects products for defects, ensuring high quality and reducing human error.

Cybersecurity— Threat Detection: AI monitors network activities to detect and respond to security threats in real time.

Vulnerability Management: AI predicts which parts of a software system are vulnerable to attacks and suggests corrective actions.

Entertainment— Content Recommendation: AI algorithms power recommendation systems in streaming services like Netflix and Spotify to suggest movies, shows, and music based on user preferences.

Game Development: AI is used to create more realistic and intelligent non-player characters (NPCs) and to enhance gaming environments.

Legal— Document Analysis: AI helps in reviewing large volumes of legal documents to identify relevant information, reducing the time and effort required for legal research.

Case Prediction: AI analyzes past legal cases to predict outcomes and provide guidance on legal strategies.

CHAPTER 2

SOLVING PROBLEMS WITH AI

2.1 Solving problems with artificial intelligence

In this course, we explore three distinct strategic domains of artificial intelligence: Searching Strategies, Constraint Satisfaction Problems (CSP), and Machine Learning.

Solving any problem first requires abstraction/problem formulation of the problem so that the problem can be tackled using an algorithmic solution. Based on that abstraction we choose a suitable strategy to solve the problem.

Real-world AI challenges are rarely straightforward. They often need to be broken down into smaller parts, with each part solved using a different strategy. For example, in creating an autonomous vehicle, informed search may help us find a route to destination, adversarial search helps us predict other drivers' actions, while machine learning helps the vehicle understand road signs.

Thankfully in this course, we'll focus on learning each strategy separately. This approach lets us dive deep into each area without worrying about combining them.

2.1.1 Searching Strategies

Informed Search

Why— Informed search strategies, such as A* and Best-First Search, utilize heuristics (we will come back to this later) to efficiently find solutions, focusing the search towards more promising paths. These strategies are essential in scenarios like real-time pathfinding for autonomous vehicles.

Local Search

Why— Local search methods are crucial for tackling optimization problems where finding an optimal solution might be too time-consuming. These methods, which include Simulated Annealing and Hill Climbing, are invaluable for tasks such as resource allocation and scheduling where a near-optimal solution is often sufficient.

Adversarial search

Why—Adversarial search techniques are essential for environments where agents compete against each other, such as in board games or market competitions. Understanding strategies like Minimax and Alpha-Beta Pruning allows one to predict and counter opponents' moves effectively.

2.1.2 Constraint Satisfaction Problems (CSP)

Constraint Satisfaction Problems (CSP)

Why—CSPs are studied to solve problems where the goal is to assign values to variables under strict constraints. Techniques like backtracking and constraint propagation are fundamental for solving puzzles, scheduling problems, and many configuration problems where all constraints must be satisfied simultaneously.

2.1.3 Machine Learning Techniques

Probabilistic Reasoning

Why—We delve into probabilistic reasoning to equip students with methods for making decisions under uncertainty. Techniques such as Bayesian Networks are vital for applications ranging from diagnostics to automated decision-making systems.

Decision Tree

Why—Decision trees are introduced due to their straightforward approach to solving classification and regression problems. They split data into increasingly precise subsets using simple decision rules, making them suitable for tasks from financial forecasting to clinical decision support.

Gradient Descent

Why—The gradient descent algorithm is essential for optimizing machine learning models, particularly in training deep neural networks. Its ability to minimize error functions makes it indispensable for developing applications like voice recognition systems and personalized recommendation engines.

2.2 Some keywords you will hear every now and then

Agent: In AI, an agent is an entity that perceives its environment through sensors and acts upon that environment using actuators. It operates within a framework of objectives, using its perceptions to make decisions that influence its actions.

Rational Agent: A rational agent acts to achieve the best outcome or, when there is uncertainty, the best expected outcome. It is "rational" in the sense that it maximizes its performance measure, based on its perceived data from the environment and the knowledge it has.

Autonomous Agent: An autonomous agent is a type of rational agent that can *learn from its own experiences and actions*. It can adjust its behavior based on new information, making up for any initial gaps or inaccuracies in its knowledge. Essentially, an autonomous agent operates independently and adapts effectively to changes in its environment.

Task Environment: In AI, the environment refers to everything external to the agent that it interacts with or perceives to make decisions and achieve its goals. It includes all factors, conditions, and entities that can influence or be influenced by the agent's actions.

2.3 Properties of Task Environment

Understanding these properties helps in the design and development of AI systems, tailoring the AI's architecture, algorithms, and decision-making processes to the specific nature of the environment it will operate in. This enables more effective, efficient, and appropriate responses to varying conditions and objectives.

Fully Observable vs. Partially Observable

Fully Observable: In these environments, the agent's sensors provide access to the complete state of the environment at all times. This allows the agent to make decisions with full knowledge of the world. For example, a chess game where the agent (player) can see the entire board and all the pieces at all times.

Partially Observable: Here, the agent only has partial information about the environment due to limitations in sensor capabilities or because some information is inherently hidden. Agents must often infer or guess the missing information to make decisions. For instance, in a poker game, players cannot see the cards of their opponents.

Deterministic vs. Stochastic

Deterministic: The outcome of any action by the agent is completely determined by the current state and the action taken. There is no uncertainty involved. For example, a tic-tac-toe game where each move reliably changes the board in a predictable way.

Stochastic: In these environments, actions have probabilistic outcomes, meaning the same action taken under the same conditions can lead to different results. Agents must deal with uncertainty and probability. For example, in stock trading, buying a stock does not guarantee profit due to market volatility.

Episodic vs. Sequential

Episodic: The agent's experience is divided into distinct episodes, where the action in each episode does not affect the next. Each episode consists of the agent perceiving and then performing a single action. An example is image classification tasks where each image is processed independently.

Sequential: Actions have long-term consequences, and thus the current choice affects all future decisions. Agents need to consider the overall potential future outcomes when deciding their actions. Navigation tasks where an agent (like a robot or self-driving car) must continuously make decisions based on past movements exemplify this.

Static vs. Dynamic

Static: The environment does not change while the agent is deliberating. This simplicity allows the agent time to make a decision without worrying about the environment moving on. An example is a Sudoku puzzle, where the grid waits inertly as the player strategizes.

Dynamic: The environment can change while the agent is considering its actions. Agents need to adapt quickly and consider the timing of actions. For example, in automated trading systems where market conditions can change in the midst of computations.

Discrete vs. Continuous

Discrete: Possible states, actions, and outcomes are limited to a set of distinct, clearly defined values. For example, a chess game has a finite number of possible moves and positions.

Continuous: The environment may change continuously, and the number of possible states or actions is infinite. For instance, driving a car involves navigating through a continuous range of positions and speeds.

Competitive vs. Cooperative

Competitive: Agents operate in environments where other agents might have conflicting objectives, like in strategic games such as in a game of chess where each player aims to defeat the other.

Cooperative: Agents work together towards a common goal, which may involve communication and shared tasks. For example, a collaborative robotics setting where multiple robots work together to assemble a product.

2.4 Types of Agents

In artificial intelligence, agents can be categorized based on their operational complexity and capabilities.

Simple Reflex Agents

These agents act solely based on the current perception, ignoring the rest of the perceptual history. They operate on a condition-action rule, meaning if a condition is met, an action is taken.

Example: A room light that turns on when it detects motion. It does not remember past movements; it only responds if it detects motion currently.

Model-Based Reflex Agents

These agents maintain some sort of internal state that depends on the percept history, allowing them to make decisions in partially observable environments. They use a model of the world to decide their actions.

Example: A thermostat that controls a heating system. It uses the history of temperature changes and the current temperature to predict and adjust the heating to maintain a set temperature.

Goal-Based Agents

These agents act to achieve goals. They consider future actions and evaluate them based on whether they lead to the achievement of a set goal.

Example: A navigation system in a car that plans routes not only based on the current location but also on the destination the user wants to reach.

Utility-Based Agents:

These agents aim to maximize their own perceived happiness or satisfaction, expressed as a utility function. They choose actions based on which outcome provides the greatest benefit according to this utility.

Example: An investment bot that decides to buy or sell stocks based on an algorithm designed to maximize the expected return on investment, weighing various financial indicators and market conditions.

2.4.1 Learning Agent

A learning agent in artificial intelligence is an agent that can improve its performance over time based on experience. This type of agent typically consists of four components: a learning element that updates knowledge, a performance element that makes decisions, a critic that assesses how well the agent is doing, and a problem generator that suggests challenging situations to learn from.

Example: Self-Driving Car

A self-driving car is a learning agent that adapts and improves its driving decisions based on accumulated driving data and experiences.

Performance Element: This part of the agent controls the car, making real-time driving decisions such as steering, accelerating, and braking based on current traffic conditions and sensor inputs.

Learning Element: It processes the data gathered from various sensors and feedback from the performance element to improve the decision-making algorithms. For example, it learns to recognize stop signs better or understand the nuances of merging into heavy traffic.

Critic: It evaluates the driving decisions made by the performance element. For instance, if a particular maneuver led to a near-miss, the critic would flag this as suboptimal.

Problem Generator: This might simulate challenging driving conditions that are not frequently encountered, such as slippery roads or unexpected obstacles, to prepare the car for a wider range of scenarios.

Over time, by learning from both successes and failures, a self-driving car improves its capability to drive safely and efficiently in complex traffic environments, demonstrating how learning agents adapt and enhance their performance based on experience.

2.5 Problem Formulation and choice of strategies

Problem formulation is introduced at the outset because it sets the stage for all AI strategies. It involves defining the problem in a way that a computer can process—identifying what the inputs are, what the desired outputs should be, and the constraints and environment within which the problem exists. This is crucial for effectively applying any AI technique, as a well-formulated problem can significantly simplify the solution process. It is foundational in areas such as robotics, where tasks need to be defined clearly before they can be automated.

2.6 Steps of Problem Formulation

1. Define the Goal

Start by clearly identifying what needs to be achieved. This involves understanding the desired outcome and what constitutes a solution to the problem. *Example:* For an autonomous vacuum cleaner, the goal might be to clean the entire floor space of a house without retracing any area unnecessarily. For a puzzle this could be the configuration of the puzzle when solved 2.1.

2. Identify the Initial State

Determine the starting point of the problem.

Example: In a chess game, the initial state is the standard starting position of all pieces on the chessboard.

For an engineering task, the current measurements of a system.

For a machine learning model, the initial data set from which the model will learn.

3. Determine the Possible Actions

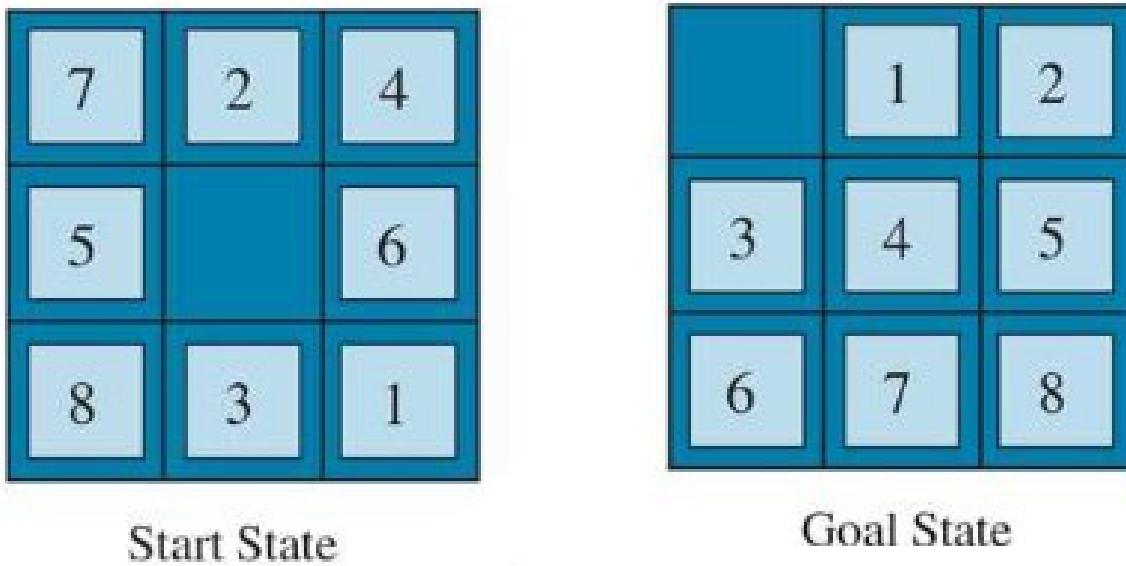


Figure 2.1: Initial and Goal State of a 8-Puzzle Game (Russel and Norvig, *Artificial Intelligence: A Modern Approach*)

List out all possible actions that can be taken from any given state.

Example: In an online booking system, actions could include selecting dates, choosing a room type, and adding guest information.

In a navigation problem, for example, these actions could be the different paths or turns one can take at an intersection.

For a sorting algorithm, actions might be the comparisons or swaps between elements.

4. Define the Transition Model

Establish how each action affects the state. The transition model describes what the new state will be after an action is taken from a current state.

Example: In a stock trading app, the transition model would define how buying or selling stocks affects the portfolio's state, including changes in cash reserves and stock quantities.

In a chess game, moving a pawn will change the state of the board.

5. Establish the Goal Test

Create a method to determine whether a given state is a goal state. This test checks if the goal has been achieved.

Example: In a puzzle like Sudoku, the goal test checks if the board is completely filled without any repeating numbers in any row, column, or grid.

In a maze-solving problem, the goal test would verify whether the current location is the exit of the maze.

6. Define the Path Cost

Decide how to measure the cost of a path. The path cost function will calculate the numerical cost of any given path from the start state to any state at any point. This is often critical in optimization problems where you want to find not just any solution, but the most cost-effective one.

Example: For a route optimization problem, the path cost could include factors like total distance, travel time, and toll costs.

7. Consider Any Constraints

Identify any constraints that must be considered. Constraints are limitations or restrictions on the possible solutions. For example, in scheduling, constraints could be the availability of resources or time slots.

Example: In university class scheduling, constraints include classroom capacities, instructor availability, and specific time blocks when certain classes can or cannot be held.

8. Select the Suitable AI Technique

Based on the problem's characteristics, such as whether the environment is deterministic or stochastic, static or dynamic, discrete or continuous, select the most appropriate AI technique. This could range from simple rule-based algorithms to complex machine learning models.

Example: For a predictive maintenance system in a factory, the suitable AI technique might involve using machine learning models like decision tree to predict equipment failures based on historical sensor data. On the other hand, in a game of chess, we use adversarial search to decide our moves by considering what the opponent might do next for certain moves.

2.7 Examples of problem formulation

1. City Traffic Navigation (Solved by Informed Search)

o Problem Formulation:

- **Goal:** To find the quickest route from a starting point (origin) to a destination (end point) while considering current traffic conditions.
- **States:** Each state represents a geographic location within the city's road network.
- **Initial State:** The specific starting location of the vehicle.
- **Actions:** From any given state (location), the actions available are the set of all possible roads that can be taken next.
- **Transition Model:** Moving from one location to another via a chosen road or intersection.

- **Goal Test:** Determines whether the current location is the destination.
 - **Path Cost:** Each step cost can be a function of travel time, which depends on factors such as distance and current traffic. The total path cost is the sum of the step costs, representing the total travel time.
- **Heuristics Used:**
- **Time Estimation:** An estimate of time from the current location to the destination, possibly using historical traffic data and real-time conditions.
 - **Distance:** Straight-line distance (Euclidean or Manhattan distance) to the goal, which helps prioritize closer locations during the search process.

2. Power Plant Operation (Solved by Local Search)

- **Problem Formulation:**
- **Goal:** To optimize the power output while minimizing fuel usage and adhering to safety regulations.
 - **States:** Each state represents a specific configuration of the power plant's operational settings (e.g., temperature, pressure levels, valve positions).
 - **Initial State:** The current operational settings of the plant.
 - **Actions:** Adjustments to the operational settings such as increasing or decreasing temperature, adjusting pressure, and changing the mix of fuel used.
 - **Transition Model:** Changing from one set of operational settings to another.
 - **Goal Test:** A set of operational conditions that meet all efficiency, safety, and regulatory requirements.
 - **Path Cost:** Typically involves costs related to fuel consumption, wear and tear on equipment, and potential safety risks. The cost function aims to minimize these while maximizing output efficiency.
- **Heuristic Used:**
- **Efficiency Metrics:** Estimations of how changes in operational settings will affect output efficiency and resource usage. This might include predictive models based on past performance data.

3. University Class Scheduling (Solved by CSP)

- **Problem Formulation:**
- **Goal:** To assign time slots and rooms to university classes in a way that no two classes that share students or instructors overlap, and all other constraints are satisfied.
 - **States:** Each state represents an assignment of classes to time slots and rooms.
 - **Initial State:** No courses are assigned to any time slots or rooms.

- **Actions:** Assign a class to a specific time slot in a specific room.
 - **Transition Model:** Changing from one assignment to another by placing a class into an available slot and room.
 - **Goal Test:** All classes are assigned to time slots and rooms without any conflicts with other classes.
 - **Path Cost:** Path cost is not typically a factor in CSP for scheduling; instead, the focus is on fulfilling all constraints.
- **Constraints:**
 - **Room Capacity:** Each class must be assigned to a room that can accommodate all enrolled students.
 - **Time Conflicts:** No instructor or student can be required to be in more than one place at the same time.
 - **Resource Availability:** Some classes require specific resources (e.g., laboratories or audio-visual equipment).
 - **Instructor Preferences:** Some instructors may have restrictions on when they can teach.

4. Disease Diagnosis (Solved by Decision Trees)

- **Problem Formulation:**
 - **Goal:** To accurately diagnose diseases based on symptoms, patient history, and test results.
 - **States:** Each state represents a set of features associated with a patient, including symptoms presented, medical history, demographic data, and results from various medical tests.
 - **Initial State:** The initial information gathered about the patient, which includes all initial symptoms and available medical history.
 - **Actions:** Actions are not typically modeled in decision trees as they are used for classification rather than processes involving sequential decisions.
 - **Transition Model:** Not applicable for decision trees since the process does not involve moving between states.
 - **Goal Test:** The diagnosis output by the decision tree, determining the most likely disease or condition based on the input features.
 - **Path Cost:** In decision trees, the cost is not typically measured in terms of path, but accuracy, specificity, and sensitivity of the diagnosis can be considered as metrics for evaluating performance.
- **Features Used:**
 - **Symptoms:** Patient-reported symptoms and observable signs.
 - **Test Results:** Quantitative data from blood tests, imaging tests, etc.

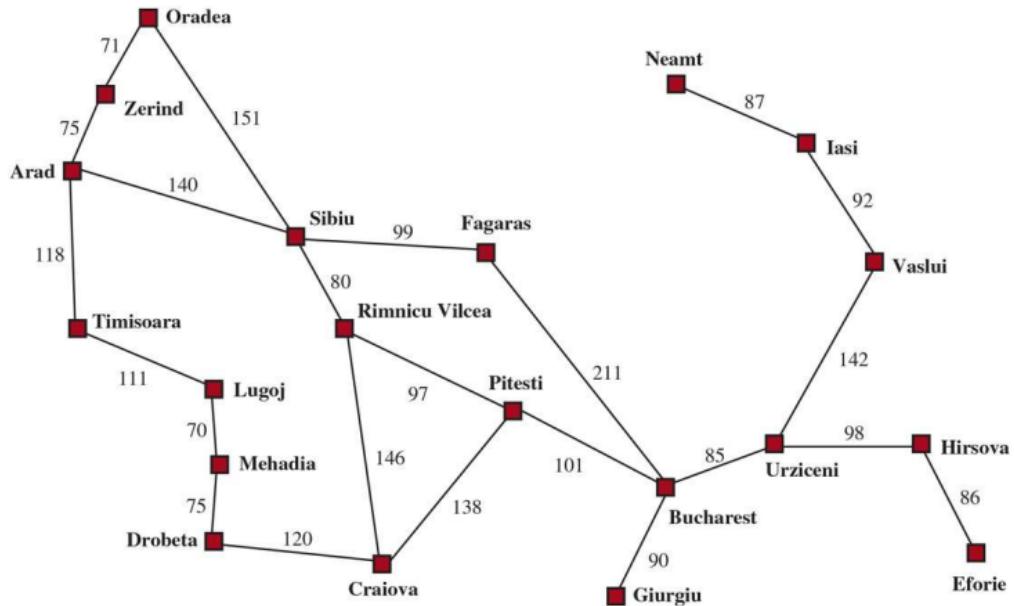
- **Demographic Data:** Age, sex, genetic information, lifestyle factors.
- **Medical History:** Previous diagnoses, treatments, family medical history.

CHAPTER 3

INFORMED SEARCH ALGORITHMS

Note: This lecture closely follows Chapter 3.6 (Heuristic Functions), 3.5 (Informed Search) to 3.5.1 (Greedy Best First Search), 3.5.2 to 3.5.4 (A* Search), 3.6.1 (Effect of Heuristic on accuracy and performance) of Russel and Norvig, *Artificial Intelligence: A Modern Approach*. The images are also borrowed from these chapters

As computer science students, you are already familiar with various search algorithms such as Breadth-First Search, Depth-First Search, and Dijkstra's/Best-First Search. These strategies fall under the category of Uninformed Search or Blind Search, which means they rely solely on the information provided in the problem definition.



A simplified road map of part of Romania, with road distances in miles.

Arad	366	Mehadia	241
Bucharest	0	Neamt	234
Craiova	160	Oradea	380
Drobeta	242	Pitesti	100
Eforie	161	Rimnicu Vilcea	193
Fagaras	176	Sibiu	253
Giurgiu	77	Timisoara	329
Hirsova	151	Urziceni	80
Iasi	226	Vaslui	199
Lugoj	244	Zerind	374

Values of h_{SLD} —straight-line distances to Bucharest.

For example, consider a map of Romania where we want to travel from Arad to Bucharest. The map indicates that Arad is connected to Zerind by 75 miles, Sibiu by 140 miles, and Timișoara by 118 miles. Using a blind search strategy, the next action from Arad would be chosen based solely on the distances to these connected cities. This approach can be slower and less efficient as it may explore paths that are irrelevant to reaching the goal efficiently.

In this course, we will focus on informed search strategies, also known as heuristic search. Informed Search uses additional information—referred to as heuristics—to make educated guesses about the most promising direction to pursue in the search space. This approach often results in faster and more efficient solutions because it avoids wasting time on less likely paths. We will study Greedy Best-First Search and A* search extensively. But first, let's explore the concept of heuristics.

3.1 Heruistic Function

In the context of informed search algorithms, a heuristic is a technique that helps the algorithm estimate the cost (often the shortest path or least costly path) from a current state (or node) to the goal state. It's essentially a function that provides guidance on which direction the search should take in order to find the most efficient path to the goal. This guidance allows informed search algorithms to perform more efficiently than uninformed search algorithms, which do not have knowledge of the goal state as they make their decisions.

3.1.1 Key Characteristics of Heuristics in Search Algorithms

Estimation: A heuristic function estimates the cost to reach the goal from a current node. This estimate does not need to be exact but should never overestimate.

Returning to the example of traveling to Bucharest from Arad: A heuristic function can estimate the shortest distance from any city in Romania to the goal. For instance, we might use the straight-line distance as a measure of the heuristic value for a city. The straight-line distance from Arad to Bucharest is 366 miles, although the optimal path from Arad to

Bucharest actually spans 418 miles. Therefore, the heuristic value for Arad is 366 miles. For each node (in this problem, city) in the state space (in this problem, the map of Romania) the heuristic value will be their straight line distance from the goal state (in this case Bucharest).

3.1.2 Why do we use heuristics?

Guidance: The heuristic guides the search process, helping the algorithm prioritize which nodes to explore next based on which seem most promising—i.e., likely to lead to the goal with the least cost.

Efficiency: By providing a way to estimate the distance to the goal, heuristics can significantly speed up the search process, as they allow the algorithm to focus on more promising paths and potentially disregard paths that are unlikely to be efficient.

3.2 Heuristic functions to solve different problems

Generating a heuristic function for use in informed search algorithms involves a process where the function must effectively estimate the cost, usually the shortest path, from any node or state in the search space to the goal. The design of the heuristic function is based on the problem domain, which requires an understanding of the rules, constraints, and ultimate goal of the problem. Here are some examples of heuristic functions for different types of problems.

8-Puzzle Game

The number of misplaced tiles (blank not included): For the figure above, all eight tiles are out of position, so the start state has $h_1 = 8$.

The sum of the distances of the tiles from their goal positions: Because tiles cannot move along diagonals, the distance is the sum of the horizontal and vertical distances- sometimes called the city-block distance or Manhattan distance.

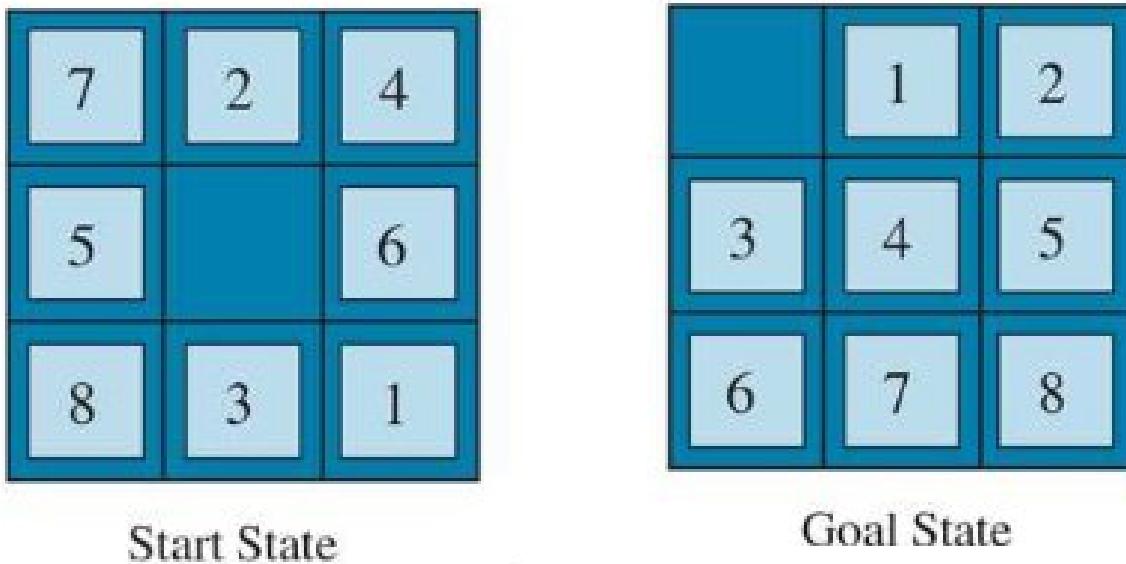


Figure 3.1: Initial and Goal State of a 8-Puzzle Game (Russel and Norvig, *Artificial Intelligence: A Modern Approach*)

Pathfinding in Maps and GPS Systems

Straight-line Distance: $d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$ where (x_1, y_1) and (x_2, y_2) are the coordinates of the current position and the goal.

Travel Time: Estimating the time needed to reach the goal based on average speeds and road types. $t = \frac{d}{v}$ where d is the distance and v is the average speed.

Traffic Patterns: Using historical or real-time traffic data to estimate the fastest route. Could involve a weighting factor, w based on traffic data, modifying the travel time: $t_{adjusted} = t \times w$.

Game AI (e.g., Chess, Go)

Material Count: Sum of the values of all pieces. For example, in chess, pawns = 1, knights/bishops = 3, rooks = 5, queen = 9.

Positional Advantage: A score based on piece positions. E.g., control of center squares in chess might be given additional points.

Mobility: Number of legal moves available M for a player at a given turn.

Web Search Engines

Keyword Frequency: The number of times a search term appears on a web-page.

$$F = \frac{\text{Number of occurrences of keyword}}{\text{Total number of words in document}}$$

Page Rank: Evaluating the number and quality of inbound links to estimate the page's importance.

Domain Authority: The reputation and reliability of the website hosting the information. Often a proprietary metric, but generally a combination of factors like link profile, site age, traffic, etc.

Robotics and Path Planning

Distance to Goal: Estimating the remaining distance to the target location. Same as straight-line distance in GPS systems.

Obstacle Proximity: Distance to the nearest obstacle to avoid collisions. $O = \min(\text{distance to each obstacle})$.

Energy Efficiency: Estimating the most energy-efficient path, important for battery-powered robots.

$$E = \sum \text{energy per unit distance} \times \text{path distance}$$

Natural Language Processing (NLP)

Word Probability: $P(w \mid \text{context})$ where w is the word and context represents the surrounding words.

Semantic Similarity: How closely words or phrases match in meaning. Often uses cosine similarity,

$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\|\|B\|}$$

where (A) , and (B) are vector representations of words or sentences.

Language Consistency: Ensuring the text follows grammatical and syntactical norms of the target language. Can be quantified using perplexity in language models.

Recommendation Systems User Behavior Tracking: Score items based on frequency and recency of user interactions. Analyzing past purchases or viewing habits to predict future interests.

Item Similarity: Recommending products similar to those a user has liked or purchased. Cosine similarity or other distance measures between item feature vectors.

Collaborative Filtering: Using preferences of similar users to recommend items. Matrix factorization techniques or neighbor-based algorithms to predict user preferences.

3.3 Greedy Best First Search- Finally an algorithm

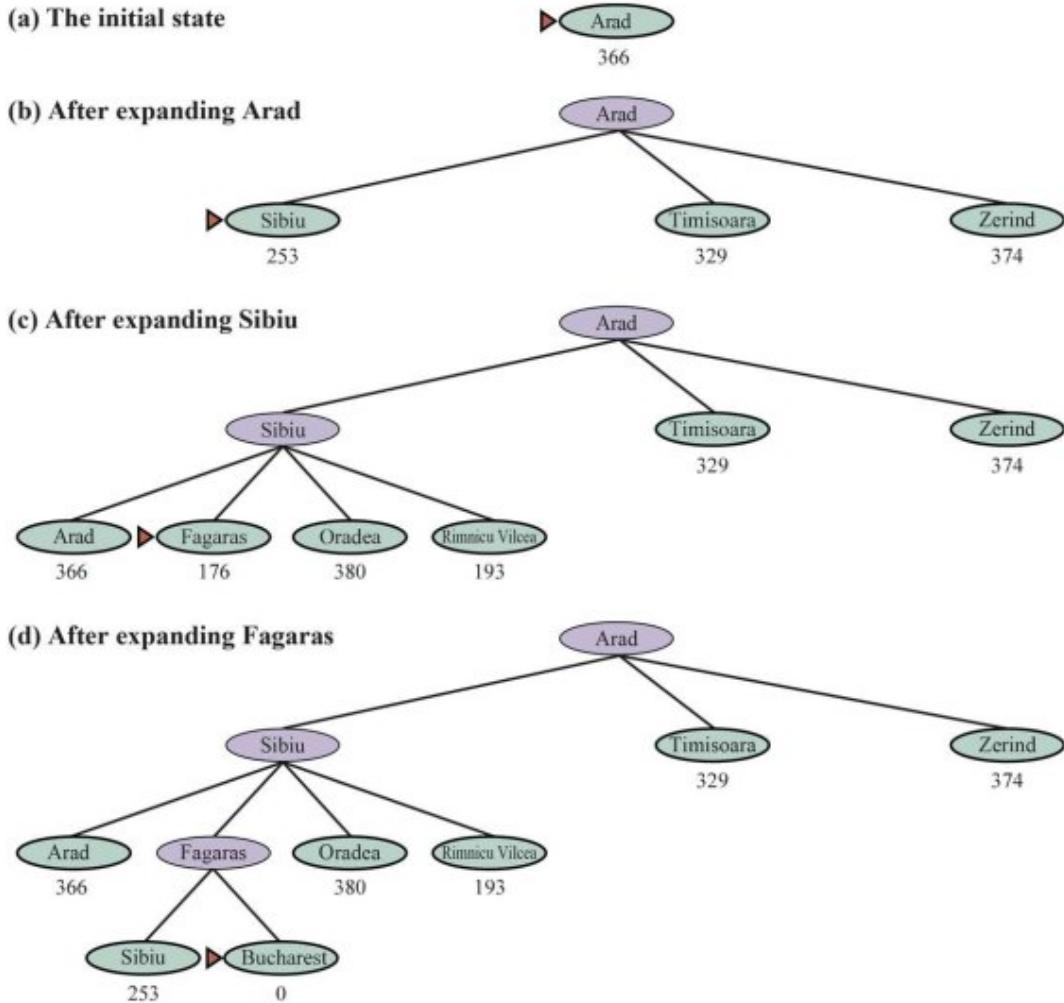
Greedy best-first search is a form of best-first search that expands first the node with the lowest $h(n)$ value—the node that appears to be closest to the goal—on the grounds that this is likely to lead to a solution quickly. So the evaluation function $f(n) = h(n)$.

Let us see how this works for route-finding problems in Romania; we use the **straight-line-distance** heuristic, which we will call h_{SLD} . If the goal is Bucharest, we need to know the straight-line distances to Bucharest, which are shown in the figure below. For example, $h_{SLD}(Arad) = 366$. Notice that the values of h_{SLD} cannot be computed from the problem description itself (that is, the ACTIONS and RESULT functions). Moreover, it takes a certain amount of world knowledge to know that h_{SLD} is correlated with actual road distances and is, therefore, a useful heuristic.

The next figure shows the progress of a greedy best-first search using h_{SLD} to find a path from Arad to Bucharest. The first node to be expanded from Arad will be Sibiu because the heuristic says it is closer to Bucharest than is either Zerind or Timisoara. The next node to be expanded will be Fagaras because it is now closest according to the heuristic. Fagaras in turn generates Bucharest, which is the goal. For this particular problem, greedy best-first search using h_{SLD} finds a solution without ever expanding a node that is not on the solution path. The solution it found does not have optimal cost, however: the path via Sibiu and Fagaras to Bucharest is 32 miles longer than the path through Rimnicu Vilcea and Pitesti. This is why the algorithm is called “greedy”—on each iteration it tries to get as close to a goal as it can, but greediness can lead to worse results than being careful.

Arad	366	Mehadia	241
Bucharest	0	Neamt	234
Craiova	160	Oradea	380
Drobeta	242	Pitesti	100
Eforie	161	Rimnicu Vilcea	193
Fagaras	176	Sibiu	253
Giurgiu	77	Timisoara	329
Hirsova	151	Urziceni	80
Iasi	226	Vaslui	199
Lugoj	244	Zerind	374

Values of h_{SLD} —straight-line distances to Bucharest.



Stages in a greedy best-first tree-like search for Bucharest with the straight-line distance heuristic h_{SLD} . Nodes are labeled with their h -values.

Greedy best-first graph search is complete in finite state spaces, but not in infinite ones. The worst-case time and space complexity is $O(|V|)$. With a good heuristic function, however, the complexity can be reduced substantially, on certain problems reaching $O(bm)$.

3.3.1 Algorithm: Greedy Best-First Search

Algorithm 1 Greedy Best First Search Algorithm

- **Input:**
 - **start:** The target node of the search
 - **goal:** The target node to reach
 - **heuristic(node):** A function that estimates the cost from node to the goal
- **Output:**
 - The path from **start** to **goal** if one exists, otherwise **None**.
- **Procedure**
 - Initialize:
 - Create a priority queue and insert the start node along with its heuristic value **heuristic(start)**.
 - Define a **visited** set to keep track of all visited nodes to avoid cycles and redundant paths.
 - Search:
 - While the priority queue is not empty:
 - * Remove the node **current** with the lowest heuristic value from the priority queue.
 - * If **current** is the goal, return the path that led to **current**.
 - * Add **current** to the **visited** set.
 - * For each neighbor **n** of **current**:
 - If **n** is not in **visited**:
 - Calculate the heuristic value **heuristic(n)**.
 - Add **n** to the priority queue with the priority set to **heuristic(n)**.
 - Failure to find the goal:
 - If the priority queue is exhausted without finding the **goal**, return **None**.

3.3.2 Key Points

Heuristic Function: This function is crucial as it determines the search behavior. A good heuristic can dramatically increase the efficiency of the search.

Completeness and Optimality: Greedy Best-First Search does not guarantee that the shortest path will be found, making it neither complete nor optimal. It can get stuck in loops or dead ends if not careful with the management of the visited set.

Data Structures: The algorithm typically uses a priority queue for the frontier and a set for the visited nodes. This setup helps in efficiently managing the nodes during the search process.

Greedy Best-First Search is particularly useful when the path's exact length is less important than quickly finding a path that is reasonably close to the shortest possible. **It is well-suited for problems where a good heuristic is available.**

3.4 A* Search algorithm

The most common informed search algorithm is A* Search (pronounced "A-star search"), a best-first search that uses the evaluation function $f(n) = g(n) + h(n)$ where $g(n)$ is the path cost from the initial state to node n , and $h(n)$ is the estimated cost of the shortest path from n to a goal state, so we have $f(n)$ estimated cost of the best path that continues from n to a goal.

3.4.1 Algorithm: A* Search

Algorithm 2 A* Search Algorithm

- **Input:**

- `start`: The starting node of the search.
- `goal`: The target node to reach.
- `neighbors(node)`: A function that returns the neighbors of `node`.
- `cost(current, neighbor)`: A function that returns the cost of moving from `current` to `neighbor`.
- `heuristic(node)`: A function that estimates the cost from `node` to the `goal`.

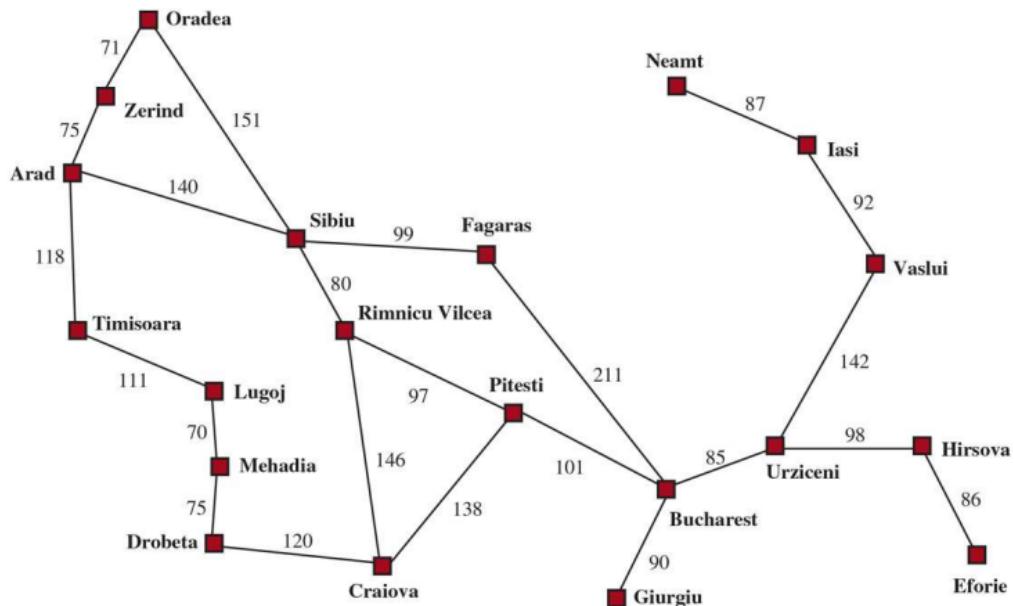
- **Output:**

- The path from `start` to `goal` if one exists, otherwise `None`.

- **Procedure**

- Initialize
 - Create a priority queue and insert the `start` node.

- Set `gScore[start]` to 0 (cost from start to `start`).
 - Set `fScore[start] = 0 + heuristic(start)` (total estimated cost from start to goal through `start`).
 - Define `cameFrom` to store the path reconstruction data.
- Search:
- While the priority queue is not empty:
 - * Remove the node `current` with the lowest `fScore` value from the priority queue.
 - * If `current` is the goal, reconstruct and return the path from start to goal using `cameFrom`.
 - * For each neighbor of `current`:
 - * For each neighbor of `current`:
 - Calculate `tentative-gScore` as `gScore[current] + cost(current, neighbor)`.
 - If `tentative-gScore` is less than `gScore[neighbor]` (or neighbor is not in priority queue):
 - Update `cameFrom[neighbor]` to `current`.
 - Update `gScore[neighbor]` to `tentative-gScore`.
 - Update `fScore[neighbor]` to `tentative-gScore + heuristic(neighbor)`.
 - If `neighbor` is not in the priority queue, add it.
- Failure to find the goal:
- If the priority queue is exhausted without reaching the `goal`, return `None`.
-



A simplified road map of part of Romania, with road distances in miles.

Arad	366	Mehadia	241
Bucharest	0	Neamt	234
Craiova	160	Oradea	380
Drobeta	242	Pitesti	100
Eforie	161	Rimnicu Vilcea	193
Fagaras	176	Sibiu	253
Giurgiu	77	Timisoara	329
Hirsova	151	Urziceni	80
Iasi	226	Vaslui	199
Lugoj	244	Zerind	374

Values of h_{SLD} —straight-line distances to Bucharest.

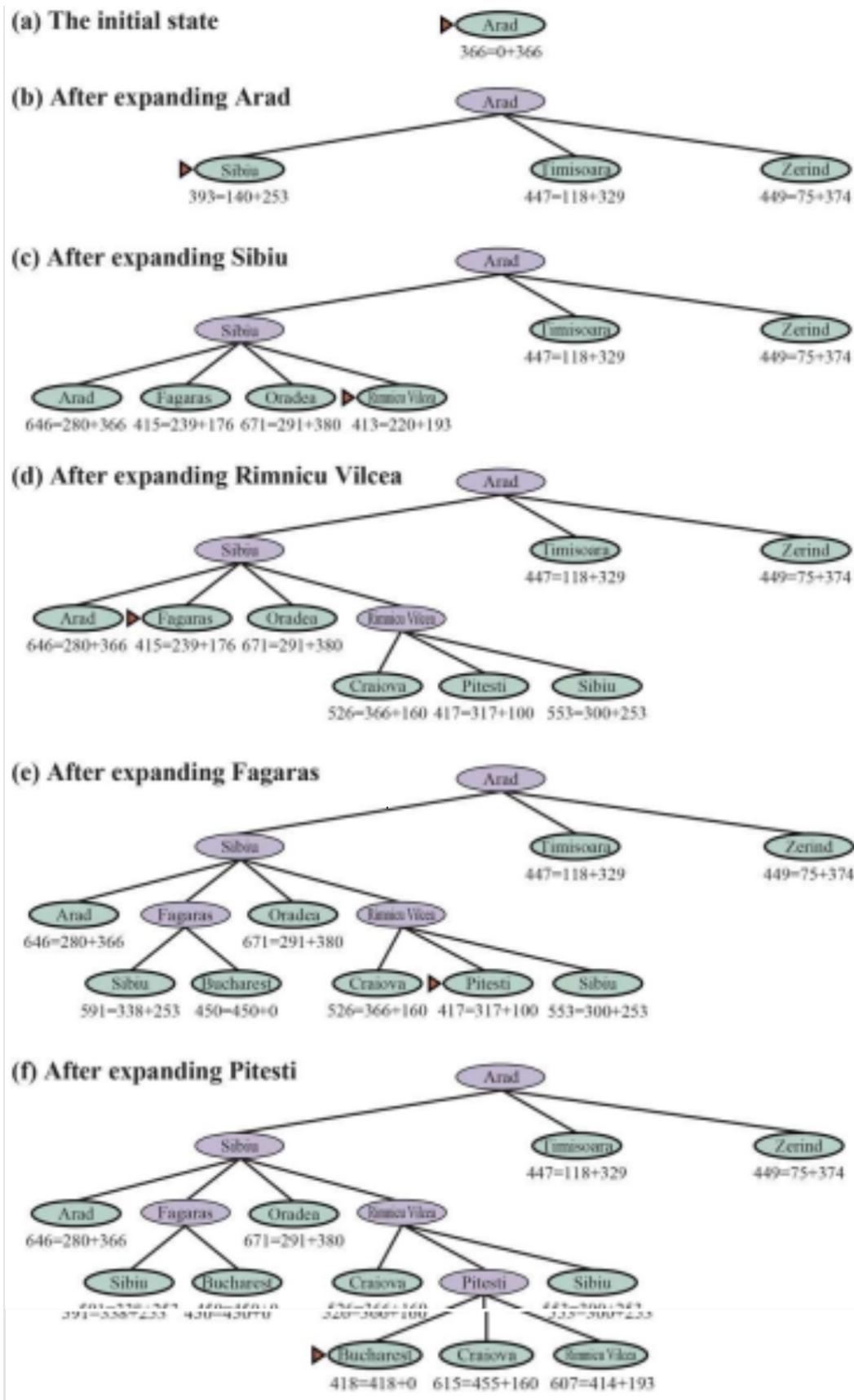


Figure 3.2: Simulation of A* Algorithm to find a path from Arad to Bucharest.
Brac University

Path Reconstruction: The path is reconstructed from the `cameFrom` map, which records where each node was reached from.

Notice that Bucharest first appears on the frontier at step (e), but isn't selected for expansion as it isn't the lowest-cost node at that moment, with a cost of 450 compared to Pitesti's lower cost of 417. The algorithm prioritizes exploring potentially cheaper routes, such as through Pitesti, before settling on higher-cost paths. By step (f), a more cost-effective path to Bucharest, costing 417, becomes available and is subsequently selected as the optimal solution.

So we can say that the A* algorithm has the following properties,

- **Cost Focus:** Prioritizes nodes with the lowest estimated total cost, $f(n) = g(n) + h(n)$
- **Guarantees Optimality:** Ensures the solution is the least expensive by expanding the cheapest node first.
- **Resource Efficiency:** Avoids exploring more expensive paths when cheaper options are available.
- **Adapts Based on New Information:** Adjusts path choices dynamically as new cost information becomes available.
- **Heuristic Importance:** Relies on the heuristic to guide the search efficiently by estimating costs.

A* Algorithm is always complete, meaning that if a solution exists then the algorithm will find the path.

However, the A* algorithm only returns optimal solutions when the heuristic has some specific properties.

3.5 Condition on heuristics for A* to be optimal:

Whether A* Search is optimal depends on two key properties of the heuristic. These are **Admissibility** and **Consistency**.

Definition 3.5.1: Admissibility

An admissible heuristic never overestimates the cost to reach the goal. This makes it optimistic about the path costs.

Proof of Cost-Optimality (via Contradiction)

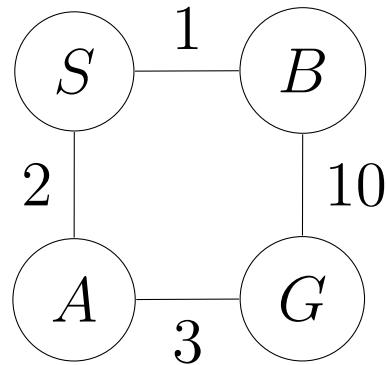
- Assume the optimal path cost is C^* but A* returns a path with a cost greater than C^* .

- There must be a node n on the optimal path that A^* did not expand.
- If $f(n)$ which is the estimated cost of the cheapest solution through n were less or equal to C^* then n would have been expanded.
- By definition and admissibility, $f(n) = g(n) + h(n)$ and should be less or equal to $g^*(n) + \text{cost}(n \text{ to goal}) = C^*$ as n is on the optimal path and $h(n)$ is less than or equal to $\text{cost}(n \text{ to goal})$ due to admissibility
- This is contradicting our assumption that $f(n) > C^*$.
- Thus, A^* must indeed return the cost-optimal path.

Example: Where Violation of admissibility Leads to a Suboptimal Solution

- **Nodes:** Start (S), A, B, Goal (G)
- **Edges with costs:**

- S to $A = 2$
- A to $G = 3$
- S to $B = 1$
- B to $G = 10$



Heuristic, $h(\text{node})$ Estimates to the Goal, G :

- $h(S) = 3$ admissible as $h(S) <$ the optimal path-cost from S to G .
- $h(A) = 10$ is in-admissible as $h(A) >$ optimal path-cost from A to G
- $h(B) = 2$ admissible.
- $h(G) = 0$ admissible.

A^* Algorithm Execution:

- **Start at S :**

$$f(S) = g(S) + h(S) = 0 + 3 = 3.$$

- **Expand S :** Adding Neighbors (A and B) to the queue:

- For A:

$$g(A) = 2, \text{ so, } f(A) = g(A) + h(A) = 2 + 10 = 12.$$

- For B:

$$g(B) = 1, \text{ so, } f(B) = g(B) + h(B) = 1 + 2 = 3$$

- Node A and Node B are currently in queue.
- **Node B selected** from the queue for expansion because of **lower f-value** despite leading to a higher cost path.
- **Expand B :** Adding Neighbor (G):

- For G via B

$$g(G) = 11, \text{ so, } f(G) = g(G) + h(G) = 11 + 0 = 11$$

- Node A and Node G are currently in queue.
- **Node G selected** from the queue for expansion because of lower f-value.
- The algorithm returns path S-B-G which is sub-optimal.

The inadmissibility of the heuristic causes the algorithm to prefer a sub-optimal path.

Definition 3.5.2: Consistency

A heuristic h is consistent if for every node n and every successor n' of n , the estimated cost from n to the goal, denoted $h(n)$, does not exceed the cost from n to n' plus the estimated cost from n' to the goal, $h(n')$. Mathematically, this is represented as: $h(n) \leq c(n, n') + h(n')$ where $c(n, n')$ is the cost to reach n' from n .

A consistent heuristic, also known as a monotonic heuristic, is a stronger condition than admissibility and plays a crucial role in ensuring that A* search finds the optimal path. Here's how it ensures that A* returns an optimal path:

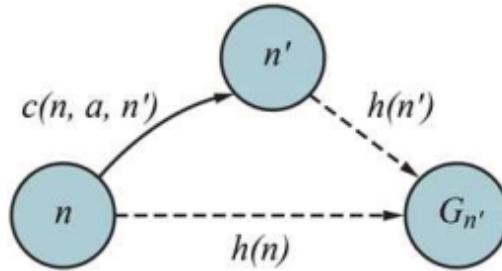
1. Path Cost Non-Decreasing:

- Consistency ensures that the f-value (total estimated cost) of a node n calculated as $f(n) = g(n) + h(n)$ does not decrease as the algorithm progresses from the start node to the goal. This is because for any node n and its successor n' : $g(n') = g(n) + c(n, n')$

- Simplifying, we find:

$$f(n) = g(n) + h(n) \leq g(n) + c(n, n') + h(n') = g(n') + h(n') = f(n')$$

- Therefore, f-values along a path do not decrease, preventing any re-exploration of nodes already deemed sub-optimal, hence streamlining the search towards the goal.



Triangle inequality: If the heuristic h is **consistent**, then the single number $h(n)$ will be less than the sum of the cost $c(n, a, a')$ of the action from n to n' plus the heuristic estimate $h(n')$.

2. Closed Set Invariance:

- When a node n is expanded, its f-value is finalized. Due to the non-decreasing nature of f-values, any path rediscovered through n will have an f-value at least as large as when n was first expanded.
- This prevents the algorithm from revisiting nodes unnecessarily, thereby ensuring efficiency in path finding.

3. Optimal Path Discovery:

- Given the consistency condition, once the goal node g is reached and its $f(g)$ calculated, there can be no other path to g with a lower f-value that has not already been considered.
- Since $h(g) = 0$ (by definition at the goal), $f(g) = g(g)$, meaning that the path cost $g(g)$ represents the total minimal cost to reach the goal from the start node.

4. Optimality Guarantee:

- The search terminates when the goal is popped from the priority queue for expansion, and due to the non-decreasing nature of f-values, this means that no other path with a lower cost can exist that has not already been evaluated.
- Thus, the path to g with cost $g(g)$ must be optimal.

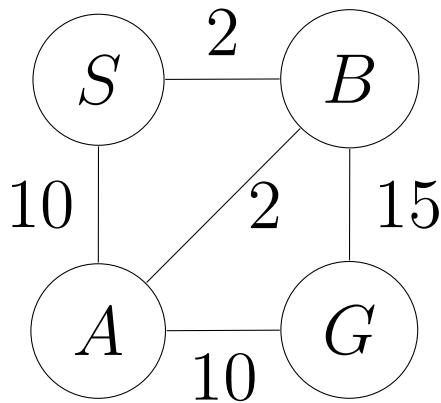
By adhering to these principles, A* search with a consistent heuristic not only finds a solution but ensures it is the optimal one.

Example: Checking inconsistency

- **Nodes:** Start (S), A, B, Goal (G)

- **Edges with costs:**

- S to A = 2
- A to G = 3
- S to B = 1
- B to G = 10



- **Heuristic (h) Estimates to the Goal (G):**

- $h(S) = 12$
- $h(A) = 7$
- $h(B) = 10$
- $h(G) = 0$

- **Checking consistency for each node:** First we find the optimal path to Goal from each node.

- Optimal path from Node S is $S - B - A - G = 14$.
- Optimal path from Node B is $B - A - G = 12$.
- Optimal path from Node A is $A - G = 10$.

- Optimal path from Node G is $G - = 0$.

For any given node, n the heuristic, $h(n)$ for that node is lower or equal than the optimal path-cost from n . So we can say that the given heuristics are admissible. An inadmissible heuristic automatically leads to inconsistent heuristic. But admissible heuristic does not guarantee consistency.

It is better to start checking consistency from the Goal node.

- We see $h(G) = 0$ is admissible and consistent.
- Next, from Node A, there is a direct path to Goal node, G with cost 10 which is the optimal path. There is also a path via B. So to be consistent,

$$\begin{aligned} h(A) &\leq \text{Cost}(A, G) + h(G) \\ &= 10 + 0 = 10 \\ h(A) &= 7 \end{aligned}$$

So, node A is consistent with node G.

$$\begin{aligned} h(A) &\leq \text{Cost}(A, B) + h(B) \\ &= 2 + 15 = 17 \\ h(A) &= 7 \end{aligned}$$

So, heuristic of node A is consistent with node B.

- Next, from Node B, there is a direct path to Goal node, G with cost 15 and a path via nod A costing 12. In this case, G and A are the child of B. So to be consistent,

$$\begin{aligned} h(B) &\leq \text{Cost}(B, G) + h(G) \\ &= 15 + 0 = 15 \\ h(B) &= 10 \end{aligned}$$

$h(B)$ is consistent with node G. However,

$$\begin{aligned} h(B) &\leq \text{Cost}(B, A) + h(A) \\ &= 2 + 7 = 9 \\ h(B) &= 10 \end{aligned}$$

So, heuristic of node B with node A is inconsistent. Similarly for node S,

$$\begin{aligned} h(S) &\leq \text{cost}(S, A) + h(A) = 17 \\ h(S) &\leq \text{cost}(S, B) + h(B) = 12 \\ h(S) &= 12 \end{aligned}$$

Making, $h(S)$ consistent with node A and B.

3.6 How to choose a better Heuristic:

In the previous lecture, we have seen that there are many possible ways to design the heuristics for a problem space. We want to know which heuristic function should be selected when we have more than one way of computing admissible and consistent heuristics.

1. **Effective Branching Factor:** The effective branching factor (**EBF**) is a measure used in tree search algorithms to provide a quantitative description of the tree's growth rate. It reflects how many children each node has, on average, in the search tree that needs to be generated to find a solution. Mathematically, the EBF is defined as the branching factor b for which:

$$N + 1 = 1 + b + b^2 + b^3 + \dots + b^d$$

where N is the total number of nodes generated in the search tree and d is the depth of the shallowest solution.

The EBF gives an insight into the efficiency of the search process, influenced heavily by the heuristic used:

- **Lower EBF:** A lower EBF suggests that the heuristic is effective, as it leads to fewer nodes being expanded. This usually indicates a more directed and efficient search.
- **Higher EBF:** A higher EBF suggests a less effective heuristic, as more nodes are being generated, indicating a broader search, which is generally less efficient.

Calculating the EBF can help evaluate the practical performance of a heuristic. An ideal heuristic would reduce the EBF to the minimum necessary to find the optimal solution, indicating a highly efficient search strategy.

d	Search Cost (nodes generated)			Effective Branching Factor		
	BFS	$A^*(h_1)$	$A^*(h_2)$	BFS	$A^*(h_1)$	$A^*(h_2)$
6	128	24	19	2.01	1.42	1.34
8	368	48	31	1.91	1.40	1.30
10	1033	116	48	1.85	1.43	1.27
12	2672	279	84	1.80	1.45	1.28
14	6783	678	174	1.77	1.47	1.31
16	17270	1683	364	1.74	1.48	1.32
18	41558	4102	751	1.72	1.49	1.34
20	91493	9905	1318	1.69	1.50	1.34
22	175921	22955	2548	1.66	1.50	1.34
24	290082	53039	5733	1.62	1.50	1.36
26	395355	110372	10080	1.58	1.50	1.35
28	463234	202565	22055	1.53	1.49	1.36

Comparison of the search costs and effective branching factors for 8-puzzle problems using breadth-first search, A^* with

h_1 (misplaced tiles), and A^* with

h_2 (Manhattan distance). Data are averaged over 100 puzzles for each solution length d from 6 to 28.

In the above figure, Stuart and Russel generated random 8-puzzle problems and solved them with an uninformed breadth-first search and with A^* search using both and reporting the average number of nodes generated and the corresponding effective branching factor for each search strategy and for each solution length. The results suggest that h_2 is better than h_1 and both are better than no heuristic at all.

2. **Dominating heuristic:** For two heuristic functions, h_1 and h_2 , we say that h_1 dominates h_2 if for every node n in the search space, the following condition holds:

$h_1(n) \geq h_2(n)$ and there is at least one node n where $h_1(n) > h_2(n)$ then we say that h_1 dominates h_2 or in other words, h_1 is the dominating heuristic.

3.6.1 Why is a dominant heuristic for efficient?

By now we should have an understanding that in A^* algorithm every node, n with

$$f(n) \leq C^*, \\ \text{or, } h(n) \leq C^* - g(n) \quad [C^* \text{ is optimal path cost}]$$

is surely expanded.

Now let us consider, two heuristics, h_1 and h_2 , both admissible and consistent. Where,

$$h_2 \geq h_1.$$

If node, n is on the optimal path, it will be expanded with A^* algorithm for both the heuristics. Meaning,

$$h_1(n) \leq h_2(n) \leq C^* - g(n).$$

If node, n is not on the optimal path, we have three possibilities,

- **Possibility 1:**

$$C^* - g(n) \leq h_1 \leq h_2,$$

in which case, node n will not be expanded at all, whichever heuristic we use.

- **Possibility 2:**

$$h_1(n) \leq h_2(n) \leq C^* - g(n).$$

in this case, node, n will be expanded for both the heuristics.

- **Possibility 3:**

$$h_1(n) \leq C^* - g(n) \leq h_2(n),$$

where, node n will be expanded when h_1 is used but not when h_2 is used.

So, we see that the using h_1 may expand the same nodes that are expanded when h_2 is used. But it may end up expanding more unnecessary nodes than those expanded by using the dominant heuristic h_2 .

You may want to go back to the last example in section 3.3 and notice that the number of nodes generated by using the dominant heuristic is significantly less for the 8-puzzle problem.

So, given that the heuristic is consistent and the computation time is not too long, it is generally a better idea to use higher valued heuristic.

CHAPTER 4

LOCAL SEARCH

Note: This lecture closely follows Chapter 4 to 4.1.1 (Local Search and Hill Climbing), 4.1.2 (Simulated Annealing) and 4.1.4 (Evolutionary algorithm) of Russel and Norvig, *Artificial Intelligence: A Modern Approach*.

4.1 Local Search

So far, we have utilized search strategies like A* search and the Greedy Best Search algorithm to navigate vast solution spaces and find sequences of actions that lead to optimal solutions. These strategies use heuristics to estimate costs from any node to the goal, improving efficiency by focusing on promising areas of the state space. However, the effectiveness of these methods depends heavily on the quality of the heuristic used. Informed search strategies are particularly useful for finding paths to the goal state. For example, solving the 8-puzzle requires a series of consecutive actions (such as moving the empty space up, down, left, or right) to reach the solution; similarly, traveling from Arad to Bucharest involves a sequence of actions moving from one connected city to another.

On the other hand, some problems focus solely on generating a goal state or a good state, regardless of the specific actions taken. For instance, in a simplified knapsack problem, the solution involves selecting a set of items that maximizes reward without exceeding the weight limit. The process does not concern itself with the order in which items are checked or selected to achieve the maximum reward. Local search strategies can be used to solve such problems. It is especially useful in problems with very large search spaces.

Local search algorithms offer a practical solution for tackling large and complex problems. Unlike global search methods that attempt to cover the entire state space, local search begins with an initial setup and gradually refines it. It makes incremental adjustments to the solution, exploring nearby possibilities through methods like Hill Climbing, Simulated Annealing, and Genetic Algorithms. These techniques adjust parts of the current solution to systematically explore the immediate area, known as the "neighborhood," for better solutions. This approach is particularly effective in environments with numerous potential solutions (local optima), helping to find an optimal solution more efficiently.

4.1.1 The Need for Local Search

Local search is particularly useful in:

- **Large or poorly understood search spaces, where exhaustive global search is impractical.** For example, in large scheduling problems where the number of potential schedules increases exponentially with the number of tasks and resources.
- **Optimization tasks, focusing on improving a measure rather than finding a path.** This is seen in machine learning hyperparameter tuning, where algorithms like Gradient Descent and its variants iteratively adjust parameters to minimize a loss function.
- **Dynamic problems, where solutions must adapt to changes without restarting the search.** A common example is in real-time strategy games, where AI opponents must continuously adjust their strategies based on the player's actions.
- Local search also complements hybrid algorithms, **combining its strengths with other strategies for enhanced performance** across various problems. For instance, Genetic Algorithms use local search within their crossover and mutation phases to fine-tune solutions.

In summary, local search provides a flexible and efficient approach for refining solutions incrementally, making it a critical tool in modern computational problem-solving. These strategies bridge the gap between the theoretical optimality of informed search and the practical needs of real-world applications.

4.1.2 Examples of problems that can be solved by local search:

Local search algorithms are versatile tools for tackling various optimization problems across many fields. Here are some example problems where local search algorithms are particularly effective:

Traveling Salesman Problem (TSP): In the TSP, the goal is to find the shortest possible route that visits a list of cities and returns to the origin city. A local search algorithm like Simulated Annealing or 2-opt (a simple local search used in TSP) can iteratively improve a given route by swapping the order of visits to reduce the overall travel distance.

Knapsack Problem: As previously mentioned, the goal here is to maximize the value of items placed in a knapsack without exceeding its weight capacity. Local search techniques can be used to iteratively add or remove items from the knapsack to find a combination that offers the highest value without breaching the weight limit.

Max Cut Problem: This is a problem in which the vertices of a graph need to be divided into two disjoint subsets to maximize the number of edges between the subsets. Local search can adjust the placement of vertices in subsets to try and maximize the number of edges that cross between them.

Protein Folding: In computational biology, protein folding simulations involve finding low-energy configurations of a chain of amino acids. Local search can be used to tweak the configuration of the protein to find the structure with the lowest possible energy state.

Layout Design: In manufacturing and architectural design, layout problems involve the placement of equipment or rooms to minimize the cost of moving materials or to maximize accessibility. Local search can rearrange the placements iteratively to improve the overall layout efficiency.

Parameter Tuning in Machine Learning: Tuning hyperparameters of a machine learning model to minimize prediction error or maximize model performance can also be approached as an optimization problem. Techniques like Grid Search, Random Search, or more sophisticated local search methods can be applied to find optimal parameter settings.

4.1.3 Some Constraint Satisfaction Problems can also be solved using local search strategies:

Graph Coloring: This problem requires assigning colors to the vertices of a graph so that no two adjacent vertices share the same color, using the minimum number of colors. Local search can explore solutions by changing the colors of certain vertices to reduce conflicts or the number of colors used.

Job Scheduling Problems: In job scheduling, the task is to assign jobs to resources (like machines or workstations) in a way that minimizes the total time to complete all jobs or maximizes throughput. Job scheduling can be viewed as a CSP when the task is to assign start times to various jobs subject to constraints such as job dependencies (certain jobs must be completed before others can start), resource availability (jobs requiring the same resource cannot overlap), and deadlines. Local search can be used to iteratively shift jobs between resources or reorder jobs to find a more efficient schedule.

Vehicle Routing Problem Similar to TSP but more complex, this problem involves multiple vehicles and aims to optimize the delivery routes from a central depot to various customers. This problem can also be modeled as a CSP, where constraints might include vehicle capacity limits, delivery time windows, and the requirement that each route must start and end at a depot. Local search can adjust routes by reassigning customers to different vehicles or changing the order of stops to minimize total distance or cost.

Local search algorithms are ideal for these and many other problems because they can provide high-quality solutions efficiently, even when the search space is extremely large and complex. They are particularly valuable when exact methods are computationally infeasible, and an approximate solution is acceptable. Local search algorithms operate by searching from a start state to neighboring states, without keeping track of the paths, nor the set of states that have been reached. That means they are **not systematic—they might never explore a portion of the search space where a solution actually resides**. However, they have two key advantages:

1. **They use very little memory; and**
2. **They can often find reasonable solutions in large or infinite state spaces for which systematic algorithms are unsuitable.**

4.2 Local Search Algorithms:

In this course we are going to learn about:

1. Hill-Climbing Algorithm / Gradient Descent
2. Problems of Hill Climbing Algorithms and Remedies
3. Simulated Annealing
4. Genetic Algorithm

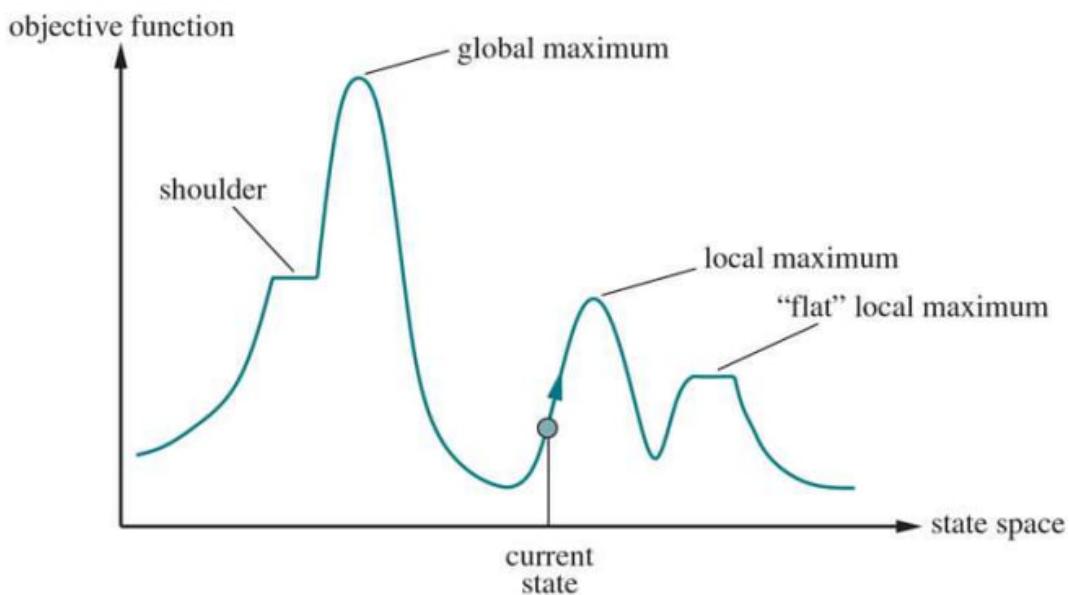
4.2.1 State-Space and Objective Function:

Local search algorithms are powerful tools for addressing optimization problems, where the objective is to find an optimal state that maximizes or minimizes a given objective function. These algorithms navigate through a metaphorical "state-space landscape," where each point or state in this landscape represents a possible solution with a specific "elevation" defined by the objective function. The concept can be more clearly understood through the following points, illustrated with examples:

1. Understanding the State-Space Landscape

Concept: Imagine the state-space of a problem as a geographical landscape where every point represents a possible solution. The elevation at each point is determined by the value of the objective function at that state.

Elevation as Objective Function: In optimization terminology, elevation could represent a value or a cost associated with each state. High elevations indicate better values in maximization problems, and lower elevations indicate lower costs in minimization problems.



A one-dimensional state-space landscape in which elevation corresponds to the objective function. The aim is to find the global maximum.

2. Objective of Local Search

Maximization (Hill Climbing): Here, the goal is to find the highest point in the landscape, akin to climbing to the peak of a hill. The algorithm iteratively moves to higher elevations, seeking to locate the highest peak, which represents the optimal solution.

Example: Maximizing sales in a retail chain by adjusting variables such as pricing, marketing spend, and store layout. Each adjustment is a "step" in the landscape, seeking higher profits (higher elevations).

Minimization (Gradient Descent): In contrast, this approach aims to find the lowest point or valley in the landscape. This is suitable for cost reduction problems where each step attempts to move to lower elevations, minimizing the objective function.

Example: Minimizing production costs in a manufacturing process by altering materials, labor, and energy usage. Each configuration change is a step toward lower costs.

3. Examples of Local Search Algorithms in Action

i. Traveling Salesman Problem (TSP):

Objective function: Minimize the total travel distance for a salesman needing to visit multiple cities and return to the starting point.

Local Search Strategy: Start with a random route and iteratively improve it by swapping the order of cities if it results in a shorter route (seeking lower valleys in terms of distance).

ii. Knapsack Problem

Objective function: Maximize the value of items placed in a knapsack without exceeding its weight capacity.

Local search Strategy: Begin with a random combination of items. Iteratively add or remove items from the knapsack to find a combination that offers the highest value without breaching the weight limit.

iii. 8-queen problem

Objective function: Minimize the number of pairs of queens that are attacking each other either horizontally, vertically, or diagonally. The ideal goal is to reduce this number to zero, indicating that no queens are threatening each other.

Local search strategy: Start with a random instance of the board and iteratively improving the placement of queens on the board.

4.3 Hill-Climbing Search:

Definition 4.3.1: Hill-Climbing Algorithm

The Hill Climbing search algorithm maintains one current state and iteratively moves to the neighboring state with the highest value, effectively following the steepest path upward. It stops when it reaches a "peak," meaning no adjacent state offers a better value. This algorithm only considers the immediate neighbors and does not plan beyond them.

Algorithm 3 Local Maximum Search

```

1: Input: initial_state, objective_function()
2: Output: state that is a local maximum
3: begin
4:   current_state ← initial_state
5:   loop do
6:     neighbor ← a state that is a neighbor of current_state
7:     If objective_function(neighbor) ≤ objective_function(current_state)
8:       Return current_state
9:     current_state ← neighbor

```

Definition 4.3.2: Gradient descent algorithm:

When the goal of the hill-climbing algorithm is to find the state with the minimum value rather than maximizing the objective, it is known as gradient descent.

4.3.1 Examples:**Simplified Knapsack Problem Setup:**

Items (value, weight):

Item 1: Value = \$10, Weight = 2 kg
 Item 2: Value = \$15, Weight = 3 kg
 Item 3: Value = \$7, Weight = 1 kg
 Item 4: Value = \$20, Weight = 4 kg
 Item 5: Value = \$8, Weight = 1 kg

Knapsack Capacity: 7 kg

Objective: Maximize the total value of the items in the knapsack such that the total weight does not exceed 7 kg.

Hill Climbing Algorithm:

Initial Solution Start with a randomly selected set of items that do not exceed the capacity. For this example, let's start with the 1st, 2nd and 3rd item in the knapsack.

The knapsack can be represented as a string of 1s and 0s, where a '1' at i 'th position indicates that the corresponding (i 'th) item has been included, and a '0' means the item has not been taken.

Initial Solution: 11100, weight: 6kg, value: 32

Neighbor Generation: Generate neighboring solutions by toggling the inclusion of each

item. For instance, if an item is not in the knapsack, consider adding it; if it is in the knapsack, consider removing it or replace an item with another item not in the solution.

Evaluate and Select: Calculate the total value and weight for each neighbor. If a neighbor exceeds the knapsack's capacity, discard it.

Choose the neighbor with the highest value that does not exceed the weight capacity.

Iteration: Repeat the process of generating and evaluating neighbors from the current solution.

If no neighbors have a higher value than the current solution, terminate the algorithm.

Termination: The algorithm stops when it finds a solution where no neighboring configurations offer an improvement.

Demonstration (for two iterations):

Step 1: Initial Solution

Knapsack: 11100

Total Value: \$32

Total Weight: 6 kg

Step 2: Generate Neighbors Add Item 4: 11110, Total Value = \$52, Total Weight = 10 kg

Add Item 5: 11101, Total Value = \$40, Total Weight = 7 kg

Replace Item 2 with Item 4: 10110, Total Value = \$37, Total weight: 7kg/

..... There can be other possible neighbors.

Step 3: Evaluate and Select

Discard Neighbor: 11110 as it exceeds knapsack weight limit.

Best neighbor: 11101 Total Value = 40, Total Weight = 7 kg

Step 4: Iteration (next steps)

Current configuration: 11101

Add Item 4: 11111, Total Value = \$60, Total Weight = 11 kg Replace items with Item 4:
All will exceed the weight

Evaluation:

All neighbors are discarded. No better valued neighbor.

Terminates as no better valued neighbor is found.

4.3.2 Key Characteristic:

Hill-Climbing Search is essentially a greedy approach, meaning it always looks for the best immediate improvement at each step, without considering long-term consequences. This can sometimes lead to suboptimal solutions.

4.3.3 Drawbacks:

Hill climbing algorithms may encounter several challenges that can cause them to get stuck before reaching the global maximum. These challenges include:

1. **Local Maxima:** A local maximum is a peak that is higher than each of its neighboring states but is lower than the global maximum. Hill climbing algorithms that reach a local maximum are drawn upward toward the peak but then have nowhere else to go because all nearby moves lead to lower values.

Example:

8-queen problem:

Background: In a 8-queens problem, our goal is to maximize the number of non-attacking pairs of queens. The maximum number of pairs possible from 8 queens is $8*7/2 = 28$. So, we want all 28 pairs to be in non-attacking positions.

Let us represent the instances as an array of 8 numbers, ranging from 1 to 8 denoting the columns and the index, ranging from 1 to 8 (sorry, programmers!) denoting the rows.

Scenario: For an instance of the 8-puzzle, 13528647, there are 5 attacking pairs or in other words, 23 non-attacking pairs. This configuration is a local maximum because it's better than all immediate neighboring configurations, but it is not the global solution since it's not conflict-free.

Problem: The Hill Climbing algorithm would stop here because all single-move alternatives lead to worse configurations, increasing the number of threats. Despite the presence of better configurations (global maxima with zero threats), the algorithm gets stuck.

2. **Plateaus:** A plateau is an area of the state-space landscape where the elevation (or the value of the objective function) remains constant. Hill climbing can become stuck on a plateau because there is no upward direction to follow. Plateaus can be particularly challenging when they are flat local maxima, with no higher neighboring states to escape to, or when they are shoulders, which are flat but eventually lead to higher areas. On a plateau, the algorithm may wander aimlessly without making any progress. **Sce-**

nario: Imagine a large part of the chessboard setup where several queens are placed in such a manner that moving any one of them doesn't change the number of conflicts—it remains constant. This flat area in the search landscape is a plateau.

For Example, the instance, 13572864 has 3 attacking pairs. Swapping the last two queens might not immediately lead to an increase in non-attacking pairs, resulting in a plateau where many configurations have an equal number of non-attacking pairs.

Problem: The Hill Climbing algorithm would find it difficult to detect any better move since all look equally non-promising.

On a plateau, every move neither improves nor worsens the situation, causing Hill Climbing to wander aimlessly without clear direction toward improvement. This lack of gradient (change in the number of conflicts) can trap the algorithm in non-productive cycles, preventing it from reaching configurations that might lead to the global maximum.

3. **Ridges:** Ridges are sequences of local maxima that make it very difficult for greedy algorithms like hill climbing to navigate effectively. Because the algorithm typically makes decisions based on immediate local gains, it struggles to cross over ridges that require a temporary decrease in value to ultimately reach a higher peak.

These challenges highlight the limitations of hill climbing algorithms in exploring complex landscapes with multiple peaks and flat areas, making them susceptible to getting stuck without reaching the best possible solution.

Scenario: Consider a situation where moving from one configuration to another better configuration involves moving through a worse one. For example, the instance 13131313 creates a ridge because small moves from this configuration typically result in fewer non-attacking pairs, requiring several coordinated moves to escape this pattern, which hill climbing does not facilitate well.

Problem: Hill Climbing may fail to navigate such transitions because it does not allow for a temporary increase in the objective function (number of conflicts in this case). Ridges in the landscape can make it very challenging for the algorithm to find a path to the optimal solution since each step must immediately provide an improvement.

4.3.4 Remedies to problems of Hill-Climbing Search Algorithm:

Hill climbing has several variations that address its basic version's limitations, such as getting stuck at local maxima, navigating ridges ineffectively, or wandering on plateaus.

4.3.5 Variants of Hill Climbing:

- **Stochastic Hill Climbing:**

How it works: This method randomly chooses among uphill moves rather than always selecting the steepest ascent. The probability of choosing a move can depend on the steepness of the ascent.

Performance: It typically converges more slowly than the standard hill climbing because it might not take the steepest path. However, it can sometimes find better solutions in complex landscapes where the steepest ascent might lead straight to a local

maximum.

Example in 8-Queens: In a landscape where queens are close to forming a solution but are stuck due to subtle needed adjustments, stochastic hill climbing can randomly explore different moves that gradually lead out of a local maximum.

- **First-Choice Hill Climbing:**

How it works: A variant of stochastic hill climbing, this strategy generates successors randomly and moves to the first one that is better than the current state.

Advantages: This is particularly effective when a state has a vast number of successors, as it reduces the computational overhead of generating and evaluating all possible moves.

Example in Knapsack Problem: With thousands of possible item combinations, generating all potential successors to evaluate the best move is impractical. First-choice hill climbing can quickly find a better solution without exhaustive comparisons.

- **Random-Restart Hill Climbing:**

How it works: This approach involves performing multiple hill climbing searches from different randomly generated initial states. This process repeats until a satisfactory solution is found.

Completeness: It is "complete with probability 1" because it is guaranteed to eventually find a goal state, assuming there is a non-zero probability of any single run succeeding. If each hill-climbing search has a probability of success, p then the expected number of restarts required is $\frac{1}{p}$.

Example in Knapsack Problem: If the algorithm gets stuck in a sub-optimal configuration due to local maxima, restarting with a new random set of items can lead to discovering better solutions.

Example in 8-Queens: Due to the complex landscape of the 8-Queens problem, random restarts can help escape from sub-optimal arrangements by exploring new configurations that might be closer to the global maximum (i.e., a solution with zero attacking pairs).

4.4 Introduction to Simulated Annealing

Simulated Annealing (SA) is a probabilistic technique for approximating the global optimum of a given function. It is particularly useful for solving large optimization problems where other methods might be too slow or get stuck in local optima.

The technique is inspired by the physical process of annealing in metallurgy, where metals are heated to a high temperature and then cooled according to a controlled schedule to achieve a more stable crystal structure.

The method was developed by S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi in 1983, based on principles of statistical mechanics.

Simulated Annealing is an optimization technique that simulates the heating and gradual cooling of materials to minimize defects and achieve a stable state with minimum energy.

It starts with a randomly generated initial solution and a high initial temperature, which allows the **acceptance of suboptimal solutions** to explore the solution space widely.

The algorithm iterates by generating small modifications to the current solution, evaluating the cost changes, and **probabilistically deciding** whether to accept the new solution based on the current temperature.

The temperature is gradually reduced according to a predetermined cooling schedule, which **decreases the likelihood of accepting worse solutions** and helps fine-tune towards an optimal solution.

The process concludes once a stopping criterion, like a specified number of iterations, a minimal temperature, or a quality threshold of the solution, is met.

Algorithm 4 Algorithm Simulated Annealing

```

1: Input:           initial_solution, initial_temperature, cooling_rate,
   stopping_temperature
2: Output: best_solution_found
Procedure:
3: current_solution ← initial_solution
4: current_temperature ← initial_temperature
5: best_solution ← current_solution
6: while: current_temperature > stopping_temperature
7:   new_solution ← generate_neighbor(current_solution)
8:   cost_difference ← cost(new_solution) - cost(current_solution)
9:   If cost_difference < 0 or exp(-cost_difference / current_temperature) >
      random(0, 1) current_solution ← new_solution
10:  current_solution ← new_solution
11:  If cost(new_solution) < cost(best_solution)
12:    best_solution ← new_solution
13:  current_temperature ← current_temperature × cooling_rate
14: return best_solution

```

Key Parameters:

- **Initial Temperature:** High enough to allow exploration.
- **Cooling Rate:** Determines how quickly the temperature decreases.
- **Stopping Temperature:** Low enough to stop the process once the system is presumed to have stabilized.

4.5 How does simulated annealing navigate the solution space:

4.5.1 Probability of Accepting a New State

The key mathematical concept in simulated annealing is the probability of accepting a new state S' from a current state S . This probability is determined using the Metropolis-Hastings algorithm, which is defined as follows:

$$P(\text{accept } S') = \min \left(1, e^{-\frac{\Delta E}{T}} \right)$$

where:

- $\Delta E = E(S') - E(S)$ is the change in the objective function (cost or energy) from the current state S to the new state S' .
- T is the current temperature.
- e is the base of the natural logarithm.

The equation $e^{-\frac{\Delta E}{T}}$ is crucial as it controls the acceptance of new solutions:

- If $\Delta E < 0$ (meaning S' is a better solution than S), then $e^{-\frac{\Delta E}{T}} > 1$, and the new solution is always accepted ($\min(1, \text{value}) = 1$).
- If $\Delta E > 0$ (meaning S' is worse), the new solution is accepted with a probability less than 1. This probability decreases as ΔE increases or as T decreases.

4.5.2 Cooling Schedule

The cooling schedule is a rule or function that determines how the temperature T decreases over time. It is typically a function of the iteration number k . A common choice is the exponential decay given by:

$$T(k) = T_0 \cdot \alpha^k$$

where:

- T_0 is the initial temperature.
- α is a constant such that $0 < \alpha < 1$, often close to 1.
- k is the iteration index.

The choice of α and T_0 influences the convergence of the algorithm. A slower cooling (higher α) allows more thorough exploration of the solution space but takes longer to converge.

4.5.3 Random Selection of Neighbors

The random selection of a neighbor S' is typically governed by a neighborhood function, which defines possible transitions from any given state S . The randomness allows the algorithm to explore the solution space non-deterministically, which is essential for escaping local optima.

4.5.4 Mathematical Convergence

Theoretically, given an infinitely slow cooling (i.e., $\alpha \rightarrow 1$ and infinitely many iterations), simulated annealing can converge to a global optimum. This stems from the ability to continue exploring new states with a non-zero probability, provided the cooling schedule allows sufficient time at each temperature level for the system to equilibrate.

Simulated annealing integrates concepts from statistical mechanics with optimization through a controlled random walk. This walk is guided by the probabilistic acceptance of new solutions, which balances between exploiting better solutions and exploring the solution space broadly, moderated by a temperature parameter that systematically decreases over time.

4.6 Example Problem:

4.6.1 Traveling Salesman Problem (TSP)

Let's consider an example of solving a small Traveling Salesman Problem (TSP) using simulated annealing. The TSP is a classic optimization problem where the goal is to find the shortest possible route that visits a set of cities and returns to the origin city, visiting each city exactly once.

Problem Setup

Suppose we have a set of five cities, and the distances between each pair of cities are given by the following symmetric distance matrix:

	A	B	C	D	E
A	0	12	10	19	8
B	12	0	3	5	6
C	10	3	0	6	7
D	19	5	6	0	4
E	8	6	7	4	0

Objective

Find the shortest path that visits each city exactly once and returns to the starting city.

Simulated Annealing Steps

1. **Initialization:** Randomly generate an initial route, say, $A \rightarrow B \rightarrow C \rightarrow D \rightarrow E \rightarrow A$.

2. **Heating:** Set an initial high temperature to allow significant exploration. For instance, start with a temperature of 100.
3. **Iteration:**
 - **Generate a Neighbor:** Create a new route by making a small change to the current route, such as swapping two cities. For instance, swap cities B and D to form a new route $A \rightarrow D \rightarrow C \rightarrow B \rightarrow E \rightarrow A$.
 - **Calculate the Change in Cost:** Determine the total distance of the new route and compare it with the current route.
 - **Acceptance Decision:** Use the Metropolis criterion to decide probabilistically whether to accept the new route based on the change in cost and the current temperature.
4. **Cooling:** Reduce the temperature based on a cooling schedule, e.g., multiply the temperature by 0.95 after each iteration.
5. **Termination:** Repeat the iteration process until the temperature is low enough or a fixed number of iterations is reached. Assume the stopping condition is when the temperature drops below 1 or after 1000 iterations.

Example Calculation

Assuming the first iteration starts with the initial route and the randomly generated neighbor as described:

- Current Route: $A \rightarrow B \rightarrow C \rightarrow D \rightarrow E \rightarrow A$ (Distance = $12 + 3 + 6 + 4 + 8 = 33$).
- New Route: $A \rightarrow D \rightarrow C \rightarrow B \rightarrow E \rightarrow A$ (Distance = $19 + 6 + 3 + 6 + 8 = 42$).

The change in cost ΔE is $42 - 33 = 9$. The probability of accepting this worse solution at temperature 100 is $e^{-9/100} \approx 0.91$. A random number is generated between 0 and 1; if it is less than 0.91, the new route is accepted, otherwise, the algorithm retains the current route.

This process is repeated, with the temperature decreasing each time, until the termination conditions are met. The result should be a route that approaches the shortest possible loop connecting all five cities.

4.7 Genetic algorithm

Genetic algorithms (GAs) are a class of optimization algorithms inspired by the principles of natural selection and genetics. These algorithms are used to solve search and optimization problems and are particularly effective for complex issues that are difficult to solve using traditional methods.

Genetic algorithms mimic the process of natural evolution, embodying the survival of the fittest among possible solutions. The core idea is derived from the biological mechanisms of

reproduction, mutation, recombination, and selection. These biological concepts are translated into computational steps that help in finding optimal or near-optimal solutions to problems across a wide range of disciplines including engineering, economics, and artificial intelligence.

4.7.1 Algorithm

Algorithm 5 Genetic Algorithm

```
1: Input: Population size, fitness function, mutation rate, crossover rate, maximum generations
2: Output: Best solution found
3: begin
4:   Initialize population with random candidates
5:   Evaluate the fitness of each candidate
6:   while termination condition not met do
7:     Select parents from the current population
8:     Perform crossover on parents to create new offspring
9:     Perform mutation on offspring
10:    Evaluate the fitness of new offspring
11:    Select individuals for the next generation
12:    best solution in new generation > best solution so far
13:      Update best solution found
14:   end while
15:   return best solution found
16: end
```

4.7.2 Explanation of the Pseudocode

Initialize Population:

A genetic algorithm begins with a population of randomly generated individuals. Each individual, or chromosome, represents a possible solution to the problem. Depending on the problem the chromosomes are encoded.

4.7.3 Evaluate Fitness:

Each solution in the population is assessed using the fitness function to determine how well it solves the problem. Depending on the problem the fitness function will be measured.

Selection:

This step involves choosing the fitter individuals to reproduce. Selection can be done in various ways, such as Truncation Selection, tournament selection, roulette wheel selection, or rank selection. In this course we are using truncation selection, where we select the fittest 3/4th of the population.

Truncation Selection

Description: Only the top-performing fraction of the population is selected to reproduce.

Procedure: Rank individuals by fitness, then select the top $x\%$ to become parents of the next generation.

Pros and Cons: Very straightforward and ensures high-quality genetic material is passed on, but can quickly reduce genetic diversity.

Crossover (Recombination):

Pairs of individuals are crossed over at a randomly chosen point to produce offspring. The crossover rate determines how often crossover will occur. Two common types of crossover techniques are single-point crossover and two-point (or double-point) crossover.

- **Single-Point Crossover:**

In single-point crossover, a single crossover point is randomly selected on the parent chromosomes. The genetic material (bits, characters, numbers, depending on the encoding of the solution) beyond that point in the chromosome is swapped between the two parents. This results in two new offspring, each carrying some genetic material from both parents.

Procedure:

Select a random point on the chromosome. The segments of the chromosomes after this point are swapped between the two parents.

Example:

Suppose we have two binary strings:

Parent 1: 110011

Parent 2: 101010

Assuming the crossover point is after the third bit, the offspring would be:

Offspring 1: 110010 (first three bits from Parent 1, last three bits from Parent 2)

Offspring 2: 101011 (first three bits from Parent 2, last three bits from Parent 1)

- **Two-Point Crossover:**

Two-point crossover involves two points on the parent chromosomes, and the genetic material located between these two points is swapped between the parents. This can

introduce more diversity compared to single-point crossover because it allows the central segment of the chromosome to be exchanged, potentially combining more varied genetic information from both parents.

Procedure:

Select two random points on the chromosome, ensuring that the first point is less than the second point.

Swap the segments between these two points from one parent to the other.

Example:

Continuing with the same parent strings:

Parent 1: 110011

Parent 2: 101010

Let's choose two crossover points, between the second and fifth bits. The offspring produced would be:

Offspring 1: 100010 (first two bits from Parent 1, middle segment from Parent 2, last bit from Parent 1)

Offspring 2: 111011 (first two bits from Parent 2, middle segment from Parent 1, last bit from Parent 2)

4.8 Mutation

With a certain probability (mutation rate), mutations are introduced to the offspring to maintain genetic diversity within the population.

The purpose of mutation is to maintain and introduce genetic diversity into the population of candidate solutions, helping to prevent the algorithm from becoming too homogeneous and getting stuck in local optima.

4.8.1 Purpose of Mutation

1. **Introduce Variation:** Mutation introduces new genetic variations into the population by altering one or more components of genetic sequences, ensuring a diversity of genes.
2. **Prevent Local Optima:** By altering the genetic makeup of individuals, mutation prevents the population from converging too early on a suboptimal solution.
3. **Explore New Areas:** It enables the algorithm to explore new areas of the solution space that may not be reachable through crossover alone.

4.8.2 How Mutation Works

Mutation operates by making small random changes to the genes of individuals in the population. In the context of genetic algorithms, an individual's genome might be represented

as a string of bits, characters, numbers, or other data structures, depending on the problem being solved.

4.8.3 Common Types of Mutation

1. Bit Flip Mutation (for binary encoding):

- **Procedure:** Each bit in a binary string has a small probability of being flipped (0 changes to 1, and vice versa).
- **Example:** A binary string ‘110010‘ might mutate to ‘110011‘ if the last bit is flipped.

2. Random Resetting (for integer or real values):

- **Procedure:** A selected gene is reset to a new value within its range.
- **Example:** In a string of integers [4, 12, 7, 1], the third element 7 might mutate to 9.

3. Swap Mutation:

- **Procedure:** Two genes are selected and their positions are swapped. This is often used in permutation-based encodings.
- **Example:** In an array [3, 7, 5, 8], swapping the second and fourth elements results in [3, 8, 5, 7].

4. Scramble Mutation:

- **Procedure:** A subset of genes is chosen and their values are scrambled or shuffled randomly.
- **Example:** In an array [3, 7, 5, 8, 2], scrambling the middle three elements might result in [3, 5, 8, 7, 2].

5. Uniform Mutation (for real-valued encoding):

- **Procedure:** Each gene has a fixed probability of being replaced with a uniformly chosen value within a predefined range.
- **Example:** In an array of real numbers [0.5, 1.3, 0.9], the second element 1.3 might mutate to 1.1.

4.8.4 Mutation Rate

The mutation rate is a critical parameter in genetic algorithms. It defines the probability with which a mutation will occur in an individual gene. A higher mutation rate increases diversity but may also disrupt highly fit solutions, whereas a lower mutation rate might not provide enough diversity, leading to premature convergence. The optimal mutation rate often depends on the specific problem and the characteristics of the population.

Evaluate New Offspring: The fitness of each new offspring is calculated.

Generation Update: The algorithm decides which individuals to keep for the next generation. This can be a mix of old individuals (elitism) and new offspring.

Termination Condition: The algorithm repeats until a maximum number of generations is reached, or if another stopping criterion is satisfied (like a satisfactory fitness level).

Return Best Solution: The best solution found during the evolution is returned.

Note: For exam problems if you are asked to simulate, unless otherwise instructed, start with 4 chromosomes in the population, select best 3 at each step, crossover between the best and the other two selected

4.8.5 Diversity in Genetic Algorithm

It is often the case that the population is diverse early on in the process, so crossover frequently takes large steps in the state space early in the search process (as in simulated annealing). After many generations of selection towards higher fitness, the population becomes less diverse, and smaller steps are typical.

It is important for the initial population to be diversified. Otherwise, similar type of chromosomes will crossover to and produce offsprings with little change in their fitness. This will lead to quick convergence of the algorithm and the chances of finding a solution with maximum fitness will be very low.

4.8.6 Advantages and Applications

Genetic algorithms are particularly useful when:

- The search space is large, complex, or poorly understood.
- Traditional optimization and search techniques fail to find acceptable solutions.
- The problem is dynamic and changes over time, requiring adaptive solutions.

4.9 Examples

4.9.1 8-Queen Problem

The 8-queens problem involves placing eight queens on an 8x8 chessboard so that no two queens threaten each other. This means no two queens can share the same row, column, or diagonal.

Chromosome Representation:

Each chromosome can be represented as a string or array of 8 integers, each between 1 and 8, representing the row position of the queen in each column represented by the index of the string or array.

Fitness Function:

Fitness is calculated based on the number of non-attacking pairs of queens. The maximum score for 8 queens is 28 (i.e., no two queens attack each other).

Calculating the number of attacking pairs: The 8-queen problem has the following conditions.

- No pairs share the same column.
- No pairs share the same row.
- No pairs share the same diagonal.

Column-wise conflict Now, due to the representation of the configuration where we have ensured that no pairs can share the same column as the indices of the array are unique.

Row-wise conflict Now, the values in each index represent the row where the queen is placed. Now, if we find the same value in multiple indices that means there are queens sharing that row.

Take for example, the configuration: [8, 7, 7, 3, 7, 1, 4, 4]. Here, the value 7 is repeated thrice and the value 4 is repeated twice. This means, there are 3 queens on the 7th row and 2 queens in the 4th row.

Counting the conflicts in 7th row: 3 queens will form ${}^3C_2 = 3(3 - 1)/2 = 3$ pairs.

Counting the conflicts in 4th row: 2 queens will form ${}^2C_2 = 2(2 - 1)/2 = 1$ pairs.

Total 4 row-wise conflicting pairs.

Diagonal Conflicts Two types of diagonals are formed in a square board we call them Major ("/" shaped) diagonal, and Minor ("\\" shaped) diagonal. **Major diagonal conflict**

If two queens, Q_1 and Q_2 share the same major diagonal conflict $\text{abs}(Q_1[\text{column}] - Q_1[\text{row}]) = \text{abs}(Q_2[\text{column}] - Q_2[\text{row}])$.

Going back to the configuration: [8, 7, 7, 3, 7, 1, 4, 4], the queen in 7th row, 2nd column is in the same diagonal as the queen in 1st row, 6th column and the queen in 7th row, 3rd column is in conflict with the queen in 4th row, 8th column. Or in other words 2 queens share the 5th ($|7-2| = |1-6|$) Major diagonal making (${}^2C_2 = 2(2-1)/2 = 1$) conflicting pair and 2 queens share the 4th ($|7-3| = |4-8|$) Major diagonal making another (1) conflicting pair.

So, there are two attacking pairs along the major diagonal.

Minor diagonal conflict If two queens, Q_1 and Q_2 share the same major diagonal conflict $Q_1[\text{column}] + Q_1[\text{row}] = Q_2[\text{column}] + Q_2[\text{row}]$.

Going back to the configuration: [8, 7, 7, 3, 7, 1, 4, 4], the queen in 8th row, 1st column is in the same diagonal as the queen in 7th row, 2nd column, the queen in 3rd row, 4th column is in conflict with the queen in 1st row, 6th column and the queen in 7th row, 5th column is in conflict with the queen in 4th row, 8th column. Or in other words 2 queens share the 9th ($|8+1|=|7+2|$) minor diagonal making 1 (${}^2C_2 = 2(2-1)/2$) conflicting pair, 2 queens share the 7th ($|3+4| = |1+6|$) minor diagonal making another (${}^2C_2 = 2(2-1)/2$) conflicting pair and 2 queens share the 12th ($|7+5| = |4+8|$) minor diagonal making another (${}^2C_2 = 2(2-1)/2$) conflicting pair.

So, there are three attacking pairs along the minor diagonal.

Summing up the number of attacking pairs:

- 4 row-wise attacking pairs.
- 2 attacking pairs on the major diagonal.
- 3 attacking pairs on the minor diagonal.
- In total ($4 + 2 + 3 = 9$) attacking pairs.

Calculating the number of non-attacking pairs: The maximum number of pairs formed from n number of queens is ${}^nC_2 = \frac{n(n-1)}{2}$. The maximum number of pairs that can be formed by 8 queens is $\frac{8(8-1)}{2} = 28$. In the ideal configuration with zero attacking pairs, we can have all 28 pairs in non-attacking mode. Thus, in any configuration,

No. of non-attacking pairs = No. of max possible non-attacking pairs - No. of attacking pairs.

For the example of [8, 7, 7, 3, 7, 1, 4, 4] configuration of the 8-queen problem, the max possible non-attacking pairs is 28 and the total no. of attack we calculated is 9. So the number of non-attacking pairs in this configuration is $28 - 9 = 19$.

So, the Fitness value of this configuration is 19.

A better way to calculate fitness is to converting the value within a range of (0 to 1) or in a $x\%$.

$$\text{fitness} = \frac{\text{no. of non-attacking pairs}}{28}$$

However, for simplicity, we are going to take the no. of non-attacking pairs as our fitness value.

Iteration 1

Initialization of Population:

We randomly generate a small population of 4 individuals for simplicity. To keep track of the individual with the best fitness so far we initially set best-so-far = null and best-fitness-so-far = 0 and update when we find better individuals.

```
best-so-far = [];
best-fitness-so-far = 0;
max-fitness = 28;
```

Step 1: Population

Initial Population

Chromosome 1: 4 2 7 3 6 8 5 1

Chromosome 2: 2 7 4 1 8 5 3 6

Chromosome 3: 5 3 1 7 2 8 6 4

Chromosome 4: 7 1 4 2 8 5 3 6

Step 2: Calculate Fitness

	Chromosomes	Fitness
Chromosome 1:	4 2 7 3 6 8 5 1	26
Chromosome 2:	2 7 4 1 8 5 3 6	24
Chromosome 3:	5 3 1 7 2 8 6 4	23
Chromosome 4:	7 1 4 2 8 5 3 6	24

As we have found a chromosome with better fitness value we have saved before, we update,

```
best-so-far = [4 2 7 3 6 8 5 1];
best-fitness-so-far = 26;
```

However, the fitness is not the highest possible value of 28, so we continue to the next step.

Step 3: Selection

Select the top 3/4th of chromosomes based on their fitness. This results in selecting Chromosome 1, 2 and 3.

	Chromosomes	Fitness	Selection
Chromosome 1:	4 2 7 3 6 8 5 1	26	selected
Chromosome 2:	2 7 4 1 8 5 3 6	24	selected
Chromosome 3:	5 3 1 7 2 8 6 4	23	not selected
Chromosome 4:	7 1 4 2 8 5 3 6	24	selected

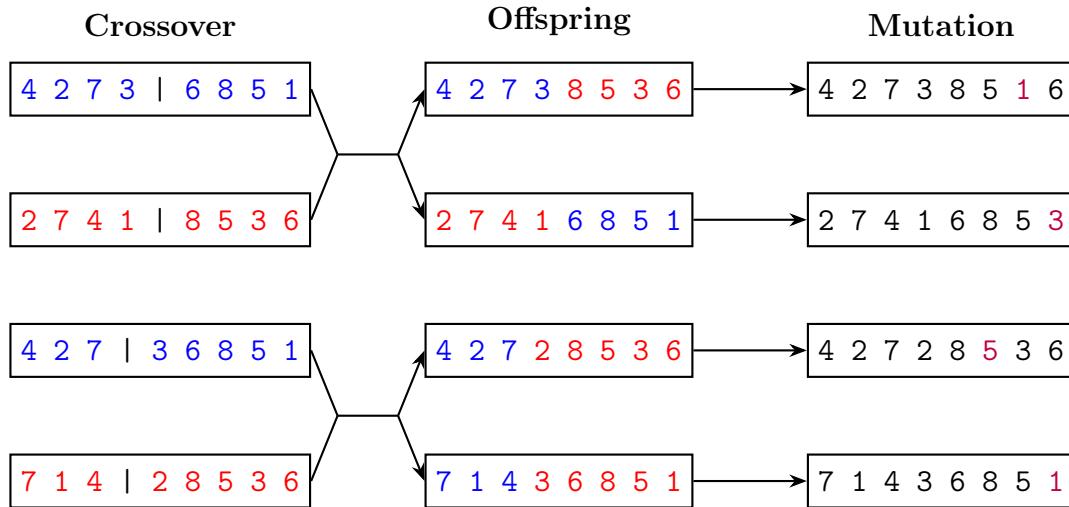
	Selected Chromosomes	Fitness
Chromosome 1:	4 2 7 3 6 8 5 1	26
Chromosome 2:	2 7 4 1 8 5 3 6	24
Chromosome 3:	7 1 4 2 8 5 3 6	24

Step 4: Crossover (Single Point)

We pair the best chromosome [4, 2, 7, 3, 8, 5, 1, 6] with the other two selected chromosomes to create new offspring.

Step 5: Mutation

For each offspring, We randomly choose an index and change the value of the row to a number chosen randomly between 1 to 8. As the number is chosen randomly it may happen that the value remains same as we see for offspring 3 and 4.



Step 6: Create new Population replacing the old population

After the mutation is complete we replace the previous population with the set of new chromosomes and repeat from step 2 with the new population until the configuration with highest fitness (non-attacking pair) is found.

New Population	
Chromosome 1:	4 2 7 3 8 5 1 6
Chromosome 2:	2 7 4 1 6 8 5 3
Chromosome 3:	4 2 7 1 8 5 3 6
Chromosome 4:	2 7 4 3 6 8 5 1

4.9.2 Traveling Salesman Problem (TSP)

Problem Setup:

We have five cities: A, B, C, D and E. The distance matrix given:

	A	B	C	D	E
A	0	2	9	10	7
B	2	0	6	5	8
C	9	6	0	12	10
D	10	5	12	0	15
E	7	8	10	15	0

Genetic Algorithm Setup

Chromosome Representation:

A permutation of the cities $[A, B, C, D, E]$.

Fitness Function:

The fitness of each chromosome is calculated as the inverse of the total route distance. The shorter the route, the higher the fitness.

$$\text{Fitness} = \frac{1}{\text{Route Distance}}$$

Example Calculation: For the chromosome $[A, B, C, D, E]$:

Distance $(A - B)$: 2

Distance $(B - C)$: 6

Distance $(C - D)$: 12

Distance $(D - E)$: 15

Distance $(E - A$ to complete the loop): 7

Total Distance: $2 + 6 + 12 + 15 + 7 = 42$

$$\text{Fitness} = \frac{1}{\text{Total Distance}} = \frac{1}{42}$$

Iteration 1

Initialization of Population:

We start with a population of four randomly generated routes as combination of the cities. We keep track of the best found route and its fitness.

```
best-so-far = [];
best-fitness-so-far = 0;
```

Initial Population

Chromosome 1:	A B C D E
Chromosome 2:	B E D C A
Chromosome 3:	E D B A C
Chromosome 4:	B D C E A

Step 2: Fitness Calculation

Lower distance means higher fitness in this problem. The fitness for each chromosome in the population is calculated.

	Chromosomes	Fitness	Selection
Chromosome 1:	A B C D E	$\frac{1}{42}$	selected
Chromosome 2:	B E D C A	$\frac{1}{45}$	selected
Chromosome 3:	E D B A C	$\frac{1}{49}$	not selected
Chromosome 4:	B D C E A	$\frac{1}{39}$	selected

Step 3: Selection

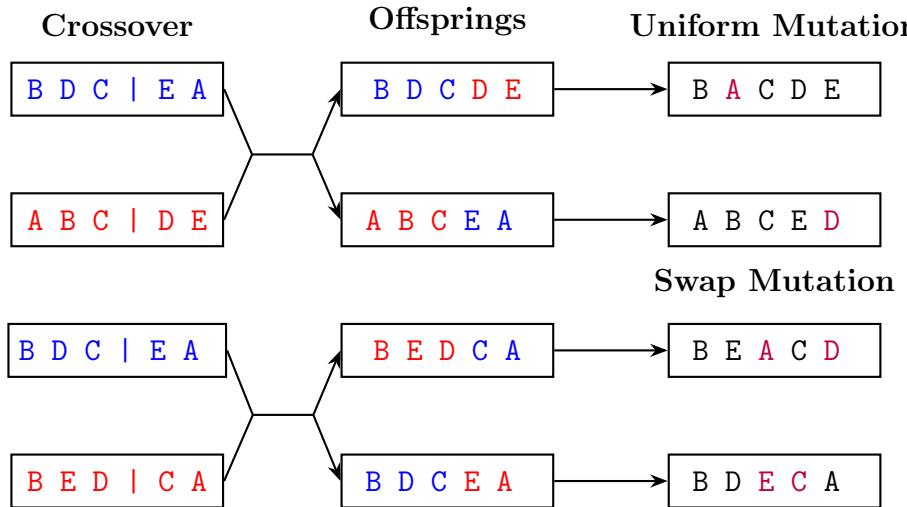
Select top 3/4th based on fitness: Chromosome 4, Chromosome 1 and Chromosome 2. Chromosome 4 has a higher fitness value than the **best-fitness-so-far** we have. So, we make an update.

```
best-so-far = [B D C E A];
best-fitness-so-far = 1/39;
```

	Selected Chromosomes	Fitness
Chromosome 1:	A B C D E	$\frac{1}{42}$
Chromosome 2:	B E D C A	$\frac{1}{45}$
Chromosome 3:	B D C E A	$\frac{1}{39}$

Step 4: Crossover (One-Point)

We make pairs between chromosome with the highest fitness, [B D C E A] with the other two chromosomes that are selected. The crossover point for the pairs are randomly generated. The crossover point in this case for both the pairs were randomly chosen to be 3.



Step 5: Mutation

Offspring 1 is mutated by randomly changing the second city from D to A. Offspring 2 is mutated by changing the fifth city from A to D.

On the other hand, Offspring 3 and Offspring 4 are mutated by swapping cities. In Offspring 3, the third and fifth cities are swapped. In Offspring 4, the third and fourth cities are swapped.

Note: It is better to use one specific type of mutation in your simulation. Always mention what type of mutation you are using.

Step 6: New Population

Replace previous population with the newly generated chromosomes.

New Population

Chromosome 1: **B A C D E**

Chromosome 2: **A B C E D**

Chromosome 3: **B E A C D**

Chromosome 4: **B D E C A**

Step 7: Repeat

This process is repeated from step 2 over several generations to find the chromosome with the highest fitness, representing the shortest possible route that visits each city exactly once and returns to the starting point.

4.9.3 0/1 Knapsack Problem

Problem Setup

Objective

Maximize the value of items packed in a knapsack without exceeding the weight limit.

Constraints:

Maximum weight the knapsack can hold: 15 kg Items:

Item	Weight (Kg)	Value (\$)
1	6	30
2	3	14
3	4	16
4	2	9

Genetic Algorithm Setup

Chromosome Representation:

Each chromosome is a string of four bits, where each bit represents whether an item is included (1) or not (0). For example, '1010' means Items 1 and 3 are included, while Items 2 and 4 are not.

Step 1: Initialization

Generate four random chromosomes (combination of the items) making the initial population. We keep track of the best found route and its fitness.

```
best-so-far = [];
best-fitness-so-far = 0;
```

Initial Population

Chromosome 1: 1 1 0 1

Chromosome 2: 1 0 1 0

Chromosome 3: 0 1 1 0

Chromosome 4: 1 0 0 1

Step 2: Fitness Evaluation

Calculate the total value of each chromosome, ensuring the weight does not exceed the capacity.

Chromosomes	Fitness
Chromosome 1:	1 1 0 1
Chromosome 2:	1 0 1 0
Chromosome 3:	0 1 1 0
Chromosome 4:	1 0 0 1

capacity.
30
39

Step 3: Selection

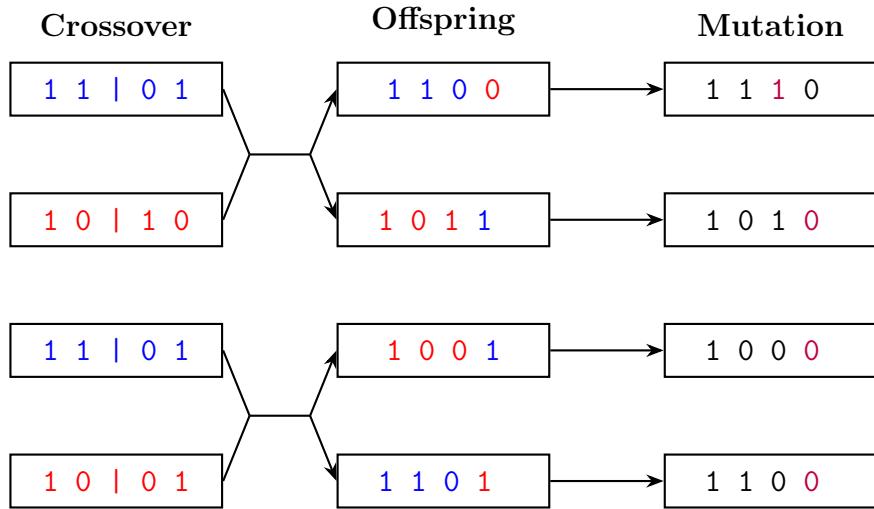
Select the top 3/4th of chromosomes based on their value.

```
best-so-far = [1101];
best-fitness-so-far = 53;
```

	Chromosomes	Fitness
Chromosome 1:	1 1 0 1	53
Chromosome 2:	1 0 1 0	46
Chromosome 4:	1 0 0 1	39

Step 4: Crossover

Perform crossover between the chromosome with the highest fitness value with the two other chromosomes from the selected chromosomes at the third bit.



Step 5: Mutation

Apply a mutation by flipping a random bit in each offspring.

Step 6: New Population

Replace the previous population with the set of new chromosomes.

New Population

Chromosome 1: **1 1 1 0**

Chromosome 2: **1 0 1 0**

Chromosome 3: **1 0 0 0**

Chromosome 4: **1 1 0 0**

Step 7: Repeat

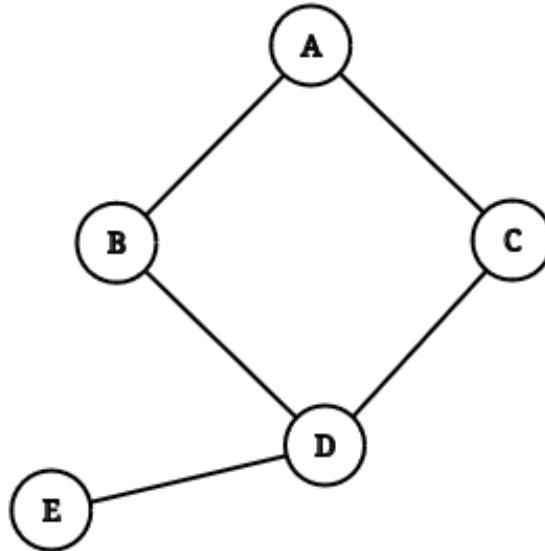
This process is repeated over several generations to find the chromosome with the highest fitness, representing the maximum value found.

4.9.4 Graph Coloring Problem:

The graph coloring problem involves assigning colors to the vertices of a graph such that no two adjacent vertices share the same color, and the goal is to minimize the number of colors used.

Graph Description

Consider a simple graph with 5 vertices (A, B, C, D, E) and the following edges: AB, AC, BD, CD, DE.



Genetic Algorithm Setup for Graph Coloring

Chromosome Representation

Each chromosome is an array where each position represents a vertex and the value at that position represents the color of that vertex. For simplicity, we'll use numerical values to represent different colors.

Initial Population

We randomly generate a small population of 4 solutions. We keep track of the best found route and its fitness.

```
best-so-far = [];
best-fitness-so-far = 0;
```

Initial Population

Chromosome 1: 1 2 3 1 2

Chromosome 2: 2 3 1 2 3

Chromosome 3: 1 2 1 3 2

Chromosome 4: 3 1 2 3 3

Fitness Function

Fitness is determined by the number of properly colored edges (i.e., edges connecting vertices of different colors). The maximum fitness for this graph is 5 (one for each edge).

Calculate the fitness based on the coloring rules.

Chromosome 1: Fitness = 5 (all edges correctly colored).

Chromosome 2: Fitness = 5.

Chromosome 3: Fitness = 4 (CD edge is incorrectly colored).

Chromosome 4: Fitness = 4 (DE edge is incorrectly colored).

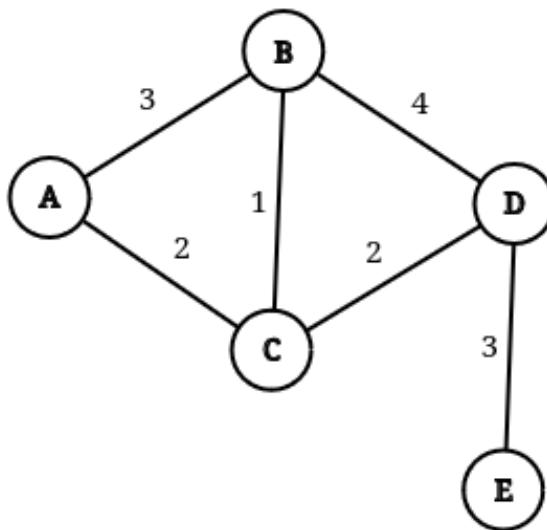
Carry on selection, crossover, mutation and population replacement as before.

4.9.5 Max-cut Problem

The Max-Cut problem is a classic problem in computer science and optimization in which the goal is to divide the vertices of a graph into two disjoint subsets such that the number of edges between the two subsets is maximized. Here, I will outline how to solve this problem using a genetic algorithm (GA).

Problem Setup

Consider a graph with 5 vertices (A, B, C, D, E) and the following edges with given weights: AB = 3, AC = 2, BC = 1, BD = 4, CD = 2, DE = 3



The goal is to find a division of these vertices into two sets that maximizes the sum of the weights of the edges that have endpoints in each set.

Genetic Algorithm Setup for Max-Cut

Chromosome Representation:

Each chromosome is a string of bits where each bit represents a vertex. A bit of '0' might represent the vertex being in set X and '1' in set Y.

Fitness Function

Fitness is determined by the sum of the weights of the edges between the two sets. For a chromosome, calculate the sum of weights for edges where one endpoint is '0' and the other is '1'.

For example, consider a configuration 10101 where the Vertices A, C, E in set Y; B, D in set X.

$$\text{Fitness} = \text{Sum of Edges (AB, BC, CD, DE)} - \text{Sum of Edges(AC)} = 3 + 1 + 2 + 3 - 2 = 7.$$

Carry on selection, crossover, mutation, and replacement as before.

4.10 Application of Genetic Algorithm in Machine Learning

4.10.1 Problem Setup: Feature Selection for Predictive Modeling

Suppose you're working with a medical dataset aimed at predicting the likelihood of patients developing a certain disease. The dataset contains hundreds of features, including patient demographics, laboratory results, and clinical parameters. Not all of these features are relevant for the prediction task, and some may introduce noise or redundancy.

Genetic Algorithm Setup for Feature Selection

Chromosome Representation:

Each chromosome in the GA represents a possible solution to the feature selection problem. Specifically, a chromosome can be encoded as a binary string where each bit represents the presence (1) or absence (0) of a corresponding feature in the dataset.

Initial Population

Generate an initial population of chromosomes randomly, where each chromosome has a different combination of features selected (1s) and not selected (0s).

Fitness Function:

The fitness of each chromosome (feature subset) is determined by the performance of a predictive model trained using only the selected features. Common performance metrics include accuracy, area under the ROC curve, or F1-score, depending on the problem specifics. Optionally, the fitness function can also penalize the number of features to maintain model simplicity.

4.10.2 Genetic Algorithm Process for Feature Selection

Step 1: Initialization

- Generate an initial population of feature subsets encoded as binary strings.

Step 2: Evaluation

- For each chromosome, train a model using only the features selected by that chromosome. Evaluate the model's performance on a validation set.

Step 3: Selection

- Select chromosomes for reproduction. Techniques like tournament selection or roulette wheel selection can be used, where chromosomes with higher fitness have a higher probability of being selected.

Step 4: Crossover

- Perform crossover between pairs of selected chromosomes to create offspring. A common method is one-point or two-point crossover, where segments of parent chromosomes are swapped to produce new feature subsets.

Step 5: Mutation

- Apply mutation to the offspring with a small probability. This could involve flipping some bits from 0 to 1 or vice versa, thus adding or removing features from the subset.

Step 6: Replacement

- Form a new generation by replacing some of the less fit chromosomes in the population with the new offspring. This could be a generational replacement or a steady-state replacement (where only the worst are replaced).

Iteration

- Repeat the process for a number of generations or until a stopping criterion is met (such as no improvement in fitness for a certain number of generations).

Example Use Case: Feature Selection in Medical Diagnosis

You have a dataset with 200 features from various medical tests. You apply a genetic algorithm with an initial population of 50 chromosomes, evolving over 100 generations. Each chromosome dictates which features are used to train a logistic regression model to predict disease occurrence.

The process might reveal that only 30 out of the 200 features significantly contribute to the prediction, eliminating redundant and irrelevant features and thus simplifying the model without sacrificing (or possibly even improving) its performance.

4.10.3 Problem Setup: Hyperparameter Optimization for a Neural Network

Suppose we are developing a neural network to classify images into categories (e.g., for a fashion item classification task). The performance of the neural network can depend heavily on the choice of various hyperparameters such as the number of layers, number of neurons in each layer, learning rate, dropout rate, and activation function.

Genetic Algorithm Setup for Hyperparameter Optimization

Chromosome Representation:

Each chromosome represents a set of hyperparameters for the neural network. For instance:

- Number of layers (e.g., 2-5)

- Neurons in each layer (e.g., 64, 128, 256, 512)
- Learning rate (e.g., 0.001, 0.01, 0.1)
- Dropout rate (e.g., 0.1, 0.2, 0.3)
- Activation function (e.g., relu, sigmoid, tanh)

Initial Population:

Generate an initial population of chromosomes, each encoding a different combination of these hyperparameters.

Fitness Function:

The fitness of each chromosome is evaluated based on the validation accuracy of the neural network configured with the hyperparameters encoded by the chromosome. Optionally, the fitness function can also include terms to penalize overfitting or excessively complex models.

Genetic Algorithm Process for Hyperparameter Optimization

Step 1: Initialization

- Generate an initial population of random but valid hyperparameter sets.

Step 2: Evaluation

- For each chromosome, construct a neural network with the specified hyperparameters, train it on the training data, and then evaluate it on a validation set. The validation set accuracy serves as the fitness score.

Step 3: Selection

- Select chromosomes for reproduction based on their fitness. High-fitness chromosomes have a higher probability of being selected. Techniques like tournament selection or rank-based selection are commonly used.

Step 4: Crossover

- Perform crossover operations between selected pairs of chromosomes to create offspring. Crossover can be one-point, two-point, or uniform (where each gene has an independent probability of coming from either parent).

Step 5: Mutation

- Apply mutation to the offspring chromosomes at a low probability. Mutation might involve changing one of the hyperparameters to another value within its range (e.g., changing the learning rate from 0.01 to 0.001).

Step 6: Replacement

- Replace the least fit chromosomes in the population with the new offspring, or use other replacement strategies like elitism where some of the best individuals from the old population are carried over to the new population.

Iteration

- Repeat the evaluation, selection, crossover, and mutation steps for several generations until the performance converges or a maximum number of generations is reached.

Example Use Case: Image Classification

You are using a dataset of fashion items where the task is to classify images into categories like shirts, shoes, pants, etc. You apply a genetic algorithm to optimize the hyperparameters of a convolutional neural network (CNN). After several generations, the GA might converge to an optimal set of hyperparameters that gives the highest accuracy on the validation dataset.

For instance, the best solution found by the GA could be:

- Number of layers: 3
- Neurons per layer: [512, 256, 128]
- Learning rate: 0.01
- Dropout rate: 0.2
- Activation function: relu

4.10.4 Genetic Algorithm in AI Games:

Example Use Case: Developing a Chess AI

Imagine you're developing an AI for a chess game. You start with 50 different strategies encoded as chromosomes. Each strategy is evaluated based on its performance in 100 games against diverse opponents. The strategies are then evolved over 100 generations, with each generation involving selection, crossover, and mutation to develop more refined and successful game strategies.

By the end of these iterations, the genetic algorithm might produce a strategy that effectively balances aggressive and defensive play, adapts to different opponent moves, and optimizes piece positioning throughout the game.

4.10.5 Genetic algorithms in Finance

Example Use Case: Diversified Investment Portfolio

Assume you manage an investment fund that considers a diverse set of assets, including stocks from various sectors, government and corporate bonds, and commodities like gold and oil. The task is to determine how much to invest in each asset class.

Assets: Stocks (technology, healthcare, finance), government bonds, corporate bonds, gold, oil.

Objective: Maximize the Sharpe ratio, considering historical returns and volatility data for each asset class.

Problem Setup:

Portfolio Optimization for Investment Management.

Suppose you are an investment manager looking to create a diversified investment portfolio. You want to determine the optimal allocation of funds across a set of available assets (e.g., stocks, bonds, commodities) to maximize returns while controlling risk, subject to various constraints like budget limits or maximum exposure to certain asset types.

Genetic Algorithm Setup for Portfolio Optimization

Chromosome Representation:

Each chromosome in the GA represents a potential portfolio, where each gene corresponds to the proportion of the total investment allocated to a specific asset.

Initial Population:

Generate an initial population of chromosomes, each encoding a different allocation strategy, ensuring that each portfolio adheres to the budget constraint (i.e., the total allocation sums to 100%).

Fitness Function:

The fitness of each chromosome (portfolio) is typically evaluated based on its expected return and risk (often quantified as variance or standard deviation). A common approach to measure fitness is to use the Sharpe ratio, which is the ratio of the excess expected return of the portfolio over the risk-free rate, divided by the standard deviation of the portfolio returns.

After running the genetic algorithm for several generations, the GA might find an optimal portfolio that, for example, allocates 20% to technology stocks, 15% to healthcare stocks, 10% to finance stocks, 20% to government bonds, 15% to corporate bonds, 10% to gold, and

10% to oil. This portfolio would have the highest Sharpe ratio found within the constraints set by the algorithm.

CHAPTER 5

LECTURE 6: ADVERSARIAL SEARCH/GAMES

Note: This lecture closely follows (sometimes, directly borrowed from) Chapter 5 to 5.2.4 of Russel and Norvig, *Artificial Intelligence: A Modern Approach*.

5.1 Introduction

Now we would like to focus on competitive environments, in which two or more agents have conflicting goals. The agents in these environments work against each other, minimizing the evaluation value of each other. To tackle such environments, we use Adversarial search. We will study Minimax search, using which we can find the optimal move for an agent in the restricted game environment. We will also study alpha-beta pruning where including pruning we can make the search more efficient by ignoring portions that do not help us find the optimal solution.

Instead of handling real life problems which have a lot of uncertainty and are difficult to handle we focus on games. Examples of such games are chess, Go, poker etc. We will use much easier games than these to demonstrate the adversarial nature.

5.1.1 Two player zero-sum games

- **Deterministic:** We can determine the outcome of an action.
- **Two players:** Only two agents working against each other. Example: Player A and Player B.
- **Turn taking:** Agents take turns alternatively. Example: After Agent (Player) A plays a move (chooses their action), Agent (Player) B will get to play their move (choose their action) and so on.
- **Perfect Information:** Fully observable environment, as in there are no parts of the game unknown to the agents.

- **Zero-Sum Games:** The total sum of points of the players in the game is zero. If one player wins (+1 points) that would mean the opponent player loses (-1 points). The total sum of the points in the game will be zero. When there is a draw, the point for both players is zero, meaning the sum will be zero.

We will use the term **move** as a synonym for **action**, and **position** as a synonym for **state**.

5.1.2 Maximizer and Minimizer

Continuing from the concept of a zero-sum game, let us refer to the two players (agents) of the game as MAX and MIN.

MAX aims to maximize its own game points. In a zero-sum game, the total payoff for all players is constant, meaning one player's gain is exactly equal to another's loss. Therefore, MAX seeks to achieve the highest possible outcome for itself, knowing that this will occur at the expense of the other player.

MAX evaluates all possible moves by considering what its opponent might choose, countering their moves. It chooses the option that leads to the highest evaluated payoff under optimal play.

On the other hand, MIN aims to win by minimizing MAX's points. While choosing its move, MIN selects the move that will reduce the good options for MAX, limiting MAX's game points and forcing it towards an outcome where MAX's game points are minimal. This strategy minimizes MIN's loss as well.

For example, in a simple game like tic-tac-toe, MAX aims to align three of its marks in a row, thereby maximizing its chance of winning and ensuring the other player's loss. Conversely, MIN attempts to block MAX's efforts to align three marks, while also seeking opportunities to create a threat that MAX must respond to, thereby steering the game towards a draw or a win for MIN.

MAX and MIN alternatively choose moves, each striving to reach their respective goals. The minimax algorithm provides a way to determine the best possible move for a player (assuming both players play optimally) by minimizing the possible loss for a worst-case (maximum loss) scenario.

These concepts are crucial in designing agents capable of thinking several steps ahead, anticipating and countering the opponent's strategies effectively.

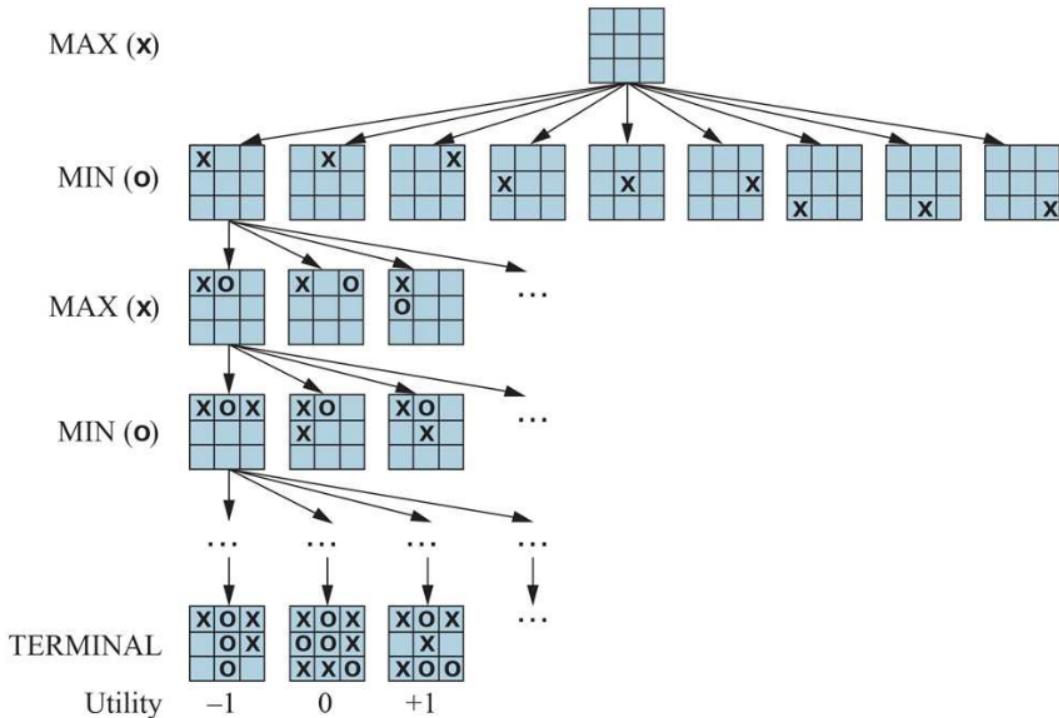
The following elements formally define the game:

- **S_0 :** The initial state, which specifies how the game is set up at the start.
- **TO-MOVE(s):** The player whose turn it is to move in state s .
- **ACTIONS(s):** The set of legal moves in state s .

- **RESULT(s, a):** The transition model, specifying the state that results from executing action a in state s .
- **IS-TERMINAL:** A test to determine if state s is terminal, returning true if the game has concluded and false otherwise. States where the game has ended are called terminal states. Example, Win, Lose, Draw.
- **UTILITY:** A utility function (also known as an objective function or payoff function), which provides the final numerical value to player p when the game concludes in a terminal state s . For instance, in chess, the result can be a win, loss, or draw, with values 1, 0, or $\frac{1}{2}$ respectively. In some games, the payoff range can be broader, such as in backgammon, where it ranges from 0 to 192.

The initial state, the action function, and the result function together define the state space graph. In this graph, the vertices represent states and the edges represent moves. A state can be reached by multiple paths. From the search tree, we can construct the game tree that traces every sequence of moves all the way to the terminal state.

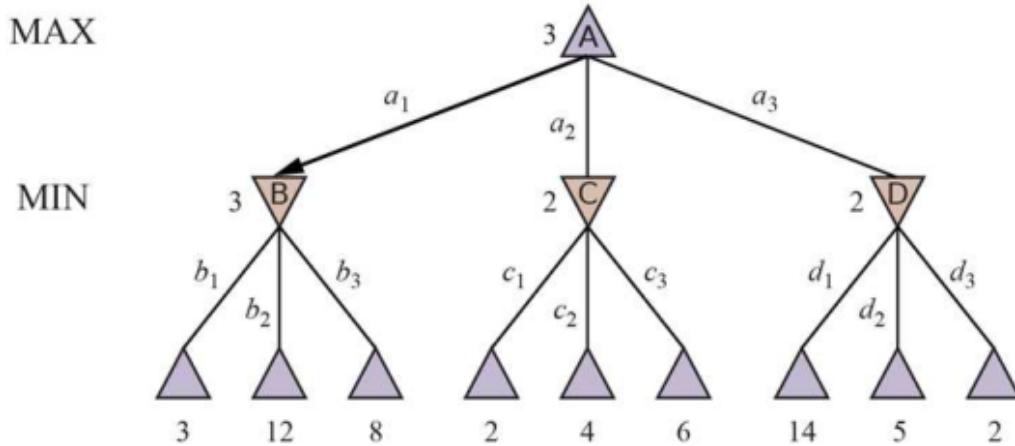
The figure below illustrates part of the game tree for tic-tac-toe (noughts and crosses). Starting from the initial state, MAX has nine possible moves. Players alternate turns, with MAX placing an X and MIN placing an O, until reaching terminal states where either one player has three in a row or the board is completely filled. The number on each leaf node represents the utility value of the terminal state from MAX's perspective; higher values benefit MAX and disadvantage MIN (hence the players' names).



A (partial) game tree for the game of tic-tac-toe. The top node is the initial state, and **MAX** moves first, placing an **x** in an empty square. We show part of the tree, giving alternating moves by **MIN (o)** and **MAX (x)**, until we eventually reach terminal states, which can be assigned utilities according to the rules of the game.

For tic-tac-toe, the game tree is relatively small—fewer than $9! = 362,880$ terminal nodes (with only 5,478 distinct states). But for chess, there are over 10^{40} nodes, so the game tree is best thought of as a theoretical construct that we cannot realize in the physical world.

5.2 Making optimal decisions using Minimax Search



A two-ply game tree. The Δ nodes are “MAX nodes,” in which it is MAX’s turn to move, and the ∇ nodes are “MIN nodes.” The terminal nodes show the utility values for MAX; the other nodes are labeled with their minimax values. MAX’s best move at the root is a_1 , because it leads to the state with the highest minimax value, and MIN’s best reply is b_1 , because it leads to the state with the lowest minimax value.

Consider the trivial game in Figure. The possible moves for MAX at the root node are labeled a_1, a_2, a_3 , and so on. The possible replies to for MIN are b_1, b_2, b_3 , and so on. This particular game ends after one move each by MAX and MIN.

NOTE: In some games, the word “move” means that both players have taken an action; therefore the word *ply* is used to unambiguously mean one move by one player, bringing us one level deeper in the game tree. The utilities of the terminal states in this game range from 2 to 14.

Given a game tree, the optimal strategy can be determined by working out the **minimax value** of each state in the tree, which we write as $\text{MINIMAX}(s)$. The minimax value is the utility (for MAX) of being in that state, *assuming that both players play optimally* from there to the end of the game. The minimax value of a terminal state is just its utility. In a non-terminal state, MAX prefers to move to a state of maximum value when it is MAX’S turn to move, and MIN prefers a state of minimum value (that is, minimum value for MAX and thus maximum value for MIN). So we have:

$$\text{MINIMAX}(s) = \begin{cases} \text{UTILITY}(s, \text{MAX}) & \text{if Is-TERMINAL}(s) \\ \max_{a \in \text{Actions}} \text{MINIMAX}(\text{RESULT}(s, a)) & \text{if To-MOVE}(s) = \text{MAX} \\ \min_{a \in \text{Actions}} \text{MINIMAX}(\text{RESULT}(s, a)) & \text{if To-MOVE}(s) = \text{MIN} \end{cases}$$

5.3 Minimax algorithm

The minimax algorithm explores the game tree using a depth-first search. If the tree’s maximum depth is m and there are b legal moves at each point, the time complexity is $\mathcal{O}(b^m)$

and the space complexity is $\mathcal{O}(bm)$.

The exponential complexity makes minimax impractical for complex games. For instance, chess has a branching factor of about 35, and the average game has a depth of about 80 plies. This means we would need to search almost 10^{123} states, which is not feasible.

Algorithm 6 MINIMAX-SEARCH

```

function MINIMAX-SEARCH(game, state) returns an action
    player  $\leftarrow$  game.To-MOVE(state)
    value, move  $\leftarrow$  MAX-VALUE(game, state)
    return move

function MAX-VALUE(game, state) returns a (utility,move) pair
    if game.IS-TERMINAL(state) then return game.UTILITY(state,player),null
     $v \leftarrow -\infty$ 
    for each a in game. ACTIONS(state) do
         $v2, a2 \leftarrow$  MIN-VALUE(game, game.RESULT(state,a))
        if  $v2 > v$  then
             $v, move \leftarrow v2, a$ 
     $v, move$ 

function MIN-VALUE(game, state) returns a (utility, move) pair
    if game.IS-TERMINAL(state) then return game.UTILITY(state,player),null
     $v \leftarrow +\infty$ 
    for each a in game.ACTIONS(state) do
         $v2, a2 \leftarrow$  MAX-VALUE(game, game.RESULT(state,a))
        if  $v2 < v$  then
             $v, move \leftarrow v2, a$ 
    return v, move

```

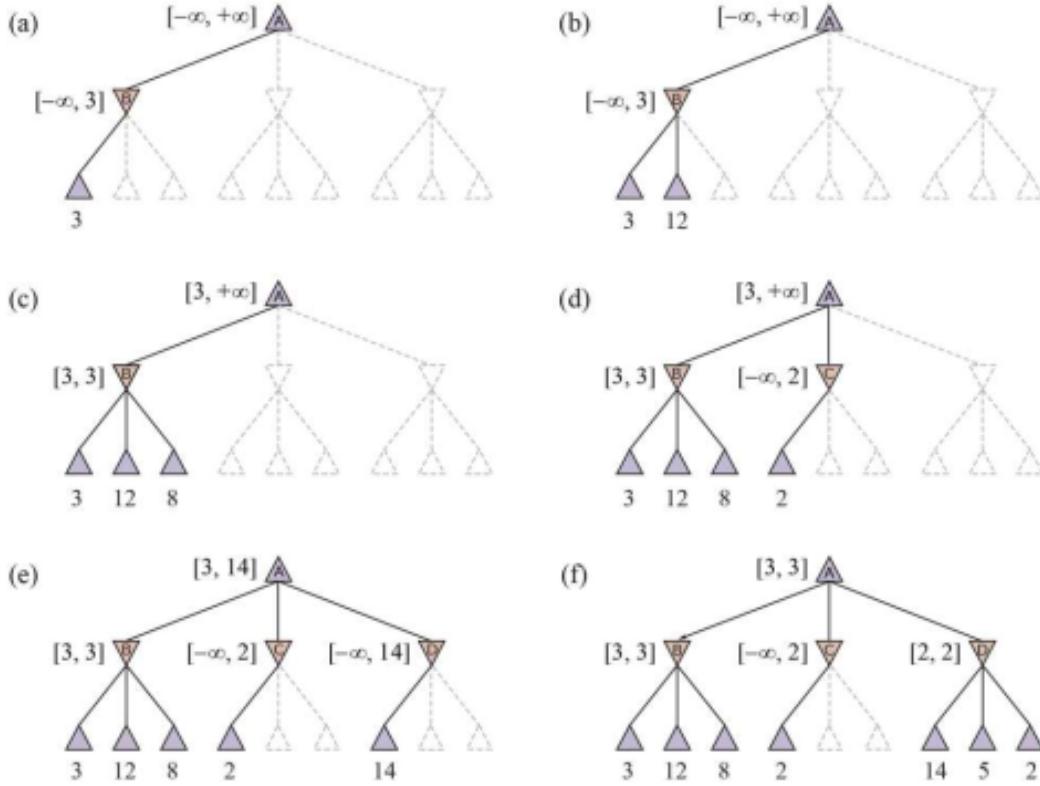
5.4 Alpha—Beta Pruning

Alpha-beta pruning is an optimization technique for the minimax algorithm, which is used in decision-making processes for two-player games like chess or tic-tac-toe. The primary goal of alpha-beta pruning is to reduce the number of nodes evaluated in the game tree, thereby improving the efficiency of the minimax algorithm.

In the minimax algorithm, every possible move and counter-move is explored to determine the optimal strategy. However, this exhaustive search can become computationally expensive, especially for complex games with large search spaces. Alpha-beta pruning addresses this issue by eliminating branches in the game tree that do not influence the final decision.

Let's revisit the two-ply game tree from the previous figure. This time, we will carefully

track what we know at each step. The steps are detailed in the figure below. The result is that we can determine the minimax decision without needing to evaluate two of the leaf nodes.



Stages in calculating the optimal decision for the game tree in the previous figure are shown below. At each point, we display the range of possible values for each node.

- The first leaf under B has a value of 3. Thus, B , a MIN node, has a value of at most 3.
- The second leaf under B has a value of 12. Since MIN would avoid this move, B 's value remains at most 3.
- The third leaf under B has a value of 8. After evaluating all of B 's successor states, B 's value is exactly 3. This means the root value is at least 3, because MAX can choose a move worth 3 at the root.
- The first leaf under C has a value of 2. Therefore, C , a MIN node, has a value of at most 2. Since B is worth 3, MAX would never choose C . Hence, there is no need to examine the other successors of C . This demonstrates alpha-beta pruning.
- The first leaf under D has a value of 14, making D worth at most 14. This is still higher than MAX's best alternative (i.e., 3), so we need to continue exploring D 's successors. Now, we also know the root's value is at most 14.

- (f) The second successor of D has a value of 5, so we continue exploring. The third successor is worth 2, so D 's exact value is 2. MAX's decision at the root is to move to B , yielding a value of 3.

Another way is to see this as a simplification of the MINIMAX formula. Consider the two unevaluated successors of node C in the figure above. Let these successors have values x and y . Then the value of the root node is given by

$$\begin{aligned}\text{MINIMAX}(\text{root}) &= \max(\min(3, 12, 8), \min(2, x, y), \min(14, 5, 2)) \\ &= \max(3, \min(2, x, y), 2) \\ &= \max(3, z, 2) \quad (\text{where } z = \min(2, x, y) \leq 2) \\ &= 3\end{aligned}$$

In other words, the value of the root and hence the minimax decision are *independent* of the values of the leaves x and y and therefore they can be pruned.

Algorithm 7 Alpha-beta Search Algorithm

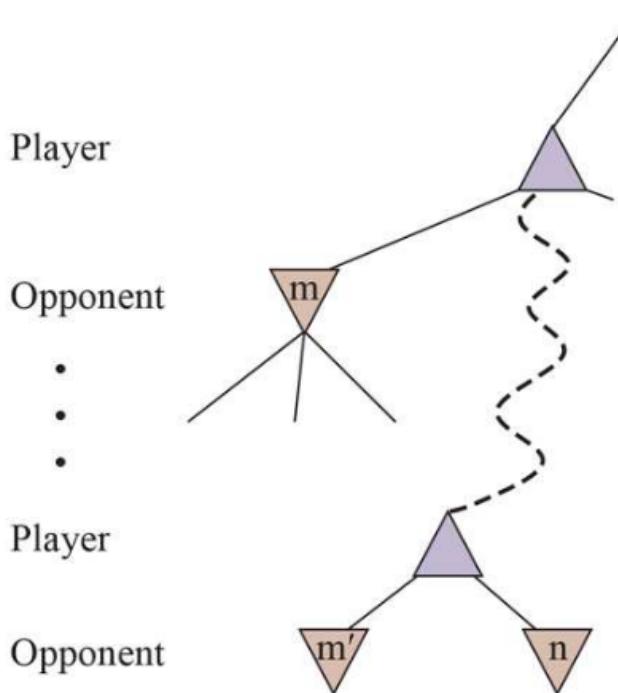
```

function ALPHA-BETA-SEARCH(game,state) returns an action
    player  $\leftarrow$  To-MOVE(state)
    value, move  $\leftarrow$  MAX-VALUE(game, state,  $-\infty$ ,  $+\infty$ )
    return move

function MAX-VALUE(game, state,  $\alpha$ ,  $\beta$ ) returns a (utility, move) pair
    if game.Is-TERMINAL(state) then return game.UTILITY(state, player), null
    v  $\leftarrow -\infty$ 
    for each a in game.ACTIONS(state) do
        v2, a2  $\leftarrow$  MIN-VALUE(game, game.RESULT(state, a),  $\alpha$ ,  $\beta$ )
        if v2  $>$  v then
            v, move  $\leftarrow$  v2, a
             $\alpha \leftarrow \text{MAX}(\alpha, v)$ 
        if v  $\geq \beta$  then return v, move
    return v, move

function MIN-VALUE(game, state,  $\alpha$ ,  $\beta$ ) returns a (utility, move) pair
    if game.Is-TERMINAL(state) then return game.UTILITY(state, player), null
    v  $\leftarrow +\infty$ 
    for each a in game.ACTIONS(state) do
        v2, a2  $\leftarrow$  MAX-VALUE(game, game.RESULT(state, a),  $\alpha$ ,  $\beta$ )
        if v2  $<$  v then
            v, move  $\leftarrow$  v2, a
             $\beta \leftarrow \text{MIN}(\beta, v)$ 
        if v  $\leq \alpha$  then return v, move
    return v, move
```

Alpha-beta pruning can be used on trees of any depth, and it often allows for pruning entire subtrees, not just leaves. The basic principle is this: consider a node n in the tree (see Figure below) where Player can choose to move to n . If Player has a better option either at the same level (e.g., m' in the figure below) or higher up in the tree (e.g., m in the figure below), Player will not move to n . Once we learn enough about n by examining some of its successor state to make this decision, we can prune n .



The general case for alpha–beta pruning. If m or m' is better than n for Player, we will never get to n in play.

Minimax search uses a depth-first approach, so we only consider the nodes along a single path in the tree at any given time. Alpha-beta pruning uses two additional parameters in $\text{MAX-VALUE}(state, \alpha, \beta)$, which set bounds on the backed-up values along the path.

α = the value of the best (i.e., highest-value) choice we have found so far at any choice point along the path for MAX. Think: α = “at least.”

β = the value of the best (i.e., lowest-value) choice we have found so far at any choice point along the path for MIN. Think: β = “at most”.

As the search progresses, these values are updated to reflect the highest and lowest scores that MAX and MIN can achieve, respectively. When a move is found that makes a branch less favorable than previously examined branches, that branch is pruned, meaning it is not explored further.

This pruning does not affect the final result of the minimax algorithm but significantly

reduces the number of nodes evaluated, making it more efficient. In the best case, alpha-beta pruning can reduce the time complexity from $\mathcal{O}(b^m)$ to $\mathcal{O}(b^{\frac{m}{2}})$ where b is the branching factor and m is the maximum depth of the tree.

Alpha-beta pruning enables more effective decision-making in complex games, allowing for deeper searches and better strategic planning within practical time limits.

5.4.1 Worst Case Scenario for Alpha-Beta Pruning

In the worst-case scenario of alpha-beta pruning, the algorithm does not effectively prune any branches, and it ends up exploring the entire game tree, similar to a standard minimax search. This situation can arise when the nodes are ordered in the least optimal way for pruning.

Key Points

- **Branch Ordering:** The worst case occurs when the best moves are always considered last. Alpha-beta pruning relies on evaluating the most promising moves first to maximize pruning. If the least promising moves are evaluated first, fewer branches are pruned.
- **No Pruning:** When the algorithm does not prune any branches, it evaluates every possible move at each depth level, leading to the maximum number of nodes being explored.
- **Time Complexity:** In the worst case, the time complexity of alpha-beta pruning is the same as that of the minimax algorithm without pruning, which is $\mathcal{O}(b^m)$. Here, b is the branching factor (the number of legal moves at each point), and m is the maximum depth of the tree.

Improving alpha-beta pruning

In complex games with large search spaces, to improve the efficiency and effectiveness of the alpha-beta pruning these following techniques can be used.

- **Move-Ordering:** Move-ordering prioritizes evaluating the most promising moves first, increasing the chances of effective pruning in alpha-beta search. This is usually achieved by sorting moves based on heuristic evaluations or previous search results, aiming to maximize the number of branches pruned.
- **Killer Moves:** Killer moves are specific moves that have previously caused significant pruning in similar positions and are tried early in the search to maximize pruning opportunities. These moves are stored and reused, assuming that moves which were effective in the past are likely to be effective again.
- **Iterative Deepening:** Iterative deepening involves repeatedly running depth-limited searches with increasing depth limits, combining the benefits of depth-first search (using

less memory) and breadth-first search (finding the optimal solution). It provides a way to use the best move from shallow searches to improve move-ordering in deeper searches.

- **Heuristic Alpha-Beta:** Heuristic alpha-beta uses evaluation functions to estimate the value of non-terminal nodes, guiding the search and pruning process more effectively based on heuristics. These evaluation functions consider factors like material balance, positional advantages, and other domain-specific knowledge to approximate the true minimax value of a position.

Part II

Modern AI

CHAPTER 1

MACHINE LEARNING BASICS

NOTE: The following chapters closely follows the textbook by Peter/Norvig and the previous slides of CSE422.

1.1 Introduction

Machine learning is a subfield of artificial intelligence (AI) that focuses on developing algorithms and models that enable computers to learn patterns from data and make decisions or predictions without being explicitly programmed for specific tasks.

In simple terms, machine learning allows systems to improve their performance over time by using past experiences (i.e., data) to generalize and adapt to new situations.

1.2 Why do we need machine learning

One may ask why we need machine learning, rather than simply programming machines to perform specific tasks.

- **Unpredictable situations** For some tasks it is impossible to know all possible scenarios beforehand.

Machine learning handles unpredictable situations by learning patterns from past data and making generalizations. Even if it has not seen the exact situation before, it can make a best guess based on similar examples.

Example: Email Spam Detection

- **Training data:** Emails labeled as "spam" or "not spam".
- **New email:** The system has never seen this exact message, but it recognizes similar keywords and sender behavior from past spam.
- **Action:** Marks it as spam.

Example: Stock Investment

- **Training data:** Historical stock prices, Market reactions to past political/economic crises, Sector-wise performance during volatility.
- **New situation:** A sudden political crisis causes uncertainty in the market which is never seen before.
- **Action** Machine learning can help predict likely drop or recovery patterns, identify safer investment sectors and can recommend adjustment of portfolio based on learned behavior from similar past events.
- **Unknown Solutions** Some tasks are easy for humans but very hard to explicitly program using traditional rules. These tasks often involve complex patterns, uncertainty, or high variability—things that machine learning handles well by learning from data rather than following fixed instructions.

Example: Face Recognition

- **Hard to program:** Human faces vary by angle, lighting, age, expression. Writing rules for all possibilities is nearly impossible.
- **Machine learning approach:** Trains on thousands of face images, learning features (like eyes, nose shape) automatically.

Example: Handwriting Recognition

- **Hard to program:** Every person writes differently, even the same letter can look very different.
- **Machine learning approach:** Learns from many handwritten samples to recognize letters based on strokes and curves.

1.3 Paradigms of Machine Learning

Machine learning includes several paradigms that guide how systems learn from data.

Rote learning is simple memorization, like a model storing exact inputs and outputs—for example, a chatbot recalling predefined replies.

Induction involves generalizing from examples; a spam filter learns patterns in labeled emails to classify new ones.

Clustering is an unsupervised method that groups similar data, such as organizing customers into segments based on purchase behavior.

Analogy and discovery use representation mapping—like a system recognizing that the relationship between Earth and Moon is similar to Jupiter and its moons.

Genetic algorithms apply evolutionary search, creating and evolving solutions over generations, useful in optimizing investment strategies, parameter tuning etc.

Finally, **reinforcement learning** is reward-based; for instance, a game-playing agent learns to win by receiving points and improving its moves over time.

1.3.1 Major Paradigms of Machine Learning

Based on learning feedbacks we can define three major paradigms of machine learning.

- **Supervised learning:** Uses labeled data to train models, where each input has a known output. For example, a model trained to detect spam emails based on labeled examples ("spam" or "not spam").
- **Unsupervised learning:** Works with data that has no labels, helping to find hidden patterns—like clustering customers based on buying behavior without knowing their categories.
- **Reinforcement learning:** Involves an agent learning by interacting with an environment and receiving rewards or penalties; for instance, a robot learning to walk by trial and error.

1.4 Supervised Learning

Supervised learning is a type of machine learning where a model is trained on a labeled dataset—that is, each input has a corresponding correct output. The goal is for the model to learn a mapping from inputs to outputs, so it can accurately predict outcomes for new, unseen data.

In supervised learning, the model learns from **input-output pairs**, also known as **feature-label pairs**.

- **Input (Features):** These are the measurable variables that form the input vector, often denoted by X .
 - *Example:* For predicting house prices, features may include:
 - * Size of the house (in sq. ft)
 - * Number of bedrooms
 - * Location rating
 - * Age of the house
- **Output (Label or Target):** This is the value or category to be predicted, usually denoted by Y .

- *Example:* The price of the house in dollars.

Each training example is represented as a pair: (X, Y) where X is the input vector containing features x_1, x_2, \dots and Y is the corresponding output label.

1.5 Steps of Supervised Learning

Supervised learning involves learning a function that maps inputs to known outputs. The general workflow includes the following steps:

- 1. Collect Data:** Gather labeled data where each input has a corresponding correct output.
- 2. Preprocess Data:** Clean the data, handle missing values, normalize features, and encode categorical variables if necessary.
- 3. Split the Data:** Divide the dataset into training and test sets (optionally a validation set).
- 4. Train the Model:** Use the training set to teach the model to learn the mapping from input to output.
- 5. Evaluate the Model:** Use the test (or validation) set to measure performance using appropriate metrics.

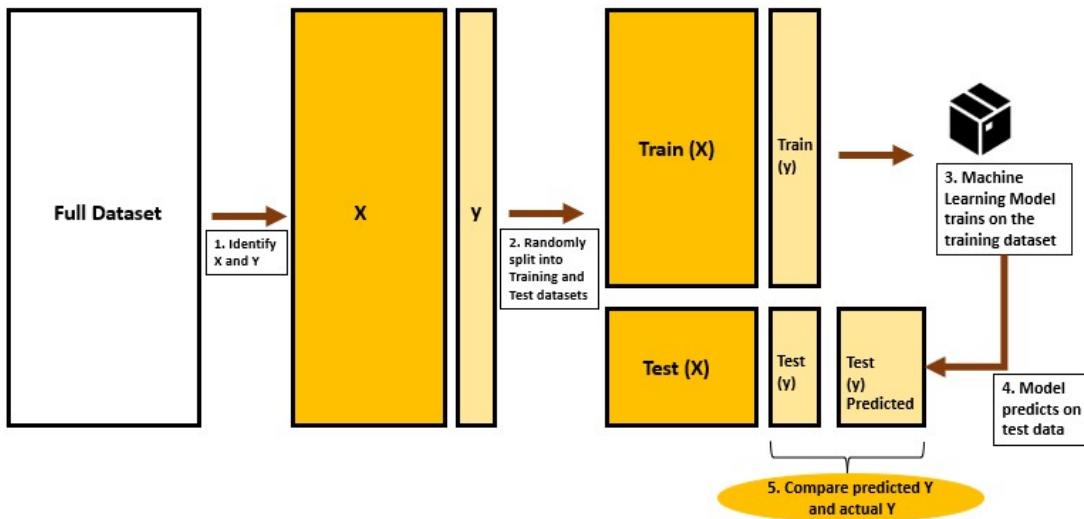


Figure 1.1: Training and evaluation of model

- 6. Tune Hyperparameters:** Adjust the model settings to optimize performance.
- 7. Deploy the Model:** Use the trained model to make predictions on new, unseen data.

1.6 Types of Supervised Learning

Supervised learning is mainly divided into two types:

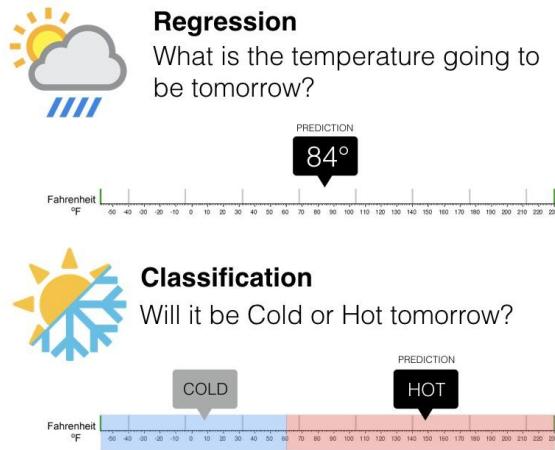


Figure 1.2: Types of supervised learning

- **Regression:**

- *Goal:* Predict continuous numeric values.
- *Example:* Predicting house prices based on size and location.

- **Classification:**

- *Goal:* Predict discrete class labels.
- *Example:* Classifying emails as “spam” or “not spam”.

1.7 Model Selection

Choosing the right model is crucial for achieving good performance. The model selection process involves:

1. **Understanding the Task:** Determine whether the problem is regression or classification.
2. **Choosing Candidate Models:** Based on task complexity, data size, and interpretability. Common models include:

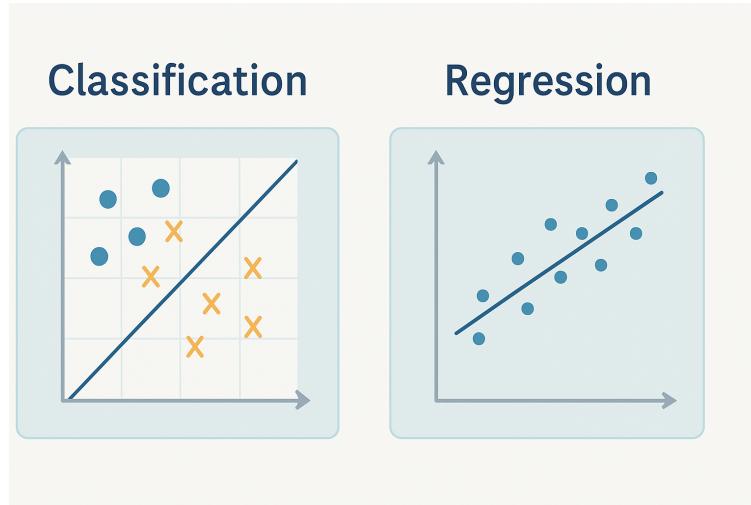


Figure 1.3: Classification Vs. Regression

- Linear/Logistic Regression
 - Decision Trees
 - k-Nearest Neighbors (k-NN)
 - Support Vector Machines (SVM)
 - Neural Networks
3. **Cross-Validation:** Evaluate models using validation techniques to estimate performance on unseen data.
 4. **Model Comparison:** Compare models based on accuracy, error rates, computation time, etc.
 5. **Select the Best Model:** Choose the model that balances performance and complexity.

1.8 Hypothesis in Supervised Learning

A **hypothesis** in supervised learning is a candidate function that maps inputs to outputs:

$$h : X \rightarrow Y$$

The goal is to find the best hypothesis h from a set called the **hypothesis space** \mathcal{H} , which includes all models that the learning algorithm can select from.

Example: In linear regression, a possible hypothesis could be:

$$h(x) = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

where w_1, w_2, \dots, w_n are weights, and b is the bias term.

1.9 Hypothesis Selection

To select the best hypothesis from \mathcal{H} , the following steps are followed:

1. **Define a loss function** – Measures the error between predicted and true outputs.
 - Example: Mean Squared Error (MSE) for regression, Cross-Entropy for classification.
2. **Train the model** – Minimize the loss function on training data using optimization techniques (e.g., gradient descent).
3. **Validate** – Use a separate validation set to test the generalization ability of the hypothesis.
4. **Select** – Choose the hypothesis with the best performance on the validation set.
5. **Test** – Evaluate the final hypothesis on unseen test data to estimate real-world performance.

Supervised learning aims to find the best hypothesis that accurately maps inputs to outputs by learning from labeled data.

Bias and Variance in Hypothesis Selection

When we select a hypothesis in machine learning, we aim for a model that generalizes well to unseen data. Two important sources of error are:

Bias

- Bias is the error due to overly simplistic assumptions in the model.
- A high-bias model may under-represent the data structure.
- Example: Using a linear model on highly nonlinear data.

Variance

- Variance is the error from excessive sensitivity to small fluctuations in the training data.
- A high-variance model captures noise along with the signal.
- Example: A high-degree polynomial fitting every point exactly.

Underfitting and Overfitting

Underfitting

- Caused by a model that is too simple to capture the underlying patterns (high bias). Example: 1.4 subplot a.
- Poor performance on both training and test data.

Overfitting

- Caused by a model that is too complex and fits the noise in the training data (high variance). Example: 1.4 subplot c and subplot d.
- Excellent performance on training data but poor generalization to new data.

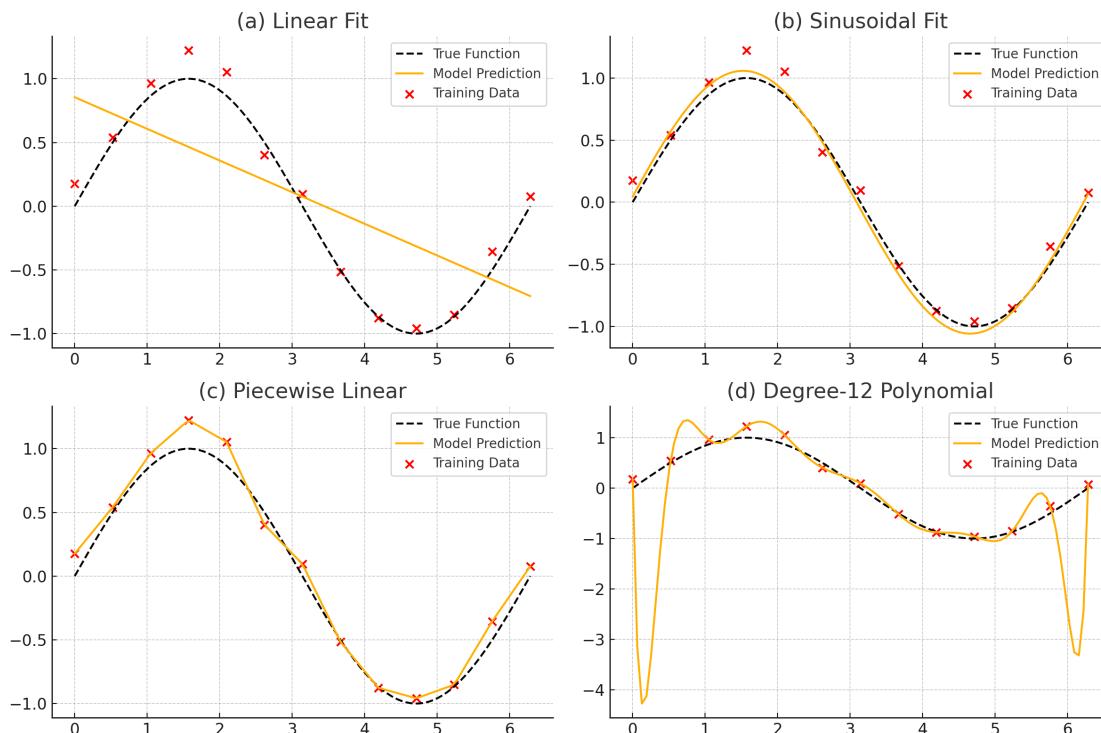


Figure 1.4: Bias-Variance Tradeoff Visualization

What is a Good Fit?

A good fit in machine learning is a model that captures the underlying structure of the data without memorizing the noise. It achieves a balance between underfitting and overfitting.

Characteristics of a Good Fit

- **Balanced bias and variance:** The model is neither too simple nor too complex.
- **Low generalization error:** It performs well on both the training set and unseen test data.
- **Appropriate complexity:** The model captures meaningful trends without reacting to random fluctuations.

Example

- A good fit curve for noisy sine-wave data 1.4 would follow the overall shape of the sine wave (subplot b) while ignoring random noise in individual data points.

A good fit is one that generalizes well — it learns the signal, not the noise.

CHAPTER 2

PROBABILITY THEORY

2.1 Probability Theory in AI

Probability theory is essential for AI systems to make decisions, predictions, and inferences under uncertainty. It underpins machine learning algorithms, from classification (e.g., Naive Bayes) to generative models (e.g., Hidden Markov Models and Variational Autoencoders).

In Machine Learning, probability theory is used in handling uncertainty, making predictions, learning from data, and updating models with observations of new data. It allows systems to make informed, data-driven decisions, even when the data is noisy or incomplete.

Models like Bayesian networks, Gaussian mixtures, and Markov models use probability theory to represent complex relationships and uncertainty in data. Techniques like Bayesian inference enable models to update predictions as new data is observed, refining the model over time.

2.2 Basic Concepts

Experiment: A process that results in one outcome from a set of possible outcomes. For example, tossing a coin or rolling a die.

Sample Space (S): The set of all possible outcomes of an experiment. For a fair coin, the sample space is

$$S = \{\text{Heads}, \text{Tails}\}.$$

Event (E): A subset of the sample space. An event may consist of one or more outcomes. For example, getting heads in a coin toss is an event.

Probability of an Event ($P(E)$): The measure of how likely an event is to occur. Probability values range from 0 (impossible event) to 1 (certain event).

Classical Theory of Probability:

$$P(E) = \frac{\text{Number of desired possible outcomes}}{\text{Number of all equally possible outcomes}}. \quad (2.1)$$

Assumption: All outcomes have equal possibility.

Example 1: Tossing a Coin

Experiment: Tossing a fair coin.

Sample Space (S): The set of all possible outcomes.

$$S = \{\text{Heads, Tails}\}.$$

Event (E): Getting "Tails" on the toss.

$$E = \{\text{Tails}\}.$$

Probability of an Event $P(E)$: The probability of getting "Tails" on the toss (assuming a fair coin).

$$P(E) = \frac{1}{2}.$$

Example 2: Rolling a Die

Experiment: Rolling a fair six-sided die.

Sample Space (S):

$$S = \{1, 2, 3, 4, 5, 6\}.$$

Event (E): Rolling a 3 or greater.

$$E = \{3, 4, 5, 6\}.$$

Probability of an Event $P(E)$: The probability of rolling a 3 or greater.

$$P(E) = \frac{4}{6} = \frac{2}{3}.$$

Statistical Theory of Probability: Run the experiment a large number of times. Create a table and collect data accordingly.

Example: Calculating the statistical the probability of getting a 3 or more if we throw a die.

Throw the die 600 times and then calculate the average number of times a 3 or greater value appears. To calculate the probability of rolling a 3 or greater using the statistical approach, we can simulate rolling a die 600 times. The results are shown in the table below:

Value	1	2	3	4	5	6
No. of times appeared	95	105	110	94	97	99

Table 2.1: Number of times each value appeared in 600 rolls of a die

$$P(E) = \frac{110 + 94 + 97 + 99}{600} = \frac{400}{600} = \frac{2}{3}.$$

Atomic Theory: For any event A , its probability will be

$$0 \leq P(A) \leq 1. \quad (2.2)$$

The sum of the probabilities of all possible events is 1:

$$\sum_{\forall} P(E_i) = 1. \quad (2.3)$$

2.3 Basic Probability Rules

Complement Rule: The probability that event E does not occur is:

$$P(E^c) = 1 - P(E) \quad (2.4)$$

where E^c denotes the complement of event E (i.e., all outcomes where E does not happen).

Addition Rule: For two events A and B , the probability that either A or B occurs is:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) \quad (2.5)$$

where $A \cap B$ is the event where both A and B occur.

Multiplication Rule: If two events A and B are independent (i.e., the outcome of one does not affect the other), the probability that both events occur is:

$$P(A \cap B) = P(A) \times P(B) \quad (2.6)$$

Conditional Probability: The probability of event A occurring given that event B has occurred is called *conditional probability* and is denoted by $P(A|B)$. It is calculated as:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad (\text{provided } P(B) > 0) \quad (2.7)$$

This concept helps in situations where we want to know how likely one event is, given that we have additional information about another event. Here, $P(B)$ acts as a normalizing constant.

Example: Conditional Probability

Suppose you are a teacher at a school and you have information about students passing or failing a math test. Out of 100 students, 40 students passed the test. Additionally, you know that 30 students who passed the test also attended an extra review session.

You want to find the probability that a student attended the extra review session given that they passed the test.

Solution: We want to calculate the **Conditional probability**

$$P(\text{Review}|\text{Pass}) = \frac{P(\text{Review and Pass})}{P(\text{Pass})}$$

From the given information we calculate,

$$P(\text{Pass}) = \frac{40}{100} = 0.4$$

$$P(\text{Review and Pass}) = \frac{30}{100} = 0.3$$

Thus,

$$P(\text{Review}|\text{Pass}) = \frac{.3}{.4} = .75$$

This means if a student has passed, there is a 75% chance that they have attended the review session.

Product Rule: From the conditional probability (2.7) we can get,

$$P(A \cap B) = P(A|B)P(B) \quad (2.8)$$

Chain Rule: We can further use the conditional probability to find the joint probability of multiple events, x_1, x_2, \dots, x_n .

$$P(x_1, x_2, \dots, x_n) = P(x_1)p(x_1|x_2)P(x_1|x_2, x_3) \dots \quad (2.9)$$

$$= \prod_i^n P(x_i|x_1, \dots, x_{i-1}) \quad (2.10)$$

2.4 Bayes' Rule: Derivation and Examples

2.4.1 Derivation of Bayes' Rule

Bayes' Rule allows us to update our beliefs based on new evidence. It is derived from the definition of conditional probability. Recall that the conditional probability of event A given event B is:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad (2.11)$$

Similarly, the conditional probability of event B given event A is:

$$P(B|A) = \frac{P(A \cap B)}{P(A)} \quad (2.12)$$

From these two definitions, we can express $P(A \cap B)$ as:

$$P(A \cap B) = P(A|B) \cdot P(B) = P(B|A) \cdot P(A) \quad (2.13)$$

Rearranging this to solve for $P(A|B)$, we get:

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} \quad (2.14)$$

This is the famous Bayes' Rule.

2.4.2 Significance of Bayes' Rule

Bayes' Rule is a powerful tool for updating the probability of an event based on new evidence. It is widely used in various fields, including weather prediction, medical diagnosis, and NLP tasks such as sentiment analysis.

Example: Weather prediction using Bayes' Rule

Consider a weather prediction scenario. We want to compute the probability that it will rain tomorrow, (A) given that the sky is cloudy, (B).

From observed data, we are given the following information:

- $P(A) = 0.30$ (the prior probability that it will rain),
- $P(B|A) = 0.80$ (the likelihood: the probability that the sky will be cloudy given that it rains),
- $P(B) = 0.60$ (the probability that the sky is cloudy, regardless of whether it rains or not).

We want to calculate $P(A|B)$ (Posterior Probability), the probability that it will rain given that the sky is cloudy.

Bayes' Rule tells us:

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} = \frac{0.80 \cdot 0.30}{0.60} = \frac{0.24}{0.60} = 0.40 \quad (2.15)$$

Thus, the probability that it will rain tomorrow given that the sky is cloudy is 0.40, or 40%.

Example of Sentiment Analysis using Bayes' Rule

Bayes' Rule can also be applied in Natural Language Processing (NLP), for instance, in sentiment analysis, where we want to classify a piece of text as either positive or negative.

For example, we want to classify a tweet as either positive or negative based on the presence of the word "good".

Let's define the following events:

- A = Positive sentiment
- B = The word "good" appears in the tweet

We want to calculate $P(A|B)$, the probability that a tweet has a positive sentiment, given that the word "good" appears.

From observed data, we know:

- $P(A) = 0.70$ (prior probability that the tweet is positive),

- $P(B|A) = 0.60$ (likelihood: the probability that the word "good" appears in a positive tweet),
- $P(B) = 0.50$ (probability that the word "good" appears in any tweet, regardless of sentiment).

Using Bayes' Rule, we calculate:

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} = \frac{0.60 \cdot 0.70}{0.50} = \frac{0.42}{0.50} = 0.84 \quad (2.16)$$

Thus, the probability that the tweet has a positive sentiment, given that the word "good" appears, is 0.84 or 84%.

2.5 Discrete Random Variables

A discrete random variable is a variable that can take on a countable number of distinct values. These values may be finite or infinite but are always countable. Discrete random variables typically arise in situations where the outcomes are distinct and can be listed, such as the number of heads in a series of coin tosses or the number of students passing an exam.

2.5.1 Notation for Discrete Random Variables

Random Variable

A discrete random variable is usually denoted by a capital letter, such as X , Y , or Z .

Possible Values

The possible values that a discrete random variable can take are denoted by lowercase letters, such as x_1, x_2, x_3, \dots

Probability Mass Function (PMF)

The probability that the random variable X takes a specific value x_i is represented as $P(X = x_i)$ or $p(x_i)$. The Probability Mass Function (PMF) gives the probability distribution of a discrete random variable.

Cumulative Distribution Function (CDF)

The cumulative probability up to a value x_i is represented as:

$$F(x_i) = P(X \leq x_i) \quad (2.17)$$

which sums the probabilities for all values less than or equal to x_i .

2.5.2 Properties of Discrete Random Variables

- **Non-Negativity:** For any x , $P(X = x) \geq 0$.
- **Normalization:** The sum of the probabilities for all possible values must equal to 1:

$$\sum_x P(X = x) = 1 \quad (2.18)$$

- **Countability:** The possible values x_1, x_2, x_3, \dots must be countable.

2.5.3 Basic Vector Notation

Vector of Random Variables

A set of n random variables can be represented as a vector. For example, a vector \mathbf{X} of random variables is often written as:

$$\mathbf{X} = (X_1, X_2, \dots, X_n) \quad (2.19)$$

where X_i are the individual random variables.

Vector of Possible Values: A random variable X can take a set of discrete values, which can also be represented as a vector of values. For example:

$$\mathbf{x} = (x_1, x_2, \dots, x_m) \quad (2.20)$$

where x_i represents a specific outcome or value of the random variable.

Probability Mass Function (PMF)

The PMF can be represented as a vector of probabilities. If X is a discrete random variable taking values x_1, x_2, \dots, x_m , the PMF is a vector of probabilities:

$$\mathbf{p} = (p(x_1), p(x_2), \dots, p(x_m)) \quad (2.21)$$

where $p(x_i) = P(X = x_i)$ is the probability that X takes the value x_i .

Example: Representation using Discrete Random Variable

Let X be a discrete random variable representing the outcome of rolling a fair 6-sided die. The possible outcomes for X are the integers from 1 to 6. The vector notation representing the values the discrete random variable can take is:

$$\mathbf{x} = (1, 2, 3, 4, 5, 6) \quad (2.22)$$

Since the die is fair, the probability of each outcome is equal. Thus, the probability of each outcome is:

$$P(X = x) = \frac{1}{6} \quad \text{for each } x \in \{1, 2, 3, 4, 5, 6\} \quad (2.23)$$

We can represent the **PMF** as a vector of probabilities:

$$\mathbf{P} = \left(\frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6} \right) \quad (2.24)$$

Cumulative Distribution Function (CDF)

The CDF gives the cumulative probability up to each value of X . The CDF vector is denoted as:

$$\mathbf{F} = (F(1), F(2), F(3), F(4), F(5), F(6)) \quad (2.25)$$

2.5.4 Discrete Random Variable with Categorical Outcomes

A discrete random variable can take categorical outcomes, where the outcomes are not numeric. For example, consider a random variable X representing the weather, where the possible outcomes are {Sunny, Cloudy, Rainy}.

Example: PMF and CDF for Categorical Weather Prediction

Let X be the random variable representing the weather on a given day. The possible outcomes are:

$$X = \{\text{Sunny}, \text{Cloudy}, \text{Rainy}\}$$

The Probability Mass Function (PMF) for X is given by:

$$P(X = \text{Sunny}) = 0.5, \quad P(X = \text{Cloudy}) = 0.3, \quad P(X = \text{Rainy}) = 0.2$$

Cumulative Distribution Function (CDF)

The CDF for X is calculated as the cumulative probability up to each outcome. Since the outcomes are ordered as "Sunny", "Cloudy", and "Rainy", the CDF is:

$$\begin{aligned} F(\text{Sunny}) &= P(X \leq \text{Sunny}) \\ &= P(X = \text{Sunny}) = 0.5 \\ F(\text{Cloudy}) &= P(X \leq \text{Cloudy}) \\ &= P(X = \text{Sunny}) + P(X = \text{Cloudy}) \\ &= 0.5 + 0.3 = 0.8 \\ F(\text{Rainy}) &= P(X \leq \text{Rainy}) \\ &= P(X = \text{Sunny}) + P(X = \text{Cloudy}) + P(X = \text{Rainy}) \\ &= 0.5 + 0.3 + 0.2 = 1 \end{aligned}$$

Thus, the CDF for X is:

$$F(X) = \begin{cases} 0.5 & \text{if } X = \text{Sunny} \\ 0.8 & \text{if } X = \text{Cloudy} \\ 1 & \text{if } X = \text{Rainy} \end{cases} \quad (2.26)$$

2.6 Data to Probability

Example: A survey is conducted on 500 people both male and female about which tv shows they liked. From the 500 people we got the following response.

Category	Male	Female	Total
GOT	80	120	200
TBBT	100	25	125
Others	50	125	175
Total	230	270	500

Table 2.2: Survey Data for TV Shows Viewership by Gender

Category	Male	Female	Total
GOT	.16	.24	.40
TBBT	.2	.05	.25
Others	.1	.25	.35
Total	.46	.54	1

Table 2.3: Probability distribution from Survey Data for TV Show Viewership by Gender

2.6.1 Marginal Probability

Marginal probability, also known as **Simple Probability** represents the probability of a particular event.

The marginal probability is obtained by summing the joint probabilities across the rows or columns.

Example: Computation using Marginal Probability

For the joint probability distribution given in 2.3 find the probability of a viewer watching TBBT.

Solution: Sum over all the values of row TBBT.

$$P(\text{TBBT}) = .2 + .05 = .25.$$

This is the value on the cell for TBBT on the Total column.

2.6.2 Marginal Probability Distribution

All the probability that occurs in the margin for a particular variable/event can be summed up to create Marginal Probability Distribution. For example, the "Total" column and "Total" row in 2.3.

2.6.3 Joint Probability

Joint Probability is the probability of two or more events occurring at the same time. For two events, this is $P(A \cap B)$. We can easily calculate it from the probability distribution.

Example: Computation Using Joint Probability Distribution

What is the joint probability of a person being Female and liking TBBT? **Solution:** We can find it directly from the probability distribution table.

$$P(T^F) = 0.05$$

2.6.4 Joint Probability Distribution

All the probability that occur jointly can be summed up to create Joint Probability Distribution (JPD). Sum of JPD is 1.

Category	Male	Female	Total
GOT	0.16	0.24	0.40
TBBT	0.20	0.05	0.25
Others	0.10	0.25	0.35
Total	0.46	0.54	1.00

Table 2.4: Joint Probability Distribution Table is highlighted

2.7 Understanding Conditional Probability using Probability Distribution Table

We can calculate conditional probability using the joint distribution table.

Example: What is the probability of a person liking GOT given that person is male?

Ans: Using values directly from table 2.3

$$\begin{aligned} P(G|M) &= \frac{P(G \cap M)}{P(M)} \\ &= \frac{0.16}{0.46} \\ &= 0.347 \end{aligned}$$

2.7.1 Absolute independence

When two random variables (events) A and B are **independent** we can write:

$$P(A \cap B) = P(A)P(B); \quad (2.27)$$

equivalently,

$$P(A) = P(A|B) \text{ and } P(B) = P(B|A) \quad (2.28)$$

2.7.2 Conditional Independence

Two events (or random variables) A and B are said to be **conditionally independent** given another event C, if:

$$P(A \cap B | C) = P(A | C) \cdot P(B | C) \quad (2.29)$$

This is denoted by:

$$A \perp B | C \quad (2.30)$$

It means: "*A is independent of B given C.*"

Intuition

Without conditioning on C, the variables A and B may be dependent. However, once C is known, learning about A gives no extra information about B, and vice versa. This means

$$P(A|B, C) = P(A|C) \text{ and } P(B|A, C) = P(B|C) \quad (2.31)$$

Now,

$$P(A, B, C) = P(A|B, C)P(B|C)P(C) \text{ using chain rule.} \quad (2.32)$$

$$\frac{P(A, B, C)}{P(C)} = \frac{P(A|C)P(B|C)P(C)}{P(C)} \text{ using 2.31} \quad (2.33)$$

$$P(A, B|C) = P(A|C)P(B|C) \text{ using conditional rule} \quad (2.34)$$

*Example: Conditional Independence Let:

- A : Person has a cough
- B : Person has a fever
- C : Person has the flu

In general, coughing and fever are correlated. But if we know the person has the flu, the presence of a cough does not tell us more about whether the person has a fever — both symptoms are explained by the flu.

$$P(\text{Cough} \cap \text{Fever} \mid \text{Flu}) = P(\text{Cough} \mid \text{Flu}) \cdot P(\text{Fever} \mid \text{Flu}) \quad (2.35)$$

Thus, Cough and Fever are conditionally independent given Flu.

Applications

Conditional independence is a fundamental concept in:

- **Naive Bayes Classifier**: Assumes features are conditionally independent given the class.
- **Bayesian Networks**: Graphical models representing dependencies via conditional independence.
- **Causal Inference**: To determine the effect of interventions and treatments.

Example: Checking Independence

Are Male viewers and GOT independent?

Ans:

$$\begin{aligned} P(M \cap \text{GOT}) &= 0.16 \\ P(M) &= 0.46 \\ P(\text{GOT}) &= 0.40 \\ P(M) * P(\text{GOT}) &= 0.46 * 0.40 = 0.184 \end{aligned}$$

Since, $P(M \cap \text{GOT}) \neq P(M) * P(\text{GOT})$ so, not independent.

$p(\text{smart} \wedge \text{study} \wedge \text{prep})$	smart		$\neg\text{smart}$	
	study	$\neg\text{study}$	study	$\neg\text{study}$
prepared	0.432	0.160	0.084	0.008
$\neg\text{prepared}$	0.048	0.160	0.036	0.072

Table 2.5: Probability distribution for Student personality and preparation

- Is smart conditionally independent of prepared, given study?
- Is study conditionally independent of prepared, given smart?
- Is smart conditionally independent of prepared, given study?
- Is study conditionally independent of prepared, given smart?

2.7.3 Example: Checking Conditional Independence

Is smart conditionally independent of prepared, given study?

$$\begin{aligned} P(Sm \cap Pr | Sd) &= P(Sm | Sd) * P(Pr | Sd) \\ P(Sm \cap Pr | Sd) &= P(Sm \cap Pr \cap Sd) / P(Sd) \\ &= 0.432 / 0.6 \\ &= 0.72 \end{aligned}$$

2.8 Law of Total Probability

The law of total probability allows us to compute the probability of an event by considering all possible conditions (or outcomes) of a related variable.

It states that the probability of an event Y is the sum of the conditional probabilities of Y given each possible outcome of another variable Z, weighted by the probability of each outcome of Z.

$$P(Y = i) = \sum_z P(Y = i | Z = z) P(z)$$

For example, if a and b are two dependent binary variables.

$$P(a) = P(a \cap b) + P(a \cap \neg b) \quad (2.36)$$

$$= P(a|b)p(b) + P(a|\neg b)P(\neg b) \text{ using conditional rule} \quad (2.37)$$

Example: Law of Total Probability

Let Y be the event of carrying an umbrella, and Z be the weather condition. The possible values of Z are:

- $Z = 1$: Sunny
- $Z = 2$: Cloudy
- $Z = 3$: Rainy

The conditional probabilities are:

$$\begin{aligned}P(Y = 1 | Z = 1) &= 0.1, \\P(Y = 1 | Z = 2) &= 0.3, \\P(Y = 1 | Z = 3) &= 0.9.\end{aligned}$$

The probabilities of the weather conditions are:

$$\begin{aligned}P(Z = 1) &= 0.4, \\P(Z = 2) &= 0.3, \\P(Z = 3) &= 0.3.\end{aligned}$$

We can compute $P(Y = 1)$ using the law of total probability:

$$\begin{aligned}P(Y = 1) &= P(Y = 1 | Z = 1)P(Z = 1) + P(Y = 1 | Z = 2)P(Z = 2) \\&\quad + P(Y = 1 | Z = 3)P(Z = 3)\end{aligned}$$

Substituting the values:

$$\begin{aligned}P(Y = 1) &= (0.1)(0.4) + (0.3)(0.3) + (0.9)(0.3) \\P(Y = 1) &= 0.04 + 0.09 + 0.27 = 0.4\end{aligned}$$

Thus, the probability that the person carries an umbrella is $P(Y = 1) = 0.4$.

CHAPTER 3

NAIVE BAYES CLASSIFICATION

Naive Bayes Classification is a probabilistic model used in supervised learning, where it learns from labeled training data to classify new, unseen data points. It applies Bayes' Theorem with the naive assumption that all input features are conditionally independent given the class label.

Despite its simplicity, it performs surprisingly well in many real-world tasks—especially in text classification problems like spam detection, sentiment analysis, and document categorization.

3.1 Key Idea

The model computes the **posterior probability** of each class given the input features:

$$P(C_k \mid x_1, x_2, \dots, x_n) = \frac{P(C_k) \cdot P(x_1, x_2, \dots, x_n \mid C_k)}{P(x_1, x_2, \dots, x_n)} \quad (3.1)$$

Using the **conditional independence assumption**:

$$P(x_1, x_2, \dots, x_n \mid C_k) = \prod_{i=1}^n P(x_i \mid C_k) \quad (3.2)$$

The posterior simplifies to:

$$P(C_k \mid x_1, x_2, \dots, x_n) \propto P(C_k) \cdot \prod_{i=1}^n P(x_i \mid C_k) \quad (3.3)$$

The predicted class is the one with the highest posterior probability.

3.2 Example: Why Independence Helps?

Problem: Classify emails as spam or not spam based on three features:

- x_1 : Presence of the word "free".

- x_2 : Presence of the word "win".
- x_3 : Presence of the word "money".

Without Independence:

- Joint probability $P(x_1, x_2, x_3 | \text{Spam})$ requires estimating $2^3 = 8$ probabilities.
- For n features, 2^n combinations need estimation.

With Independence:

$$P(x_1, x_2, x_3 | \text{Spam}) = P(x_1 | \text{Spam}) \cdot P(x_2 | \text{Spam}) \cdot P(x_3 | \text{Spam}).$$

- Only 3 probabilities, $P(x_1 | \text{Spam})$, $P(x_2 | \text{Spam})$, $P(x_3 | \text{Spam})$ need estimation.

3.3 Steps in Naive Bayes Classification

1. **Prepare Labeled Data:** Each data point has features (x_1, x_2, \dots, x_n) and a class label.
2. **Estimate Probabilities:** From training data, calculate:
 - Prior probabilities $P(C_k)$
 - Likelihoods $P(x_i | C_k)$
3. **Apply Bayes' Rule:** Use the formula to compute posterior probabilities.
4. **Predict Class:** Choose the class with the highest posterior.

3.4 Example: Probability of HIV

Problem: Calculate the probability of having HIV after a positive test result.

Given:

- HIV prevalence, $P(\text{HIV}) = 0.008$
- Test sensitivity, $P(T | \text{HIV}) = 0.95$
- Test specificity, $P(\neg T | \neg \text{HIV}) = 0.95$

Solution:

$$P(\text{HIV} | T) \propto P(T | \text{HIV}) \cdot P(\text{HIV}) = (0.95) \cdot (0.008) = 0.0076$$

$$P(\neg \text{HIV} | T) \propto P(T | \neg \text{HIV}) \cdot P(\neg \text{HIV}) = (0.05) \cdot (0.992) = 0.0496$$

Conclusion: Even with a positive result, the probability of having HIV is low due to the low prior probability.

3.5 Example: Naive Bayes for Spam Detection

Problem: Classify an email as "Spam" or "Not Spam" based on the occurrence of the words "Free", "Money" and "Win."

Dataset: The dataset consists of 50 spam emails and 50 non-spam emails. The count of the presence of words "Free", "Win", and "Money" is given in the dataset.

Word	Spam	Non-spam
"Free"	30	5
"Win"	25	10
"Money"	20	5

Table 3.1: Keyword Frequencies in Emails

Prior Probabilities:

$$P(S) = \frac{\text{Total Spam}}{\text{Total Emails}} = \frac{50}{100} = 0.5$$

$$P(\neg S) = \frac{\text{Total Non-spam}}{\text{Total Emails}} = \frac{50}{100} = 0.5$$

Likelihood Probabilities:

$$P(\text{Free}|S) = \frac{\text{Spam emails with "Free"}}{\text{Total Spam}} = \frac{30}{50} = 0.6$$

$$P(\text{Win}|S) = \frac{\text{Spam emails with "Win"}}{\text{Total Spam}} = \frac{25}{50} = 0.5$$

$$P(\text{Money}|S) = \frac{\text{Spam emails with "Money"}}{\text{Total Spam}} = \frac{20}{50} = 0.4$$

$$P(\text{Free}|\neg S) = \frac{\text{Non-spam emails with "Free"}}{\text{Total Not Spam}} = \frac{5}{50} = 0.1$$

$$P(\text{Win}|\neg S) = \frac{\text{Non-spam emails with "Win"}}{\text{Total Non-spam}} = \frac{10}{50} = 0.2$$

$$P(\text{Money}|\neg S) = \frac{\text{Non-spam emails with "Money"}}{\text{Total Non-spam}} = \frac{5}{50} = 0.1$$

Posterior Calculation:

$$\begin{aligned} P(\text{Spam}|\text{Free, Win, Money}) &\propto P(\text{Free}|S) \cdot P(\text{Win}|S) \cdot P(\text{Money}|S) \cdot P(S) \\ &= 0.6 \cdot 0.5 \cdot 0.4 \cdot 0.5 \\ &= 0.06 \end{aligned}$$

$$\begin{aligned}
 P(\neg S | \text{Free, Win, Money}) &\propto P(\text{Free}|\neg S) \cdot P(\text{Win}|\neg S) \cdot P(\text{Money}|\neg S) \cdot P(\neg S) \\
 &= 0.1 \cdot 0.2 \cdot 0.1 \cdot 0.5 \\
 &= 0.001
 \end{aligned}$$

Comparing the results, we can say that there is higher probability that the email is spam if all three words "Free", "Win" and "Money" are present in the email.

3.6 Example: Predicting Tennis Play from dataset

Dataset:

Sl. No.	Outlook	Temp.	Humidity	Windy	Play Tennis?
1	Sunny	Hot	High	False	No
2	Sunny	Hot	High	True	No
3	Overcast	Hot	High	False	Yes
4	Rainy	Mild	High	False	Yes
5	Rainy	Cool	Normal	False	Yes
6	Rainy	Cool	Normal	True	No
7	Overcast	Cool	Normal	True	Yes
8	Sunny	Mild	High	False	No
9	Sunny	Cool	Normal	False	Yes
10	Rainy	Mild	Normal	False	Yes
11	Sunny	Mild	Normal	True	Yes
12	Overcast	Mild	High	True	Yes
13	Overcast	Hot	Normal	False	Yes
14	Rainy	Mild	High	True	No

Table 3.2: Weather Dataset for Tennis Prediction

Computing Joint Probabilities

Predict whether to play tennis given: Outlook = Sunny, Temperature = Hot, Humidity = High, Windy = False

Learning Phase: Computing Probabilities

Step 1: Compute Class Priors

Class	Count	Prior Probability
Yes	9	$P(Y) = \frac{9}{14} = 0.64$
No	5	$P(N) = \frac{5}{14} = 0.35$

Table 3.3: Class Priors for the Tennis Dataset

Step 2: Compute Conditional Probabilities

For each feature and class, calculate $P(\text{Feature}|\text{Class})$.

Conditional Probabilities

Outlook

Outlook	Class: Yes	Class: No	Total
Sunny	$\frac{2}{9}$	$\frac{3}{5}$	5
Overcast	$\frac{4}{9}$	$\frac{0}{5}$	4
Rainy	$\frac{3}{9}$	$\frac{2}{5}$	5

Table 3.4: Conditional Probabilities for Outlook

Temperature

Temperature	Class: Yes	Class: No	Total
Hot	$\frac{2}{9}$	$\frac{2}{5}$	4
Mild	$\frac{4}{9}$	$\frac{2}{5}$	6
Cool	$\frac{3}{9}$	$\frac{1}{5}$	4

Table 3.5: Conditional Probabilities for Temperature

Humidity

Humidity	Class: Yes	Class: No	Total
High	$\frac{3}{9}$	$\frac{4}{5}$	7
Normal	$\frac{6}{9}$	$\frac{1}{5}$	7

Table 3.6: Conditional Probabilities for Humidity

Windy

Windy	Class: Yes	Class: No	Total
False	$\frac{6}{9}$	$\frac{2}{5}$	8
True	$\frac{3}{9}$	$\frac{3}{5}$	6

Table 3.7: Conditional Probabilities for Windy

Summary of Learned Probabilities

Feature	$P(\text{Feature} Y)$	$P(\text{Feature} N)$
Outlook = Sunny	0.22	0.6
Temperature = Hot	0.22	0.4
Humidity = High	0.33	0.8
Windy = False	0.66	0.4

Table 3.8: Likelihoods for Tennis Prediction

Final Decision

We compute the posterior probabilities using Bayes' Rule:

$$\begin{aligned} P(Y|X) &\propto P(Y) \cdot P(S|Y) \cdot P(\text{Hot}|Y) \cdot P(\text{High}|Y) \cdot P(\text{False}|Y) \\ &= (0.64)(0.22)(0.22)(0.33)(0.66) = 0.007 \end{aligned}$$

$$\begin{aligned} P(N|X) &\propto P(N) \cdot P(S|N) \cdot P(\text{Hot}|N) \cdot P(\text{High}|N) \cdot P(\text{False}|N) \\ &= (0.36)(0.6)(0.4)(0.8)(0.4) = 0.046 \end{aligned}$$

Result: The model predicts that we should **not** play tennis.

3.7 Example: Bayesian Diagnosis with another HIV Example

Single Positive Test

Given:

- $P(\text{HIV}) = 0.008$, $P(\text{Positive}|\text{HIV}) = 0.95$
- $P(\text{Positive}|\neg\text{HIV}) = 0.05$, $P(\neg\text{HIV}) = 0.992$

Using Bayes' theorem:

$$P(\text{HIV}|\text{Positive}) = \frac{0.95 \cdot 0.008}{0.95 \cdot 0.008 + 0.05 \cdot 0.992} \approx 0.156 \quad (3.4)$$

Conclusion: Despite a positive test, the chance of having HIV is only 15.6%.

Two Positive Tests

Assuming independence:

$$\begin{aligned} P(\text{HIV}|T_1, T_2) &= \frac{P(\text{HIV})P(T_1|\text{HIV})P(T_2|\text{HIV})}{P(T_1)P(T_2)} \\ &= \frac{0.008 \cdot 0.95 \cdot 0.95}{(0.008 \cdot 0.95 + 0.992 \cdot 0.05)^2} \approx 0.744 \end{aligned} \quad (3.5)$$

Conclusion: Two positive tests increase the probability to 74.4%.

3.8 Types of Naive Bayes Classifiers

- **Gaussian Naive Bayes:** - For continuous data. - Assumes features follow a Gaussian distribution:

$$P(x_i|C) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \quad (3.6)$$

- **Multinomial Naive Bayes:** - For count data (e.g., word counts). - Used in text classification.

- **Bernoulli Naive Bayes:** - For binary data (e.g., presence/absence of a word).

3.9 Advantages and Limitations

3.9.1 Advantages

- Simple, efficient, and easy to implement.

- Works well on small data and high dimensions.
- Handles categorical and numerical data.

3.9.2 Limitations

- Assumes independence of features.
- Performs poorly when features are correlated.
- Sensitive to zero-frequency problems which has to be solved using Laplace smoothing.

CHAPTER 4

DECISION TREE

4.1 Introduction

A Decision Tree is a supervised learning model used for both classification and regression tasks. It is tree-structured with:

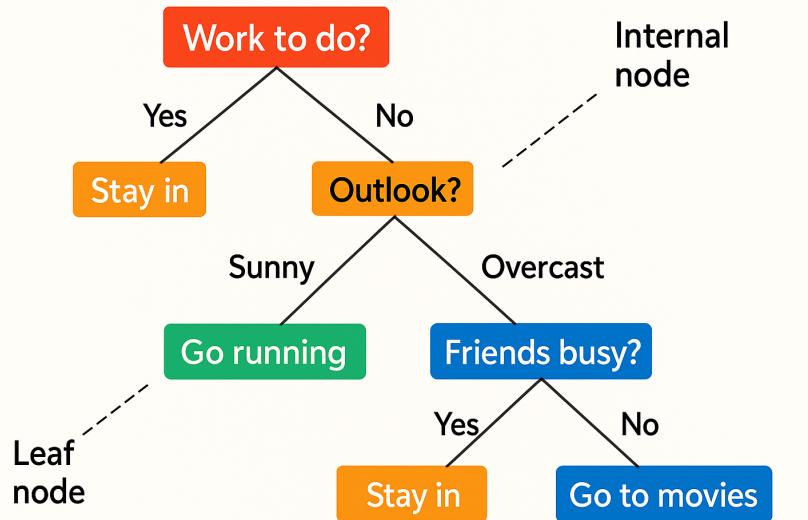


Figure 4.1: An Example of Decision Tree

- Internal nodes representing decisions or tests on features. Example: 4.1 Outlook, Friends Busy
- Branches representing the outcomes of those tests. Example: 4.1 Yes, No, Sunny, Overcast etc.
- Leaf nodes representing class labels (classification) or numerical predictions (regression). Example: 4.1 Stay in, Go running, Go to Movies.

The objective is to split the dataset in a way that reduces impurity or prediction error.

4.2 Entropy: A Measure of Uncertainty

Entropy, from Shannon's Information Theory, measures the impurity or disorder in a dataset D with k class labels, $C_1, C_2 \dots C_k$:

$$H(D) = - \sum_{i=1}^k P_i \log_2 P_i, \quad (4.1)$$

where P_i is the relative frequency (or probability) of class C_i .

This formula originates from Claude Shannon's foundational work on information theory. Entropy quantifies the expected amount of information (or surprise) from observing the outcome of a random variable, such as the class label of a data point.

If all classes are equally probable:

$$P_i = \frac{1}{k}$$

$$k = \frac{1}{P_i}$$

$$\log_2 k = -\log_2 P_i$$

Extending this idea to the general case where the classes C_1, C_2, \dots, C_k have arbitrary probabilities P_1, P_2, \dots, P_k , we can measure the total entropy as:

$$H(D) = - \sum_{i=1}^k P_i \log_2 P_i$$

We are measuring information in units of **bits**, which is why we use base 2 in the logarithm.

Consider a binary classification problem with 2 class labels. To distinguish between these two labels, we need 1 bit. More generally, if we have k classes, we need $\log_2 k$ bits to represent them.

Other common units for measuring information include:

- **Nats**, which use the natural logarithm (\log_e or \ln), commonly used in physics and information theory when working with continuous distributions.
- **Hartleys**, which use base 10 logarithms (\log_{10}), sometimes used in communication systems.

In the context of classification, entropy reflects the unpredictability of the class label.

- High entropy (close to 1 bit for binary classification) means classes are equally mixed, making the outcome uncertain.
- Low entropy (0) means all samples belong to one class, hence fully predictable.

High entropy in the training data indicates a rich diversity of examples across different classes. This is good for training because it provides the model with sufficient information to learn meaningful decision boundaries.

However, when splitting data at a node, high entropy is undesirable—it means the split has not made the data more pure or homogeneous. The goal of each split in a decision tree is to reduce entropy (increase purity), so a good split is one that creates low-entropy subsets.

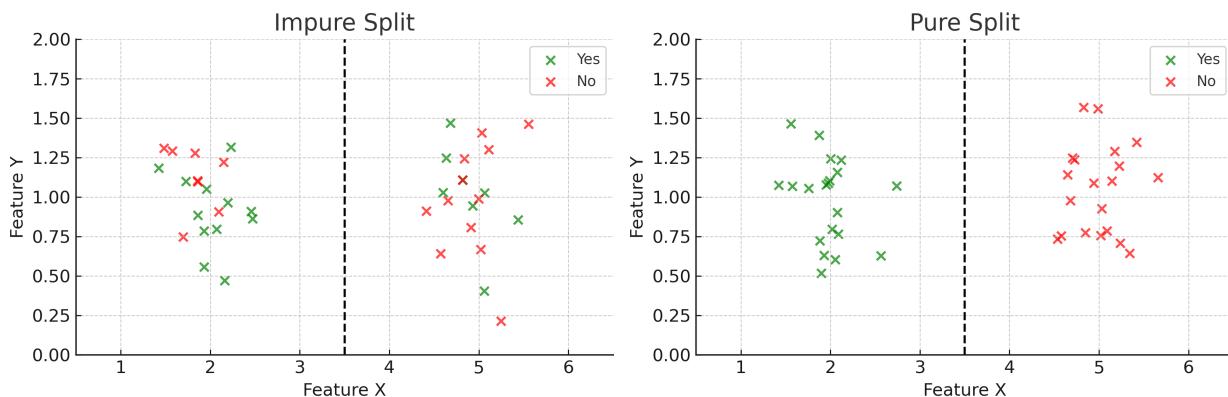


Figure 4.2: Visual comparison of an impure split (left) vs a pure split (right).

Example in fig: 4.2 we can see, in the impure split, both child nodes contain a mix of classes (Yes in green, No in red), indicating high entropy. In the pure split, the data is cleanly separated by class, resulting in low entropy and a more informative split.

- **High entropy in training data** ensures the model encounters all class types.
- **Low entropy in split subsets** ensures the decision tree is making clear distinctions that improve classification.

4.3 Conditional Entropy and Information Gain

Let's say a feature X can have m values (subsets) X_1, X_2, \dots, X_m after a split. We define the **Conditional Entropy** of the feature for the dataset as:

$$H(D|X) = \sum_{j=1}^m P(X = X_j) H(D|X = X_j) \quad (4.2)$$

$$H(D|X) = \sum_{j=1}^m \frac{\text{#of instances where } X = X_j}{\text{Total instances in } X} H(D|X = X_j) \quad (4.3)$$

Here, $H(D | X = X_j)$ also written as $H(D | X_j)$ or $H(D_{X_j})$ is the **Specific Conditional Entropy** for the subset where the feature X has value X_j .

$$H(D | X_j) = - \sum_{i=1}^k P(C_i | X_j) \log_2 P(C_i | X_j). \quad (4.4)$$

Here, $P(C_i | X_j)$ is the probability of the i 'th class, C_i in the subset where feature $X = X_j$.

Using all these, we can compute the **Information Gain (IG)** which measures reduction in entropy.

$$IG(D, X) = H(D) - H(D|X) \quad (4.5)$$

Features with higher IG are preferred for splitting.

4.4 Example: Toy Dataset

Instance	Weather	Play Tennis?
1	Sunny	Yes
2	Sunny	No
3	Rainy	Yes
4	Sunny	Yes
5	Rainy	No

Table 4.1: Example Dataset

Entropy of Entire Dataset

$$H(D) = - \left(\frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5} \right) \approx 0.971$$

Conditional Entropy given Weather

- Sunny: 3 instances (2 Yes, 1 No)

$$H(D_{Sunny}) = - \left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right) \approx 0.918$$

- Rainy: 2 instances (1 Yes, 1 No)

$$H(D_{Rainy}) = - \left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right) = 1$$

$$H(D|\text{Weather}) = P(\text{sunny}) \cdot H(D_{\text{sunny}}) + P(\text{rainy})H(D_{\text{rainy}})$$

$$H(D|\text{Weather}) = \frac{3}{5} \cdot 0.918 + \frac{2}{5} \cdot 1 = 0.9508$$

Information Gain

$$IG(D, \text{Weather}) = 0.971 - 0.9508 = 0.0202$$

4.5 The ID3 Algorithm

Pseudocode for ID3 Decision Tree Construction

Input: Dataset D , Feature Set F

Output: Decision Tree T

Procedure ID3(D, F):

1. **Compute** entropy $H(D)$
2. **For each** feature $X \in F$:
 - a. Compute Information Gain $IG(D, X)$
3. **Select** feature X^* with highest IG
4. **Split** D into subsets D_1, D_2, \dots, D_m based on the values v_1, v_2, \dots, v_m of X^*
5. **For each** subset D_j corresponding to $X^* = v_j$:
 - a. **If** D_j is pure (all same class) or F is empty:
 - Return leaf node with majority class label in D_j
 - b. **Else:** Recurse: $ID3(D_j, F - \{X^*\})$

4.6 Example: Full ID3 Walkthrough

Dataset: Play Tennis Example

Day	Outlook	Temp	Humidity	Wind	Play?
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Table 4.2: Play Tennis Dataset

Step 1: Compute Entropy of the Entire Dataset

9 instances are labeled Yes, and 5 are labeled No.

$$\begin{aligned}
 H(D) &= - \left(\frac{9}{14} \log_2 \frac{9}{14} + \frac{5}{14} \log_2 \frac{5}{14} \right) \\
 &= - (0.643 \cdot (-0.643) + 0.357 \cdot (-1.485)) \\
 &\approx 0.940
 \end{aligned}$$

Step 2: Compute Information Gain for Each Attribute

Compute Information Gain For Outlook

Compute Conditional Entropy $H(D|\text{Outlook})$

Sunny: 5 instances (2 Yes, 3 No)

$$\begin{aligned}
 H(D_{\text{Sunny}}) &= - \left(\frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5} \right) \\
 &= - (0.4 \cdot (-1.322) + 0.6 \cdot (-0.737)) \\
 &\approx 0.971
 \end{aligned}$$

Overcast: 4 instances (4 Yes, 0 No) $\Rightarrow H = 0$.

Rain: 5 instances (3 Yes, 2 No)

$$\begin{aligned} H(D_{Rain}) &= - \left(\frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5} \right) \\ &= - (0.6 \cdot (-0.737) + 0.4 \cdot (-1.322)) \\ &\approx 0.971 \end{aligned}$$

Weighted Conditional Entropy:

$$\begin{aligned} H(D|\text{Outlook}) &= \frac{5}{14} \cdot 0.971 + \frac{4}{14} \cdot 0 + \frac{5}{14} \cdot 0.971 \\ &= \frac{10}{14} \cdot 0.971 \\ &\approx 0.694 \end{aligned}$$

Compute Information Gain for Outlook

$$\begin{aligned} IG(D, \text{Outlook}) &= H(D) - H(D|\text{Outlook}) \\ &= 0.940 - 0.694 = 0.246 \end{aligned}$$

We do similar computation to find $IG(D, \text{Temp})$, $IG(D, \text{Wind})$, $IG(D, \text{Humidity})$.

Computing Information Gain of Temperature

- Hot: 4 instances \rightarrow 2 Yes, 2 No $\rightarrow H = 1$
- Mild: 6 instances \rightarrow 4 Yes, 2 No $\rightarrow H = -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \approx 0.918$
- Cool: 4 instances \rightarrow 3 Yes, 1 No $\rightarrow H = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \approx 0.811$

$$\begin{aligned} H(D|\text{Temp}) &= \frac{4}{14} \cdot 1 + \frac{6}{14} \cdot 0.918 + \frac{4}{14} \cdot 0.811 \\ &\approx 0.286 + 0.393 + 0.232 = 0.911 \end{aligned}$$

$$IG(D, \text{Temp}) = 0.940 - 0.911 = \boxed{0.029}$$

Computing Information Gain for Humidity

- High: 7 instances \rightarrow 3 Yes, 4 No $\rightarrow H = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \approx 0.985$
- Normal: 7 instances \rightarrow 6 Yes, 1 No $\rightarrow H = -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} \approx 0.592$

$$\begin{aligned}
 H(D|\text{Humidity}) &= \frac{7}{14} \cdot 0.985 + \frac{7}{14} \cdot 0.592 \\
 &= 0.5 \cdot (0.985 + 0.592) \\
 &= 0.789
 \end{aligned}$$

$$IG(D, \text{Humidity}) = 0.940 - 0.789 = [0.151]$$

Computing Information Gain of Wind

- Weak: 8 instances \rightarrow 6 Yes, 2 No $\rightarrow H = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \approx 0.811$
- Strong: 6 instances \rightarrow 3 Yes, 3 No $\rightarrow H = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = 1$

$$\begin{aligned}
 H(D|\text{Wind}) &= \frac{8}{14} \cdot 0.811 + \frac{6}{14} \cdot 1 \\
 &= 0.463 + 0.429 = 0.892
 \end{aligned}$$

$$\begin{aligned}
 IG(D, \text{Wind}) &= 0.940 - 0.892 \\
 &= [0.048]
 \end{aligned}$$

Summary of Information Gain

- $IG(\text{Outlook}) = 0.246$
- $IG(\text{Humidity}) = 0.151$
- $IG(\text{Wind}) = 0.048$
- $IG(\text{Temperature}) = 0.029$

Therefore, **Outlook** has the highest information gain and is selected as the **root node**.

Next Step: ID3 on Subset with Outlook = Sunny

Subset with Outlook = Sunny:

Day	Temp	Humidity	Wind	Play?
1	Hot	High	Weak	No
2	Hot	High	Strong	No
8	Mild	High	Weak	No
9	Cool	Normal	Weak	Yes
11	Mild	Normal	Strong	Yes

Table 4.3: Subset D_{Sunny}

Class counts: 2 Yes, 3 No $\Rightarrow H(D_{Sunny}) \approx 0.971$.

Compute IG for remaining features (Temp, Humidity, Wind):

- **Humidity:**

- High (3 instances, all No): $H = 0$
- Normal (2 instances, both Yes): $H = 0$
- $H(D_{Sunny}|Humidity) = \frac{3}{5} \cdot 0 + \frac{2}{5} \cdot 0 = 0$
- $IG = 0.971 - 0 = 0.971$

- **Temp and Wind** can be computed similarly.

ID3 on Subset with Outlook = Rainy

Subset D_{Rain} :

Day	Temp	Humidity	Wind	Play?
4	Mild	High	Weak	Yes
5	Cool	Normal	Weak	Yes
6	Cool	Normal	Strong	No
10	Mild	Normal	Weak	Yes
14	Mild	High	Strong	No

Table 4.4: Subset D_{Rain}

Class distribution: 3 Yes, 2 No $\Rightarrow H(D_{Rain}) = 0.971$ (previously computed).

Compute IG for remaining features:

- **Wind:**

- Weak: 3 instances (all Yes) $\Rightarrow H = 0$
- Strong: 2 instances (both No) $\Rightarrow H = 0$
- $H(D_{Rain}|Wind) = \frac{3}{5} \cdot 0 + \frac{2}{5} \cdot 0 = 0$
- $IG = 0.971 - 0 = 0.971$

- **Humidity:**

- High: 2 instances (1 Yes, 1 No) $\Rightarrow H = 1$

- Normal: 3 instances (2 Yes, 1 No)

$$H = - \left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right) \approx 0.918$$

- Weighted entropy:

$$H = \frac{2}{5} \cdot 1 + \frac{3}{5} \cdot 0.918 \approx 0.951$$

- $IG = 0.971 - 0.951 = 0.020$

- **Temp:**

- Mild: 3 instances (2 Yes, 1 No) $\Rightarrow H \approx 0.918$

- Cool: 2 instances (1 Yes, 1 No) $\Rightarrow H = 1$

- Weighted entropy:

$$H = \frac{3}{5} \cdot 0.918 + \frac{2}{5} \cdot 1 \approx 0.951$$

- $IG = 0.971 - 0.951 = 0.020$

Conclusion: Wind gives the highest information gain. So we split on **Wind**:

- $Wind = Weak \Rightarrow 3 \text{ Yes} \Rightarrow \text{Leaf Node: Yes}$
- $Wind = Strong \Rightarrow 2 \text{ No} \Rightarrow \text{Leaf Node: No}$

ID3 on Subset with Outlook = Overcast

Subset D_{Overcast} :

Day	Temp	Humidity	Wind	Play?
3	Hot	High	Weak	Yes
7	Cool	Normal	Strong	Yes
12	Mild	High	Strong	Yes
13	Hot	Normal	Weak	Yes

Table 4.5: Subset D_{Overcast}

All instances are labeled **Yes** $\Rightarrow H(D_{\text{Overcast}}) = 0$

This is a pure leaf node and doesn't require further splitting.

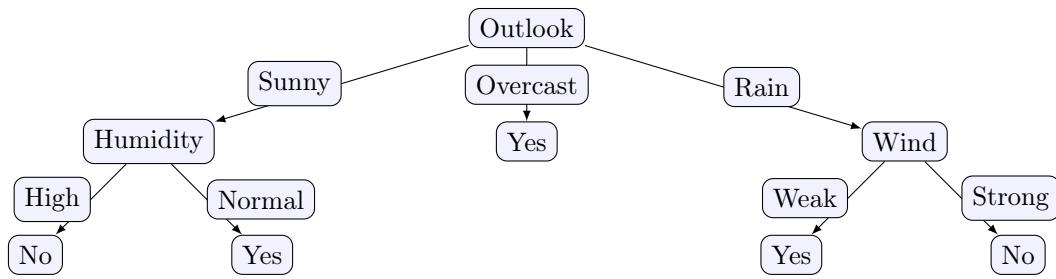


Figure 4.3: Final ID3 Decision Tree for Play Tennis Dataset

Special Example: Computing Entropy With Non-Binary Class

Consider a small dataset with 3 class labels: A, B, and C. The class distribution is as follows:

Class	Frequency
A	3
B	2
C	1

Table 4.6: Small Dataset with 3 Classes

Total instances: $N = 6$

Class probabilities:

$$\begin{aligned} P_A &= \frac{3}{6} = 0.5 \\ P_B &= \frac{2}{6} \approx 0.333 \\ P_C &= \frac{1}{6} \approx 0.167 \end{aligned}$$

Entropy is calculated as:

$$\begin{aligned} H(D) &= -(P_A \log_2 P_A + P_B \log_2 P_B + P_C \log_2 P_C) \\ &= -(0.5 \cdot \log_2 0.5 + 0.333 \cdot \log_2 0.333 + 0.167 \cdot \log_2 0.167) \\ &= 1.46 \text{ bits (approx)} \end{aligned}$$

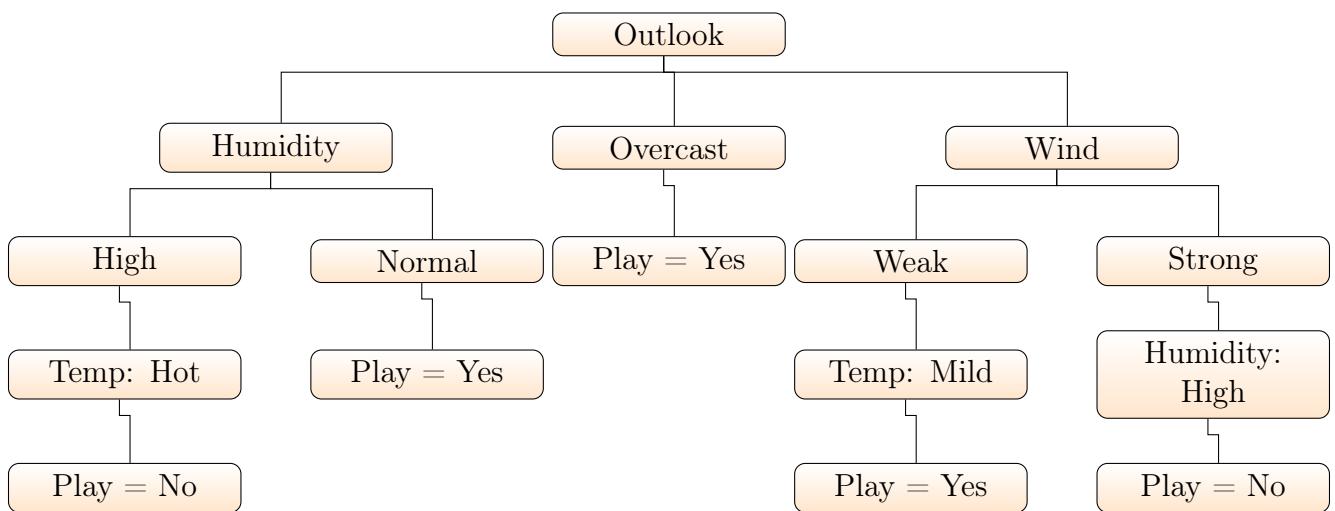
For, k number of class labels, we will have to use probability distribution $[P_1 \dots P_k]$, P_i representing the probability of the i 'th class label.

4.7 Overfitting and Pruning

Overfitting

A decision tree becomes overfitted when it grows too deep and begins to memorize noise or minor fluctuations in the training data, rather than learning the true underlying patterns. This often leads to excellent performance on the training set but poor generalization to unseen data.

Example: Overfit Tree



In this tree, the model makes decisions based on very specific combinations of attributes such as *Temp = Hot* and *Humidity = High* or *Wind = Strong* and *Humidity = High*, rather than broader, generalizable patterns. This level of detail may capture noise in the training dataset rather than meaningful trends.

Suppose a training dataset has a few instances where

- When *Outlook = Sunny*, *Humidity = High*, and *Temp = Hot*, the player didn't play.
- But for *Outlook = Sunny*, *Humidity = High*, and *Temp = Mild*, the player did play.

If the tree tries to fit such fine distinctions, it may become too sensitive to slight variations and overfit.

Underfitting

Underfitting occurs when a decision tree is too shallow or too simple to capture the underlying structure of the data. It fails to learn the relationships between features and the target class, resulting in poor performance on both the training and test sets.

An underfit model makes overly broad generalizations and may return the same prediction

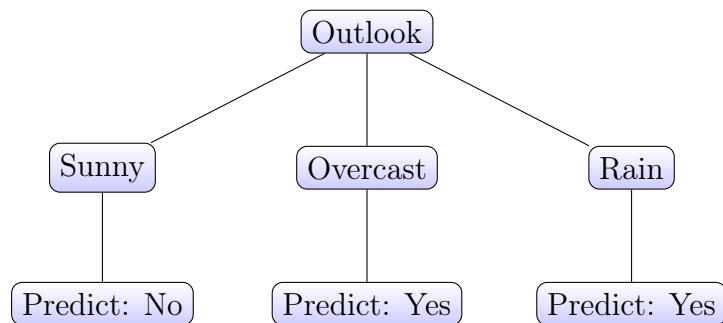
across many different inputs. This can happen due to early stopping, excessive pruning, or not including enough splits to isolate relevant patterns.

Example of Underfitting

Imagine trying to model a dataset using only the root node feature, without allowing the tree to explore deeper levels. For instance, if we split solely on *Outlook*, but ignore important distinctions made by *Humidity* or *Wind*, the tree might generalize:

- All "Overcast" days result in "Yes"
- All "Sunny" days result in "No"
- All "Rainy" days result in "Yes"

While this may cover dominant patterns, it ignores the nuances (e.g., differences based on humidity or wind strength), leading to high error on both known and new data.



Common Causes of Underfitting

- Tree depth restricted too early (pre-pruning)
- High minimum sample split thresholds
- Limited features available or used

Solution

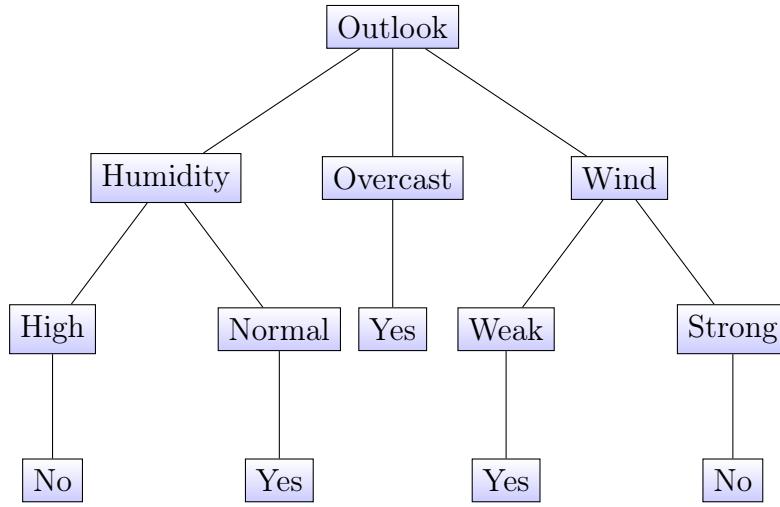
Allow the tree to grow deeper and use more features until training performance improves, followed by pruning to improve generalization.

Post-Pruning

To combat overfitting, we apply pruning strategies:

- **Post-pruning:** Build the full tree and then remove branches that do not improve accuracy on a validation set.

- **Pre-pruning (early stopping):** Stop tree construction early if further splits do not add significant value.



4.8 Handling Special Cases

4.8.1 Missing Data

Handling missing data is essential for building robust decision trees. There are several strategies to deal with missing attribute values:

- **Ignore the instance:** Remove records with missing values, especially if they form a small portion of the dataset.
- **Impute the missing value:** Replace with the mean (for numerical attributes), mode (for categorical attributes), or a more advanced method like KNN or regression-based imputation.

4.8.2 Continuous Attributes

Decision trees can handle continuous (numeric) features by choosing optimal split points:

- Determine possible thresholds (e.g., average of adjacent values).
- Evaluate each threshold using information gain (or Gini impurity).
- Select the threshold that maximizes the split quality.

Example: Given temperatures [60, 70, 80] and their associated play outcomes, try a threshold like $\text{Temp} \leq 65$ to split. If this improves the purity of child nodes, the threshold is retained.

4.8.3 Regression Trees: Handling Numerical Target Variables

When the target variable is numerical rather than categorical, we use regression trees instead of classification trees.

- These are suitable for tasks like predicting house prices, exam scores, or temperature.
- The tree is constructed by splitting the data at each node based on a feature and a threshold that minimizes prediction error — typically measured using metrics like Mean Squared Error (MSE) or Mean Absolute Error (MAE).
- Unlike classification trees that output a class label, regression tree leaves output a real-valued prediction (usually the average value of target variables in that node).

Example: Predicting House Prices

- Suppose we're predicting house prices based on area, number of bedrooms, and location score.
- A node might split on the rule " $\text{Area} \leq 2000 \text{ sq ft}$ ", separating smaller and larger homes.
- The leaf nodes would return values such as \$150,000 for smaller homes and \$320,000 for larger ones, based on the average price within each group.

CHAPTER 5

LINEAR REGRESSION AND GRADIENT DESCENT

NOTE: This chapter is based on 19.6 of Artificial Intelligence a Modern Approach (4th Ed.) by Stuart Russel and Peter Norvig

5.1 Introduction

Linear regression is a supervised learning algorithm in machine learning and statistics. It is used to model the relationship between a dependent variable (Y) and one or more independent variables ($x_1, x_2\dots$) by fitting a linear equation to observed data.

5.2 Mathematical Model

Linear regression assumes that the output variable y is a linear combination of the input features $\mathbf{x} = [x_1, x_2, \dots, x_n]$, with an additive error term ϵ . The general form of the linear regression model is:

$$y = w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n + \epsilon \quad (5.1)$$

Or in vectorized form:

$$y = \mathbf{w}^\top \mathbf{x} + \epsilon \quad (5.2)$$

Here,

- $\mathbf{x} \in \mathbb{R}^n$ is the input feature vector.
- $\mathbf{w} \in \mathbb{R}^n$ is the weight vector (parameters).
- w_0 is the intercept.
- ϵ is the noise or error term.

5.3 Learning Objective

In supervised learning, the core goal is to find a model that makes accurate predictions on unseen data. To achieve this, we define a **loss function**, which measures how well our model's predictions match the actual outputs. In linear regression, a widely used loss function is the **Mean Squared Error (MSE)**.

Mean Squared Error as a Loss Function

For a dataset of m training examples $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$, where $\mathbf{x}_i \in \mathbb{R}^n$ and $y_i \in \mathbb{R}$, the model predicts outputs as:

$$\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i \quad (5.3)$$

The Mean Squared Error (MSE) loss function is defined as:

$$\mathcal{L}(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2 \quad (5.4)$$

Where:

- $\mathcal{L}(\mathbf{w})$ is the loss function quantifying the average squared difference between predicted values and actual values.
- m is the number of training samples.
- y_i is the actual output for the i -th sample.
- \hat{y}_i is the predicted output: $\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i$.

Our objective is to **minimize this loss function** with respect to the weight vector \mathbf{w} . Minimizing the loss helps the model generalize well and make accurate predictions.

To find the optimal weights \mathbf{w} , we compute the gradient of $\mathcal{L}(\mathbf{w})$ with respect to \mathbf{w} and set it to zero:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 \quad (5.5)$$

this provides the condition for the minimum.

Solving this system (either analytically or using optimization methods like gradient descent) yields the weight vector \mathbf{w} that minimizes the prediction error across the training data.

5.4 Examples of Linear Regression

Linear regression is widely used in real-world tasks such as:

- Predicting housing prices from features like size and location.
- Estimating a student's exam score based on study hours.
- Modeling sales as a function of advertising budget.

5.5 Univariate Linear Regression

Univariate linear regression, also known as "fitting a straight line", is the simplest case of regression. Here, we consider the output variable (Y) to be dependent of one input variable (x). The output variable takes the form, $Y = w_0 + w_1x$.

Let us define the weight vector, $W = \langle w_0, w_1 \rangle$ and define the hypothesis (linear model) with these weights,

$$h_w(x) = w_1x + w_0 \quad (5.6)$$

can be used to predict the value of the output variable. We are given the following data of students' exam scores based on the number of hours they studied:

Example: Univariate Linear Regression

Hours Studied (x)	Exam Score (y)
1	52
2	55
3	61
4	66
5	71
6	75

Table 5.1: Dataset: Hours Studied vs Exam Score

Let us fit a linear regression model of the form:

$$y = w_0 + w_1x$$

We can determine the value of the weights solving 5.3 analytically.

Analytical solution

$$\frac{\partial \mathcal{L}}{\partial w_1} = 0 \text{ and } \frac{\partial \mathcal{L}}{\partial w_0}$$

$$\frac{\partial}{\partial w_0} \sum_{i=1}^m (y - (w_1x + w_0))^2 = 0 \text{ and } \frac{\partial}{\partial w_1} \sum_{i=1}^m (y - (w_1x + w_0))^2 = 0.$$

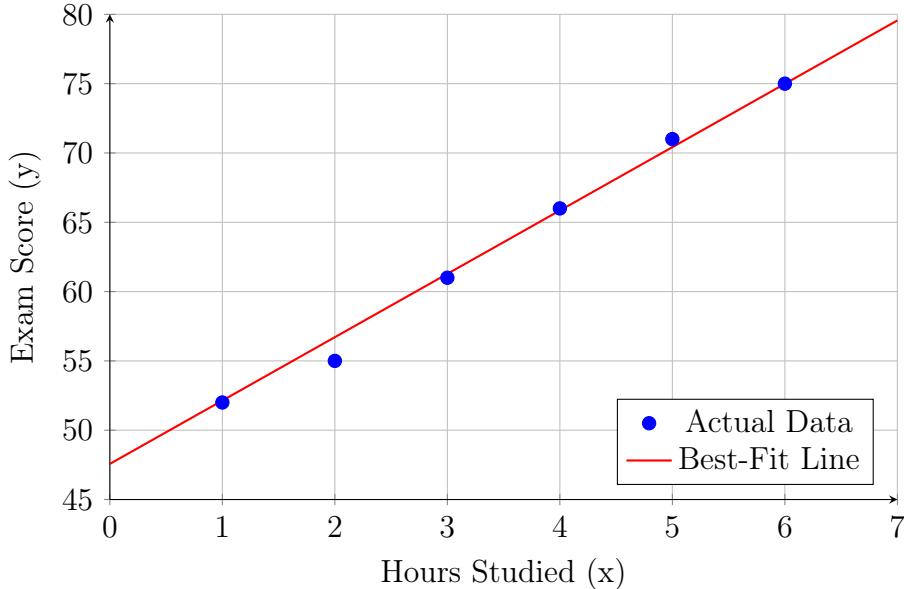
Solving these equations, we get:

$$w_1 = \frac{m(\sum x_j y_j) - (\sum x_j)(\sum y_j)}{m(\sum x_j^2) - (\sum x_j)^2}; \quad w_0 = \frac{1}{m} \left(\sum y_j - w_1 (\sum x_j) \right) \quad (5.7)$$

and found that the best-fit line for this data is:

$$\hat{y} = 47.57 + 4.57x.$$

Data Plot with Regression Line



Example: Computing Loss using MSE

Using mean squared error in (5.3) we find the Loss function for the weights, $w_0 = 47.57$ and $w_1 = 4.57$ for the dataset given in table: 5.1.

x_i	y_i	$\hat{y}_i = 47.57 + 4.57x_i$	$y_i - \hat{y}_i$	$(y_i - \hat{y}_i)^2$
1	52	52.14	-0.14	0.0196
2	55	56.71	-1.71	2.9241
3	61	61.28	-0.28	0.0784
4	66	65.85	0.15	0.0225
5	71	70.42	0.58	0.3364
6	75	74.99	0.01	0.0001
Total Sum of Squared Errors				3.3811
MSE = Total / 6				0.5635

Table 5.2: Step-by-Step Computation of Mean Squared Error (MSE)

5.5.1 Multivariable Linear Regression

In multivariable linear regression, where the output depends on n independent input variables x_1, x_2, \dots, x_m , we aim to find n corresponding weights w_1, w_2, \dots, w_m . These weights are

considered optimal when the partial derivatives of the loss function L with respect to each weight are zero:

$$\frac{\partial \mathcal{L}}{\partial w_1} = 0, \quad \frac{\partial \mathcal{L}}{\partial w_2} = 0, \quad \dots, \quad \frac{\partial \mathcal{L}}{\partial w_m} = 0 \quad (5.8)$$

Linear regression has a closed-form solution using the normal equation:

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (5.9)$$

This solution directly computes the optimal weights that minimize the loss function.

However, computing this closed form solution to find the optimal solution can be complex specially in high dimensional data.

While the analytical solution is faster for small, simple problems, we introduce **gradient descent algorithm** for modern, large-scale, and complex machine learning models due to its flexibility and efficiency.

5.6 Gradient Descent Algorithm

Gradient descent moves iteratively through this space to find the point where the loss is minimized. Starting from an initial point $\mathbf{w}^{(0)}$, we compute the gradient:

$$\nabla \mathcal{L}(\mathbf{w}) = \left[\frac{\partial \mathcal{L}}{\partial w_0}, \frac{\partial \mathcal{L}}{\partial w_1}, \dots \right] \quad (5.10)$$

We update the parameters in the direction of steepest descent:

$$w_i^{(t+1)} = w_i^{(t)} - \alpha \frac{\partial \mathcal{L}}{\partial w_i} \Big|_{w_i=w_i^t} \quad (5.11)$$

Here, α is the learning rate that controls the step size, t represents the iteration number and i refers to the specific parameter in the weight vector being updated.

Algorithm 8 Gradient Descent Algorithm

- 1: **Input:** Learning rate α , initial weights \mathbf{w} , loss function $L(\mathbf{w})$
 - 2: **Output:** Optimized weights \mathbf{w}
 - 3: **while** not converged **do**
 - 4: Compute gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{w}}$
 - 5: Update weights: $w \leftarrow w - \alpha \frac{\partial \mathcal{L}}{\partial w}$
 - 6: **end while**
 - 7: **return** \mathbf{W}
-

Learning Parameter, α

The parameter α in Equation 5.11 is called the learning rate or step size. It controls how quickly the algorithm updates the weights and how fast it converges to a minimum.

Choosing an appropriate learning rate is important. If α is too large, the algorithm may overshoot or oscillate around the minimum and fail to converge. If it is too small, the algorithm will take a long time to reach the minimum.

The learning rate can be a constant or a decaying parameter. A decaying learning rate helps the algorithm explore more widely in the early stages and fine-tune near the minimum.

Computing Gradient of Loss ($\frac{\partial \mathcal{L}}{\partial w_k}$) in Univariate Linear Regression with Mean Square Error

Let us consider a univariate linear regression with mean square error as the loss function. We already know,

$$\hat{y} = h_w(x) = w_0 + w_1 x.$$

and

$$\mathcal{L} = \frac{1}{m} \sum_{i=1}^m (y - \hat{y})^2$$

$$\frac{\partial \mathcal{L}}{\partial w_0} = \frac{\partial}{\partial w_0} \frac{1}{m} \sum_{i=1}^m (y - \hat{y})^2$$

Using chain rule of differentiation,

$$\frac{\partial}{\partial w_0} \frac{1}{m} \sum_{i=1}^m (y - \hat{y})^2 = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_0} (y - \hat{y})^2 \quad (5.12)$$

$$= \frac{2}{m} \sum_{i=1}^m (y - \hat{y}) \frac{\partial}{\partial w_0} (y - \hat{y})$$

$$= -\frac{2}{m} \sum_{i=1}^m (y - \hat{y}) \frac{\partial \hat{y}}{\partial w_0}$$

$$= -\frac{2}{m} \sum_{i=1}^m (y - \hat{y}) \frac{\partial}{\partial w_0} (w_0 + w_1 x)$$

$$\frac{\partial \mathcal{L}}{\partial w_0} = -\frac{2}{m} \sum_{i=1}^m (y - \hat{y}) \quad (5.13)$$

$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial}{\partial w_1} \frac{1}{m} \sum_{i=1}^m (y - \hat{y})^2$$

Using chain rule of differentiation,

$$\begin{aligned}
 \frac{\partial}{\partial w_1} \frac{1}{m} \sum_{i=1}^m (y - \hat{y})^2 &= \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} (y - \hat{y})^2 \\
 &= \frac{2}{m} \sum_{i=1}^m (y - \hat{y}) \frac{\partial}{\partial w_1} (y - \hat{y}) \\
 &= -\frac{2}{m} \sum_{i=1}^m (y - \hat{y}) \frac{\partial \hat{y}}{\partial w_1} \\
 &= -\frac{2}{m} \sum_{i=1}^m (y - \hat{y}) \frac{\partial}{\partial w_1} (w_0 + w_1 x) \\
 \frac{\partial \mathcal{L}}{\partial w_1} &= -\frac{2}{m} \sum_{i=1}^m (y - \hat{y}) x_1
 \end{aligned} \tag{5.14}$$

Computing Gradient of Loss ($\frac{\partial \mathcal{L}}{\partial w_k}$) in Multivariate Linear Regression with Mean Square Error

$$\frac{\partial \mathcal{L}}{\partial w_0} = -\frac{2}{m} \sum_{i=1}^m (y - \hat{y}) \tag{5.15}$$

$$\frac{\partial \mathcal{L}}{\partial w_k} = -\frac{2}{m} \sum_{i=1}^m (y - \hat{y}) x_k \tag{5.16}$$

Here, x_k denotes the k 'th input variable and w_k is its corresponding weight.

Parameter Update Rule with Gradient Descent

$$w_0^{(t+1)} = w_0^{(t)} + \alpha \frac{1}{m} \sum_{i=1}^m (y - \hat{y}). \tag{5.17}$$

Since, 2 is a constant it can be absorbed into the learning parameter α .

$$w_k^{(t+1)} = w_k^{(t)} + \alpha \frac{1}{m} \sum_{i=1}^m (y - \hat{y}) x_k. \quad (5.18)$$

Example: Parameter Update with Loss of Gradient

We compute the gradients of the MSE loss with respect to w_0 and w_1 as follows:

$$\frac{\partial \mathcal{L}}{\partial w_0} = \frac{-2}{m} \sum_{i=1}^m (y_i - \hat{y}_i), \quad \frac{\partial \mathcal{L}}{\partial w_1} = \frac{-2}{m} \sum_{i=1}^m (y_i - \hat{y}_i) x_i$$

At initial guess $w_0 = 0$, $w_1 = 0$, the prediction $\hat{y}_i = 0$.

Table 5.3: Gradient of Loss with Respect to w_0 and w_1

x_i	y_i	$\hat{y}_i = 0$	$y_i - \hat{y}_i$	$y_i - \hat{y}_i$	$(y_i - \hat{y}_i)x_i$
1	52	0	52	52	52
2	55	0	55	55	110
3	61	0	61	61	183
4	66	0	66	66	264
5	71	0	71	71	355
6	75	0	75	75	450
Total:			380		1414

From the table:

$$\begin{aligned} \sum_{i=1}^6 (y - \hat{y}) &= 380 \\ \sum_{i=1}^6 (y - \hat{y}) x_i &= 1414 \end{aligned}$$

Let's take $\alpha = .1$, and with $m = 6$, the after the first iteration the updated weights will be

$$\begin{aligned} w_0^{(2)} &= w_0^{(1)} + \alpha \frac{\partial \mathcal{L}}{\partial w_0} \\ &= 0 + (.1) \frac{1}{6} \cdot 380 \end{aligned}$$

and

$$\begin{aligned} w_1^{(2)} &= w_1^{(1)} + \alpha \frac{\partial \mathcal{L}}{\partial w_1} \\ &= 0 + (.1) \frac{1}{6} \cdot 1414 \end{aligned}$$

In the second iteration, we will use the updated weights, $w_0^{(2)}$ and $w_1^{(2)}$ to predict \hat{y} and compute the gradient.

The gradient descent algorithm will repeat this process until convergence.

Understanding Gradient Descent in Weight Space

In supervised learning models like linear regression, we represent our model's parameters as a vector $\mathbf{w} = [w_0, w_1, \dots, w_n]^\top$. Each point in this multi-dimensional space represents a particular configuration of the model.

Weight Parameter Space

Consider a model with two parameters w_0 (intercept) and w_1 (slope). The space formed by all possible values of these parameters is called the **weight parameter space**. Every point (w_0, w_1) in this space corresponds to a different hypothesis or model.

Loss Surface Over Parameter Space

We define a loss function $\mathcal{L}(\mathbf{w})$, such as Mean Squared Error (MSE), that measures how well a model with parameters \mathbf{w} fits the training data. The loss function can be visualized as a surface defined over the weight space. The height of the surface at each point \mathbf{w} corresponds to the loss for that model.

Geometric Intuition

Gradient descent can be viewed as a ball rolling downhill on the loss surface:

- The gradient points in the direction of steepest ascent.
- The negative gradient is the direction of steepest descent.

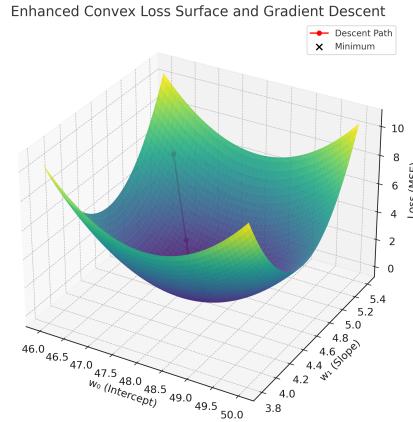


Figure 5.1: Loss vs Weights

- Each iteration moves the parameter vector \mathbf{w} closer to the minimum of the loss function.

Example: Gradient Descent Algorithm Applied

Let us apply gradient descent algorithm to the example dataset in 5.1. We have plotted the Loss against the parameter space, w_0 and w_1 in fig:5.1.

We can see that it has taken a convex form because of the quadratic nature of the Loss function. We can identify the minima (point where the loss is minimum) at 47.57 and 4.57.

Benefits of This Perspective

Understanding gradient descent in weight space helps to:

- Visualize optimization as navigation across a surface.
- Interpret convergence behavior (e.g., overshooting, slow descent).
- Tune hyperparameters like the learning rate.
- Understand curvature, flat regions, and saddle points.

Variations of Gradient Descent

Gradient Descent has several variations, depending on how much data is used to compute the gradient at each step. The three most common types are:

1. Batch Gradient Descent

In batch gradient descent, the gradient is computed using the entire training dataset:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial}{\partial \mathbf{w}} \mathcal{L}_{\text{full batch}}(\mathbf{w})$$

Pros: Accurate direction of descent.

Cons: Can be slow and memory-intensive for large datasets.

2. Stochastic Gradient Descent (SGD)

In SGD, the gradient is estimated using only a single training example at each iteration:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial}{\partial \mathbf{w}} \mathcal{L}_i(\mathbf{w})$$

Pros: Fast and can escape local minima.

Cons: Noisy updates, may not converge smoothly.

3. Mini-batch Gradient Descent

This is a compromise between batch and stochastic gradient descent. The gradient is computed on a small batch of examples (e.g., 32 or 64 samples):

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial}{\partial \mathbf{w}} \mathcal{L}_{\text{mini-batch}}(\mathbf{w})$$

Pros: Efficient and stable. Widely used in practice.

Cons: Requires tuning batch size for best performance.

Why Use Gradient Descent?

Gradient descent is often preferred when:

- The dataset is too large to store or invert the matrix $\mathbf{X}^\top \mathbf{X}$
- The model is nonlinear or does not have a closed-form solution
- We want to train on streaming data in real-time (online learning)
- Regularization techniques are used (e.g., L1, L2)

5.7 Advantages of Linear Regression

- It is **interpretable** and easy to implement.
- It can be solved analytically using closed-form solutions or optimized using gradient descent.
- It forms the basis of more complex models like logistic regression and neural networks.

CHAPTER 6

LINEAR CLASSIFIER AND LOGISTIC REGRESSION

NOTE: This chapter is written based on chapter 19.6.4 and 19.6.5 of Artificial Intelligence A Modern Approach by Stuart Russel and Peter Norvig

6.1 Linear Functions with Thresholds for Classification

Linear functions can also be used for classification tasks. In this case, the linear equation is used to define a boundary that separates the two classes. This linear equation is called the **Decision Boundary**. A linear decision boundary, also known as a linear separator, is used to separate data into classes. If a dataset can be perfectly divided by such a boundary, it is called linearly separable.

Our reference book shows an example of classification between two types of seismic events:

- Earthquakes, which interest seismologists, labeled 0.
- underground explosions, which concern arms control experts, labeled 1.

The classification is dependent on two input features

- The magnitudes of body, x_1 .
- surface waves from the seismic signal, x_2

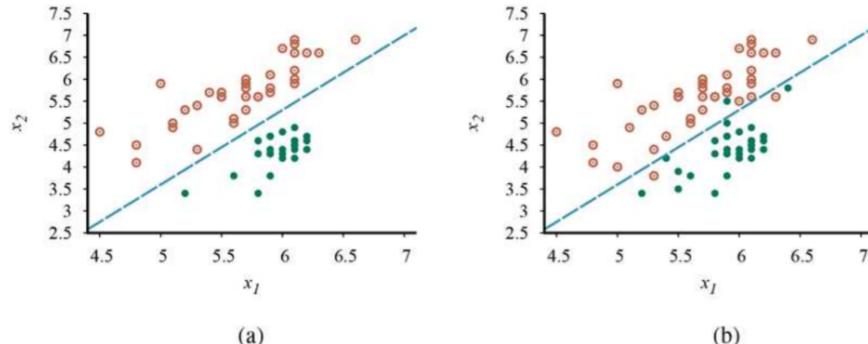
The goal is to learn a hypothesis, h that can correctly classify new points, (x_1, x_2) , labeling earthquakes as 0 and explosions as 1. The linear separator for this example dataset is

$$-4.9 + 1.7x_1 + x_2 = 0$$

From the plots we can see that the data points for explosions are below the line, meaning

$$-4.9 + 1.7x_1 + x_2 > 0$$

Figure 19.15



(a) Plot of two seismic data parameters, body wave magnitude x_1 and surface wave magnitude x_2 , for earthquakes (open orange circles) and nuclear explosions (green circles) occurring between 1982 and 1990 in Asia and the Middle East (Kebeasy *et al.*, 1998). Also shown is a decision boundary between the classes. (b) The same domain with more data points. The earthquakes and explosions are no longer linearly separable.

and the data points for earthquakes are above the line, meaning

$$-4.9 + 1.7x_1 + x_2 < 0$$

In Linear Classification with threshold we want to develop an equation for the decision boundary.

$$Z = \mathbf{W} \cdot \mathbf{X} + W_0 \quad (6.1)$$

by predicting the weights, \mathbf{W} and W_0 . This equation is then used to construct the hypothesis,

$$h_{\mathbf{W}}(\mathbf{x}) = \begin{cases} 1 & \text{if } Z > 0; \\ 0 & \text{if } Z < 0 \end{cases} \quad (6.2)$$

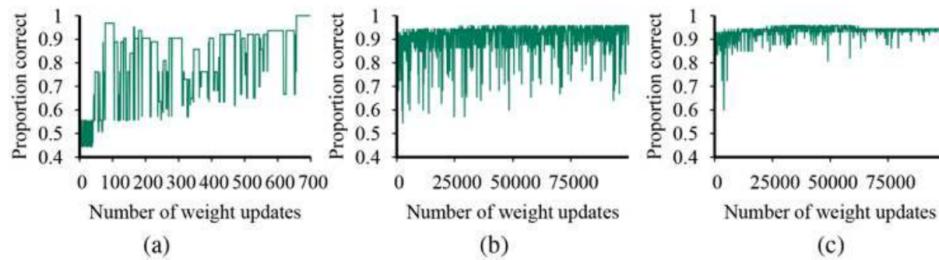
here,

- $h_{\mathbf{W}}(\mathbf{x})$ represents the hypothesis for the classification.
- Z is the equation of decision boundary.
- \mathbf{W} is the weight vector.
- \mathbf{X} is the value of the input features, $X_1, X_2 \dots X_m$.
- W_0 is the bias term.

6.1.1 Perceptron Learning

Since, the result of $h_{\mathbf{W}}(\mathbf{x})$ is either 0 or 1, we cannot choose the \mathbf{W} minimizing loss. The gradient in the weight space is 0 almost everywhere except the point where $W \cdot X = 0$. So

Figure 19.16



(a) Plot of total training-set accuracy vs. number of iterations through the training set for the perceptron learning rule, given the earthquake/explosion data in Figure 19.15(a). (b) The same plot for the noisy, nonseparable data in Figure 19.15(b); note the change in scale of the x -axis. (c) The same plot as in (b), with a learning rate schedule $\alpha(t) = 1000/(1000 + t)$.



gradient descent will not work in this case. To solve this problem we introduce **perceptron learning**. In this process the weights are only adjusted when the prediction is incorrect.

$$W_0 = W_0 + \alpha \cdot (y - \mathbf{h}_w(\mathbf{X})). \quad (6.3)$$

$$W_j = W_j + \alpha \cdot (y - \mathbf{h}_w(\mathbf{X})) X_j. \quad (6.4)$$

Typically, perceptron rule uses a randomly chosen example from the data to compute the updated weight.

A **training curve** plots how well the classifier is doing on the same training data as it keeps learning, one update at a time.

When Perceptron Learning Fails

Perceptron rule does not work well with linearly non-separable noisy data.

When the data are not linearly separable, there is no perfect line that can divide the two classes without mistakes. The perceptron algorithm is designed to adjust its weights whenever it makes a mistake. Since mistakes can never be fully avoided in this case, the perceptron keeps finding errors and keeps adjusting the weights forever, even if it often comes close to the best possible solution. This stops the perceptron rule from converging.

Plot **a** of figure ?? shows how the perceptron learning rule improves over time when it is trained on linearly separable data.

In this case, the curve shows that the perceptron eventually finds a perfect (zero-error) straight-line separator between the two classes.

However, the learning process is not smooth—it jumps around before finally getting it right.

The number of updates can vary between different training runs. Plot **b** and **c** shows that

the perceptron rule fails to converge even after 10,000 updates when dealing with noisy data.

Learning Parameter for Perceptron Rule

A decaying learning rate helps because, in the beginning, the algorithm can take big steps to explore quickly, and later, it takes smaller, finer steps to carefully settle near the best solution without overshooting. This gradual slowdown allows the perceptron to stabilize instead of endlessly bouncing around.

Plot c of ?? shows the training process when using a decaying learning rate. It still doesn't reach perfect convergence even after 100,000 updates, but it does much better than when using a fixed learning rate.

Normally, the perceptron rule may not converge to a stable solution if the learning rate stays fixed.

However, if we make the learning rate smaller over time — for example, using $\alpha = \frac{1000}{1000+t}$ (where t is the number of updates) — then the perceptron can be shown to converge to a minimum-error solution, as long as the training examples are presented in a random order.

Still, finding the true minimum-error solution is a very hard problem (NP-hard), so we expect that it will take many passes over the data before convergence happens.

6.1.2 Linear Classifier with Logistic Regression

In linear classification with logistic regression, instead of using hard thresholds for classification, we use a sigmoid function as the classifier.

$$h_w(Z) = \frac{1}{1 + e^{-Z}}, \quad (6.5)$$

where,

$$Z = \mathbf{W} \cdot \mathbf{X} + W_0. \quad (6.6)$$

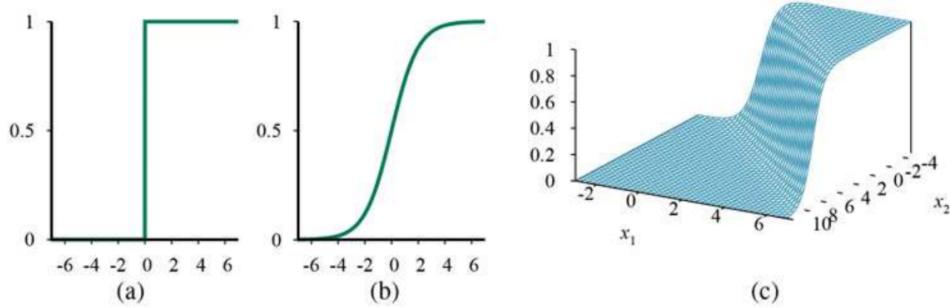
Here,

- W_0 is the bias term.
- X is the input vector with inputs $X_1, X_2 \dots X_m$.
- W is the weight vector with $W_1, W_2 \dots W_m$

The sigmoid function outputs the probability that an example belongs to Class 1.

The probability of belonging to Class 0 is simply the complement of this probability.

Unlike the hard threshold in 6.2, logistic regression uses a continuous function for its model, $h_w(\mathbf{X})$.



(a) The hard threshold function $\text{Threshold}(z)$ with 0/1 output. Note that the function is nondifferentiable at $z = 0$. (b) The logistic function, $\text{Logistic}(z) = \frac{1}{1+e^{-z}}$, also known as the sigmoid function. (c) Plot of a logistic regression hypothesis $h_w(\mathbf{x}) = \text{Logistic}(\mathbf{w} \cdot \mathbf{x})$ for the data shown in Figure 19.15(b).

Figure 6.1

Therefore, logistic regression minimizes a continuous loss function (MSE or binary cross-entropy), which is smooth and differentiable.

This allows the use of gradient descent, leading to more stable and gradual updates.

In contrast to perceptron learning, Logistic regression is able to find best-fit boundary even with noisy linearly non-separable data.

6.1.3 Gradient Descent for Logistic Regression

Since, we now use a continuous function for our hypothesis, we can implement the gradient descent algorithm to optimize the weights. The update rules will be as follows:

$$w_0^{(t+1)} = w_0^{(t)} - \alpha \frac{\partial \mathcal{L}}{\partial w_0} \quad (6.7)$$

$$w_1^{(t+1)} = w_1^{(t)} - \alpha \frac{\partial \mathcal{L}}{\partial w_1} \quad (6.8)$$

6.1.4 Computing Gradient

We will show the computation of gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{W}}$ for two types of Loss function.

- Mean squared Error, $L_2 = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$.
- Binary Cross-Entropy, $L = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \ln \hat{y}^{(i)} + (1 - y^{(i)}) \ln (1 - \hat{y}^{(i)})]$.

Here, $\hat{y} = \mathbf{h}_{\mathbf{w}}(\mathbf{Z}) = \frac{1}{1+e^{-z}}$

6.1.5 Logistic regression: Computing Gradient for Mean Squared Error

Remember,

$$Z = W_0 + W_1X_1 + W_2X_2 \dots W_mX_m$$

So,

$$\frac{\partial Z}{\partial W_0} = 1 \quad (6.9)$$

and

$$\frac{\partial Z}{\partial W_j} = X_j \quad (6.10)$$

Where, X_j is the j'th input and W_j is the associated weight.

Computing $\frac{\partial \hat{y}}{\partial W_0}$:

$$\begin{aligned} \frac{\partial \hat{y}}{\partial W_0} &= \frac{\partial}{\partial W_0} \left(\frac{1}{1 + e^{-z}} \right) \\ &= \frac{-1}{(1 + e^{-z})^2} \frac{\partial}{\partial W_0} (1 + e^{-z}) \\ &= \frac{-1}{(1 + e^{-z})^2} \frac{\partial (e^{-z})}{\partial W_0} \\ &= \frac{e^{-z}}{(1 + e^{-z})^2} \frac{\partial z}{\partial W_0} \\ &= \frac{1}{(1 + e^{-z})} \frac{e^{-z}}{(1 + e^{-z})} \frac{\partial z}{\partial W_0} \\ &= \frac{1}{(1 + e^{-z})} \frac{e^{-z} + 1 - 1}{(1 + e^{-z})^2} \\ &= \frac{1}{(1 + e^{-z})} \left(1 - \frac{1}{1 + e^{-z}} \right) \frac{\partial z}{\partial W_0} \\ &= \frac{1}{(1 + e^{-z})} \left(1 - \frac{1}{1 + e^{-z}} \right) \quad [\text{using 6.9}] \end{aligned}$$

$$\frac{\partial \hat{y}}{\partial W_0} = \hat{y}(1 - \hat{y}) \quad (6.11)$$

By similar calculation and using 6.10:

$$\frac{\partial \hat{y}}{\partial W_j} = \hat{y}(1 - \hat{y})X_j \quad (6.12)$$

Computing the gradient of Loss, Mean Squared Error

$$\frac{\partial \mathcal{L}}{\partial W_0} = \frac{\partial}{\partial W_0} \sum_{i=1}^m \frac{1}{m} (y^{(i)} - \hat{y}^{(i)})^2$$

using chain rule of differentiation,

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W_0} &= -\frac{2}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) \frac{\partial}{\partial W_0} (y^{(i)} - \hat{y}^{(i)}) \\ &= -\frac{2}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) \frac{\partial \hat{y}}{\partial W_0} \end{aligned}$$

Similarly,

$$\frac{\partial \mathcal{L}}{\partial W_j} = -\frac{2}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) \frac{\partial \hat{y}}{\partial W_j} \quad (6.13)$$

using 6.11 and 6.12, we get,

$$\frac{\partial \mathcal{L}}{\partial W_0} = -\frac{2}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) \hat{y}^{(i)} (1 - \hat{y}^{(i)}) \quad (6.14)$$

and

$$\frac{\partial \mathcal{L}}{\partial W_j} = -\frac{2}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) \hat{y}^{(i)} (1 - \hat{y}^{(i)}) X_j \quad (6.15)$$

Finally, we can write the update rule in the gradient descent algorithm for logistic regression with mean squared error loss.

$$W_0^{(t+1)} = W_0^{(t)} + \alpha \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) \hat{y}^{(i)} (1 - \hat{y}^{(i)}) \quad (6.16)$$

and

$$W_j^{(t+1)} = W_j^{(t)} + \alpha \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) \hat{y}^{(i)} (1 - \hat{y}^{(i)}) X_j \quad (6.17)$$

6.1.6 Logistic Regression: Computing gradient of Binary Cross-Entropy

The Loss function with Binary Cross-Entropy is given by

$$\mathcal{L} = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \ln \hat{y}^{(i)} + (1 - y^{(i)}) \ln (1 - \hat{y}^{(i)})] \quad (6.18)$$

Now,

$$\frac{\partial \mathcal{L}}{\partial W_0} = -\frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial W_0} [y^{(i)} \ln \hat{y}^{(i)} + (1 - y^{(i)}) \ln (1 - \hat{y}^{(i)})] \quad (6.19)$$

Computing $\frac{\partial}{\partial W_0}(y \ln \hat{y})$

$$\begin{aligned} \frac{\partial}{\partial W_0}(y \ln \hat{y}) &= \frac{y}{\hat{y}} \frac{\partial \hat{y}}{\partial W_0} \\ &= \frac{y}{\hat{y}} \cdot \hat{y}(1 - \hat{y}) \quad [\text{Using 6.11}] \\ \frac{\partial}{\partial W_0}(y \ln \hat{y}) &= y - y\hat{y} \end{aligned} \quad (6.20)$$

Computing $\frac{\partial}{\partial W_0}((1 - y) \ln (1 - \hat{y}))$

$$\begin{aligned} \frac{\partial}{\partial W_0}((1 - y) \ln (1 - \hat{y})) &= \frac{1 - y}{1 - \hat{y}} \frac{\partial(1 - \hat{y})}{\partial W_0} \\ &= -\frac{1 - y}{1 - \hat{y}} \frac{\partial \hat{y}}{\partial W_0} \\ &= -\frac{1 - y}{1 - \hat{y}} \cdot \hat{y}(1 - \hat{y}) \quad [\text{Using 6.11}] \\ \frac{\partial}{\partial W_0}((1 - y) \ln (1 - \hat{y})) &= y\hat{y} - \hat{y} \end{aligned} \quad (6.21)$$

Using 6.20 and 6.21 in 6.19 we get

$$\frac{\partial \mathcal{L}}{\partial W_0} = -\frac{1}{m} \sum_{i=1}^m y^{(i)} - \hat{y}^{(i)} \quad (6.22)$$

A similar calculation using 6.12 we get,

$$\frac{\partial \mathcal{L}}{\partial W_j} = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) X_j \quad (6.23)$$

Finally, using 6.22 and 6.23 we get the update rule for gradient descent algorithm in logistic regression with binary cross-entropy.

$$W_0^{(new)} = W_0^{(old)} + \alpha \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) \quad (6.24)$$

$$W_j^{(new)} = W_j^{(old)} + \alpha \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) X_j \quad (6.25)$$

6.2 Example: Logistic Regression with Binary Cross-Entropy

We consider the dataset:

	x_1	x_2	y
	0.5	1.2	0
	1.0	0.8	0
	1.5	1.3	1
	2.0	1.7	1

The logistic regression model is:

$$z = w_0 + w_1 x_1 + w_2 x_2, \quad \hat{y}_i = \frac{1}{1 + e^{-z}}$$

Initial weights: $w_0 = 0$, $w_1 = 0$, $w_2 = 0$. from the table:

Table 6.1: Gradient Computation at Initial Weights

x_1	x_2	y	z	\hat{y}	$\hat{y} - y$	$(\hat{y} - y)x_1$	$(\hat{y} - y)x_2$
0.5	1.2	0	0	0.5	0.5	0.25	0.60
1.0	0.8	0	0	0.5	0.5	0.50	0.40
1.5	1.3	1	0	0.5	-0.5	-0.75	-0.65
2.0	1.7	1	0	0.5	-0.5	-1.00	-0.85
Total:				0.0	-1.00	-0.50	

$$\begin{aligned} \sum_{i=1}^6 (y - \hat{y}) &= 0; \\ \sum_{i=1}^6 (y - \hat{y})x_1 &= -1; \\ \sum_{i=1}^6 (y - \hat{y})x_2 &= -0.50. \end{aligned}$$

The gradients are:

$$\begin{aligned} W_0^{new} &= W_0^{old} + \alpha \frac{\partial \mathcal{L}}{\partial w_0} \\ &= 0 + \alpha \cdot \frac{1}{4} \cdot 0 = [0] \\ W_1^{new} &= W_1^{old} + \alpha \frac{\partial \mathcal{L}}{\partial w_1} \\ &= 0 + \alpha \cdot \frac{1}{4} \cdot (-1.00) = [-.25] \\ W_2^{new} &= W_2^{old} + \alpha \frac{\partial \mathcal{L}}{\partial w_2} \\ &= 0 + \alpha \cdot \frac{1}{4} \cdot (-0.50) = [-.125] \end{aligned}$$

where α is the learning rate.

CHAPTER 7

NEURAL NETWORKS

7.1 Introduction

Neural networks are a foundational component of modern machine learning and artificial intelligence. Inspired by the structure of the human brain, a neural network consists of layers of interconnected units called neurons that process data in stages. These models are particularly effective in tasks like image classification, natural language processing, and regression analysis.

7.2 Perceptron: The Basic Unit

The perceptron is the simplest type of artificial neuron, introduced by Frank Rosenblatt in 1958. It forms the building block of a neural network and functions similarly to linear classifiers that use hard thresholding, as you may recall from the previous chapter.

The linear equation for the weighted sum is given by:

$$z = \sum_{i=1}^n w_i x_i + b, \quad y = f(z) \tag{7.1}$$

Here:

- x_i are the input features
- w_i are the corresponding weights
- b is a bias term
- f is an activation function, a step function in this case.

$$y = f(z) \text{ (step)} = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases} = 0 \tag{7.2}$$

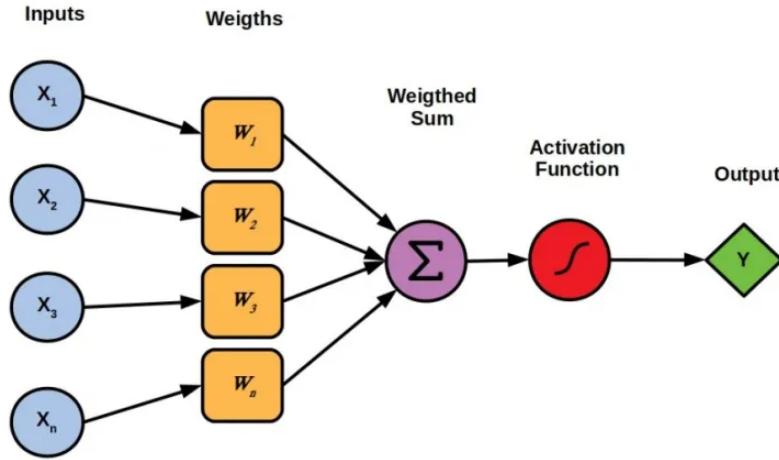


Figure 7.1: A neural network with single neuron perceptron

Image credit: <https://www.linkedin.com/pulse/perceptron-the-basic-building-block-neural-networks-ayush-meharkure-p46if/>

A perceptron takes several input values, applies weights to them, sums them together with a bias term, and then passes the result through an **activation function** — in our case, a step function in the original model. It can be used to separate linearly separable data.

Example: Let $x = [2, 3]$, $w = [0.4, -0.7]$, $b = 0.2$:

$$\begin{aligned} z &= (0.4) \cdot 2 + (-0.7) \cdot 3 + 0.2 \\ &= -1.1 \\ y &= f(z) = 0 \text{ using step function.} \end{aligned}$$

Although simple, the perceptron cannot handle tasks that are not linearly separable (such as the XOR function). To address this limitation, we combine multiple perceptrons in layers to form more powerful models known as multi-layer neural networks.

7.3 Network Layers

7.3.1 Input Layer

The input layer is the first layer in a neural network. It serves as the interface between raw data and the computational structure of the network. Each neuron in this layer represents one feature or attribute of the input data. For instance, in image data, each pixel could be one input node. This layer does not perform any computations beyond passing the inputs to the next layer.

7.3.2 Hidden Layers

Hidden layers are intermediate layers that lie between the input and output layers. These layers are called "hidden" because their outputs are not part of the final network output but are used internally.

Weight Matrix: The full set of weights between two layers can be written as a matrix $W^{(l)}$, where each row corresponds to the weights leading into a single neuron in the current layer. If the previous layer has n neurons and the current layer has m neurons, then $W^{(l)}$ is an $m \times n$ **matrix** and the bias vector will have m entries each associated with one neuron.

$$W^{(l)} = \begin{bmatrix} w_{11}^{(l)} & w_{12}^{(l)} & \cdots & w_{1n}^{(l)} \\ w_{21}^{(l)} & w_{22}^{(l)} & \cdots & w_{2n}^{(l)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1}^{(l)} & w_{m2}^{(l)} & \cdots & w_{mn}^{(l)} \end{bmatrix}_{m \times n}; \quad b^{(l)} = \begin{bmatrix} b_1^{(l)} \\ b_2^{(l)} \\ \vdots \\ b_m^{(l)} \end{bmatrix}_{m \times 1}$$

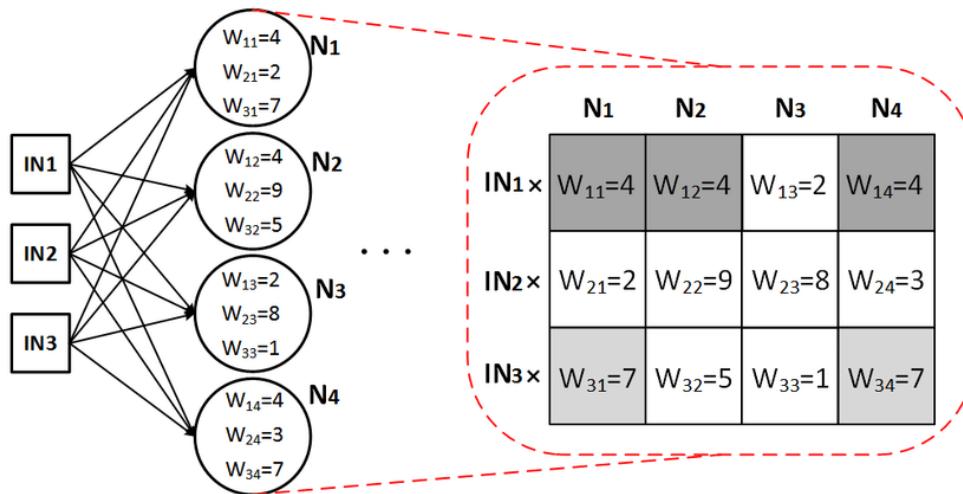


Figure 7.2: Weight Matrix for a hidden layer

Image credit: <https://medium.com/coinmonks/the-mathematics-of-neural-network-60a112dd3e05>

Each neuron in a hidden layer performs the following computation:

$$z_j^{(l)} = \sum_{i=1}^n w_{ji}^{(l)} a_i^{(l-1)} + b_j^{(l)}, \quad a_j^{(l)} = f(z_j^{(l)}) \quad (7.3)$$

Where:

- $a_i^{(l-1)}$ is the activation from neuron i in the previous layer
- $w_{ji}^{(l)}$ is the weight from neuron i in layer $l - 1$ to neuron j in layer l

- $b_j^{(l)}$ is the bias for neuron j in layer l
- f is the activation function applied element-wise

In a vector representation:

$$z^{(l)} = \underbrace{\begin{bmatrix} w_{11}^{(l)} & w_{12}^{(l)} & \cdots & w_{1n}^{(l)} \\ w_{21}^{(l)} & w_{22}^{(l)} & \cdots & w_{2n}^{(l)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1}^{(l)} & w_{m2}^{(l)} & \cdots & w_{mn}^{(l)} \end{bmatrix}}_{W^{(l)} \in \mathbb{R}^{m \times n}} \cdot \underbrace{\begin{bmatrix} a_1^{(l-1)} \\ a_2^{(l-1)} \\ \vdots \\ a_n^{(l-1)} \end{bmatrix}}_{a^{(l-1)} \in \mathbb{R}^{n \times 1}} + \underbrace{\begin{bmatrix} b_1^{(l)} \\ b_2^{(l)} \\ \vdots \\ b_m^{(l)} \end{bmatrix}}_{b^{(l)} \in \mathbb{R}^{m \times 1}} \in \mathbb{R}^{m \times 1} \quad (7.4)$$

more consisely,

$$z^{(l)} = W^{(l)} \cdot a^{(l)} + b^{(l)} \quad (7.5)$$

- $W^{(l)}$: Weight matrix of shape $m \times n$, where each row corresponds to a neuron in the current layer, and each column to a neuron in the previous layer.
- $a^{(l-1)}$: Activation vector from the previous layer, of shape $n \times 1$. For the first layer the activation vector is the original input vector, X with n attributes.
- $b^{(l)}$: Bias vector for the current layer, with one bias per neuron, of shape $m \times 1$.
- $z^{(l)}$: Pre-activation vector of the current layer, of shape $m \times 1$, computed as the weighted sum of previous activations plus bias.

7.3.3 Activation Function

An activation function is a non-linear function applied to the output of each neuron after the weighted sum and bias addition.

$$a^{(l)} = f(z^{(l)}) = \begin{bmatrix} f(z_1^{(l)}) \\ f(z_2^{(l)}) \\ \vdots \\ f(z_m^{(l)}) \end{bmatrix} \in \mathbb{R}^{m \times 1} \quad (7.6)$$

- f : The activation function (e.g., sigmoid, ReLU, tanh), applied element-wise.
- $z^{(l)}$: The pre-activation vector computed from the weighted sum and bias.
- $a^{(l)}$: The output activations from the current layer, used as input to the next layer (or final prediction if it's the output layer).

It introduces non-linearity into the model, allowing it to learn complex patterns. Common activation functions include:

- **Step**: $\text{step}(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases}$

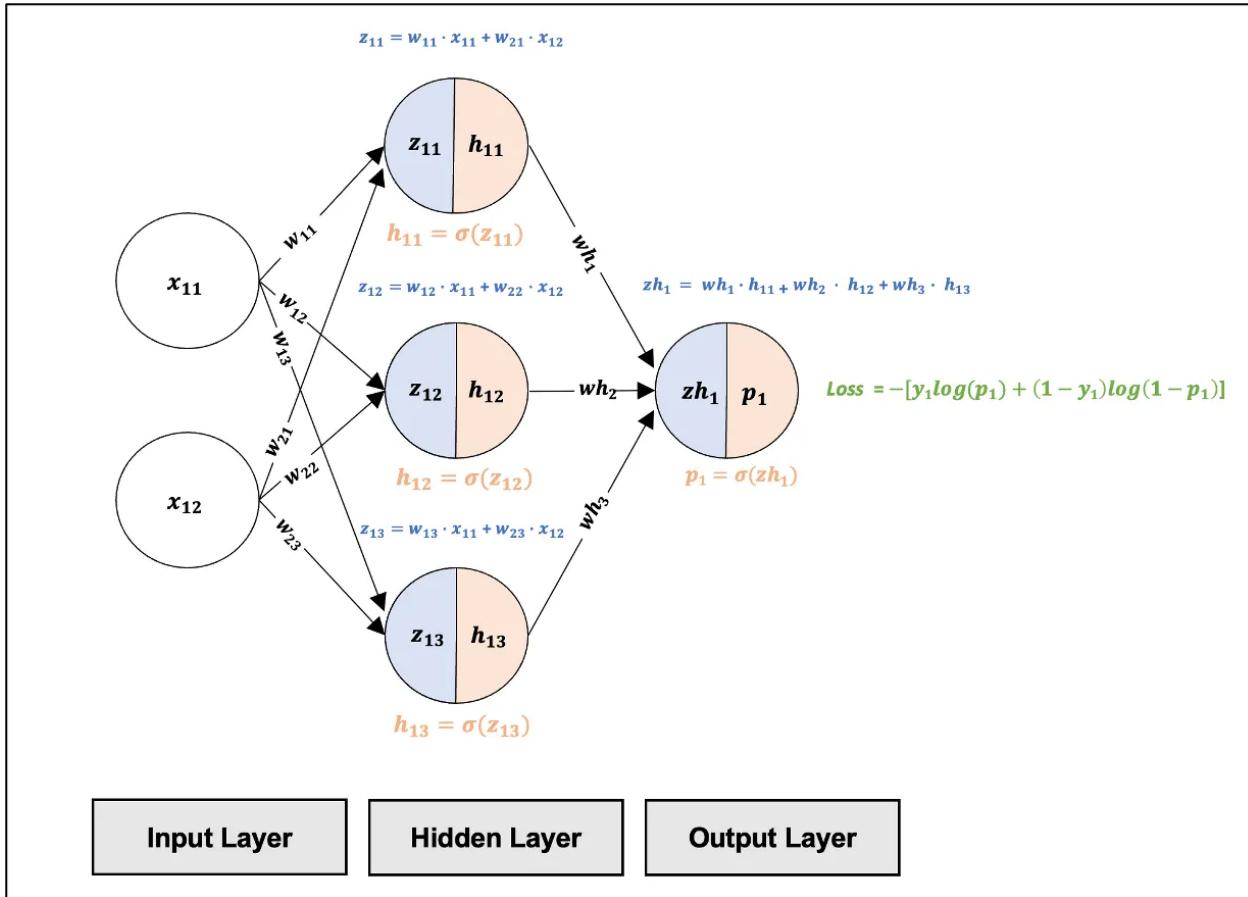


Figure 7.3: Hidden Layer of Neural Network

Image credit: <https://pub.towardsai.net/the-multilayer-perceptron-built-and-implemented-from-scratch-70d6b30f1964>

- **Sigmoid:** $\sigma(z) = \frac{1}{1+e^{-z}}$ maps input to (0,1)
- **Tanh:** $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$.
- **ReLU:** $\text{ReLU}(z) = \max(0, z)$.

Each has different properties affecting training dynamics and convergence.

7.3.4 Output Layer

The output layer produces the final result of the neural network. The type of activation function used in this layer depends on the task:

- **Sigmoid:** Used for binary classification problems, maps output to a probability in (0,1)
- **Softmax:** Used for multi-class classification, gives a probability distribution over classes
- **Identity:** Used in regression tasks, outputs real-valued predictions

7.4 Understanding the Weight Matrix and Activation Function

The weight matrix $W^{(l)}$ plays a crucial role in how information flows and transforms across layers of a neural network. It is responsible for linearly combining the activations from the previous layer and shaping how those activations influence the neurons in the current layer.

For a given layer l , the weight matrix $W^{(l)}$ has dimensions $m \times n$, where:

- n is the number of neurons in the previous layer $l - 1$
- m is the number of neurons in the current layer l

The output of the linear transformation is:

$$z^{(l)} = W^{(l)}a^{(l-1)} + b^{(l)} \quad (7.7)$$

This vector $z^{(l)}$ is then passed through an activation function f to introduce non-linearity:

$$a^{(l)} = f(z^{(l)}) \quad (7.8)$$

Role of the Weight Matrix

- Each row in $W^{(l)}$ contains the weights associated with one neuron in the current layer.
- Each column in $W^{(l)}$ corresponds to the contribution from one neuron in the previous layer.
- The dot product $W^{(l)}a^{(l-1)}$ produces a vector of weighted sums, one per neuron in layer l .

Relation to Activation Functions

The activation function f operates element-wise on the result of the matrix multiplication $z^{(l)}$. Without f , the entire network would be a stack of linear transformations—essentially equivalent to a single matrix multiplication. The use of f makes it possible for the network to model complex non-linear relationships in data.

For example, in the step function:

$$f(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (7.9)$$

The activation function thresholds the result of each neuron's computation, allowing the network to model logical decision boundaries like AND, OR, or XOR (with enough layers).

Thus, the weight matrix defines how signals propagate and interact, and the activation function defines how these signals are interpreted at each neuron.

Detailed Example: Two-Layer Feedforward Network

Input:

$$x = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

Weights and Biases:

$$\begin{aligned} W^{(1)} &= \begin{bmatrix} 0.5 & -0.6 \\ 0.3 & 0.8 \end{bmatrix}, & b^{(1)} &= \begin{bmatrix} 0.1 \\ -0.2 \end{bmatrix} \\ W^{(2)} &= \begin{bmatrix} 1.2 & -0.5 \end{bmatrix}, & b^{(2)} &= 0.3 \end{aligned}$$

Step 1: First Linear Transformation (Hidden Layer)

$$\begin{aligned} z^{(1)} &= W^{(1)}x + b^{(1)} \\ &= \begin{bmatrix} 0.5 \cdot 1 + (-0.6) \cdot 2 + 0.1 \\ 0.3 \cdot 1 + 0.8 \cdot 2 - 0.2 \end{bmatrix} \\ &= \begin{bmatrix} -0.6 \\ 1.7 \end{bmatrix} \end{aligned}$$

Step 2: Activation Function (Step Function)

$$a^{(1)} = f(z^{(1)}) = \text{ReLU}(z^{(1)}) = \begin{bmatrix} 0 \\ 1.7 \end{bmatrix}$$

Step 3: Second Linear Transformation (Output Layer)

$$z^{(2)} = W^{(2)}a^{(1)} + b^{(2)} = 1.2 \cdot 0 + (-0.5) \cdot 1 + 0.3 = -0.2$$

Step 4: Final Activation (Step Function)

$$\hat{y} = f(z^{(2)}) = \text{sigmoid}(-0.2) = \frac{1}{1 + e^{-(-0.2)}}$$

In this example:

- The first layer transformed the input using weights and biases.
- The step activation introduced non-linearity, turning negative values to 0 and non-negatives to 1.
- The second layer used the resulting activations to compute a final output.

7.5 Feedforward Neural Networks

A feedforward neural network (FNN) is the most basic type of artificial neural network where connections between the nodes do not form cycles. The network processes input data layer-by-layer, from input to output, applying transformations at each stage.

Feedforward networks have layers where data flows one-way:

$$\begin{aligned} a^{(1)} &= f^{(1)}(W^{(1)}x + b^{(1)}) \\ a^{(2)} &= f^{(2)}(W^{(2)}a^{(1)} + b^{(2)}) = \hat{y} \end{aligned}$$

Each activation $a^{(l)}$ is passed to the next layer until the final output is produced.

General Structure

Assume a network with L layers (excluding the input layer), where each layer l has $n^{(l)}$ neurons.

Given input $x \in \mathbb{R}^{n^{(0)}}$ (e.g., a feature vector), the feedforward steps are as follows:

Step 1: Weighted Sum at Each Neuron

Each neuron computes a weighted sum of its inputs plus a bias:

$$z^{(l)} = W^{(l)}a^{(l-1)} + b^{(l)} \quad (7.10)$$

Here:

- $W^{(l)} \in \mathbb{R}^{n^{(l)} \times n^{(l-1)}}$ is the weight matrix of layer l
- $a^{(l-1)} \in \mathbb{R}^{n^{(l-1)}}$ is the activation output from the previous layer (or input x if $l = 1$)
- $b^{(l)} \in \mathbb{R}^{n^{(l)}}$ is the bias vector of layer l

Step 2: Activation Function

Each $z^{(l)}$ is passed through a non-linear activation function σ to produce activations for the current layer:

$$a^{(l)} = \sigma(z^{(l)}) \quad (7.11)$$

Common choices for σ include:

- Sigmoid: $\sigma(z) = \frac{1}{1+e^{-z}}$
- ReLU: $\sigma(z) = \max(0, z)$
- Tanh: $\sigma(z) = \tanh(z)$

Final Output Layer

The output layer's activation, $a^{(L)}$, gives the network's prediction:

$$\hat{y} = a^{(L)} = \sigma(z^{(L)}) \quad (7.12)$$

Example 1: Single Layer Network (Result Prediction)

Consider a single-layer network (with sigmoid activation function) used to predict whether a student passes an exam based on two features:

- x_1 : Number of study hours
- x_2 : Number of hours of sleep

Input vector:

$$x = \begin{bmatrix} 5 \\ 7 \end{bmatrix}, \quad W = \begin{bmatrix} 0.4 & 0.6 \end{bmatrix}, \quad b = -3$$

Compute the linear combination:

$$z = Wx + b = 0.4 \cdot 5 + 0.6 \cdot 7 - 3 = 2 + 4.2 - 3 = 3.2$$

Apply the sigmoid activation function:

$$\hat{y} = \sigma(3.2) = \frac{1}{1 + e^{-3.2}} \approx 0.9608$$

The model predicts a 96% chance the student will pass based on study and sleep hours.

Example 2: Single-Layer Network (Spam Classifier)

Let us consider a single layer neural network with sigmoid activation function. Suppose we want to predict whether an email is spam based on two features:

- x_1 : Number of links in the email

- x_2 : Number of exclamation marks

Input vector:

$$x = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$$

Let the weights and bias be:

$$w = [0.7 \ 1.5], \quad b = -2$$

Step 1: Compute the weighted sum:

$$\begin{aligned} z &= 0.7 \cdot 3 + 1.5 \cdot 1 - 2 \\ &= 2.1 + 1.5 - 2 = 1.6 \end{aligned}$$

Step 2: Apply sigmoid function:

$$y = \sigma(z) = \frac{1}{1 + e^{-1.6}} \approx 0.832$$

The probability of spam is approximately 83.2%.

Example 3: 3-Layer Network (Loan Approval) — with 3 Neurons per Layer

Considering a 3-layer neural network where each hidden layer has 3 neurons for the previous example. We are using sigmoid activation function in every layer. We are still using three input features:

- x_1 : Applicant's income (normalized)
- x_2 : Credit score (normalized)
- x_3 : Number of existing loans

Input vector:

$$x = \begin{bmatrix} 0.8 \\ 0.6 \\ 2 \end{bmatrix}$$

Layer 1 (Hidden Layer 1):

$$W^{(1)} = \begin{bmatrix} 0.2 & 0.4 & -0.5 \\ -0.3 & 0.1 & 0.6 \\ 0.5 & -0.2 & 0.3 \end{bmatrix}, \quad b^{(1)} = \begin{bmatrix} 0.1 \\ 0.2 \\ -0.1 \end{bmatrix}$$

Pre-activation vector, $z^{(1)} = W^{(1)}x + b^{(1)}$

$$= \begin{bmatrix} 0.16 + 0.24 - 1.0 + 0.1 \\ -0.24 + 0.06 + 1.2 + 0.2 \\ 0.4 - 0.12 + 0.6 - 0.1 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 1.22 \\ 0.78 \end{bmatrix}$$

Apply activation (sigmoid): $a^{(1)} = \sigma(z^{(1)}) = \begin{bmatrix} 0.3775 \\ 0.772 \\ 0.6859 \end{bmatrix}$

Layer 2 (Hidden Layer 2):

$$W^{(2)} = \begin{bmatrix} 0.3 & 0.7 & -0.2 \\ -0.1 & 0.4 & 0.5 \\ 0.2 & -0.3 & 0.6 \end{bmatrix}, \quad b^{(2)} = \begin{bmatrix} -0.2 \\ 0.3 \\ 0.1 \end{bmatrix}$$

Pre-activation vector: $z^{(2)} = W^{(2)}a^{(1)} + b^{(2)}$

$$= \begin{bmatrix} 0.1133 + 0.5404 - 0.1372 - 0.2 \\ -0.0378 + 0.3088 + 0.343 + 0.3 \\ 0.0755 - 0.2316 + 0.4115 + 0.1 \end{bmatrix} \approx \begin{bmatrix} 0.3165 \\ 0.914 \\ 0.354 \end{bmatrix}$$

Apply activation: $a^{(2)} = \sigma(z^{(2)}) = \begin{bmatrix} 0.5785 \\ 0.7138 \\ 0.5875 \end{bmatrix}$

Output Layer:

$$W^{(3)} = [1.0 \quad -1.5 \quad 0.6], b^{(3)} = 0.1$$

Pre-activation vector, $z^{(3)} = 1.0 \cdot 0.5785 - 1.5 \cdot 0.7138 + 0.6 \cdot 0.5875 + 0.1$
 $= -0.0397$

Apply sigmoid activation: $\hat{y} = \sigma(z^{(3)}) = \frac{1}{1 + e^{-0.0397}} \approx 0.4901$

The model predicts a 49.01% probability of loan approval.

7.6 Activation and Loss Functions

In this section, we define different activation function and the equation of Loss for them.

7.6.1 Sigmoid (Binary Classification)

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

$$\mathcal{L}_{\text{BCE}} = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})]$$

7.6.2 Softmax (Multiclass Classification)

$$\hat{y}_j = \frac{e^{z_j}}{\sum_k e^{z_k}}$$

$$\mathcal{L}_{\text{CCE}} = - \sum_j y_j \log(\hat{y}_j)$$

7.6.3 Linear (Regression)

$$\hat{y} = W \cdot X + b$$

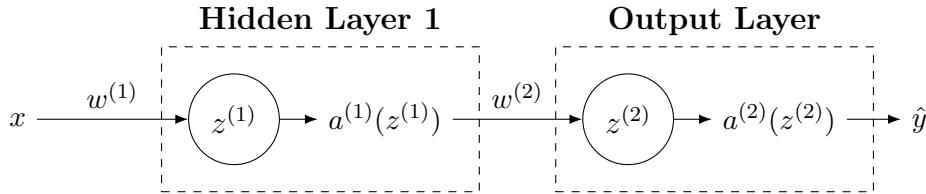
$$\mathcal{L}_{\text{MSE}} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

7.7 Training: Gradient Descent and Backpropagation

7.7.1 Backpropagation in Feedforward Neural Networks

Backpropagation is the core algorithm used to train feedforward neural networks by updating the weights and biases to minimize a loss function. It computes the gradient of the loss function with respect to each parameter using the chain rule of calculus.

Assume a network with L layers, where each layer l has activations $a^{(l)}$, weights $W^{(l)}$, biases $b^{(l)}$, and pre-activations $z^{(l)}$.



For simplicity, let us consider $L = 2$, with activation functions, $a^{(1)}$ for the hidden layer and $a^{(2)}$ for the output layer. The weight matrix for the hidden layer is $w^{(1)}$ and for the output layer is $w^{(2)}$.

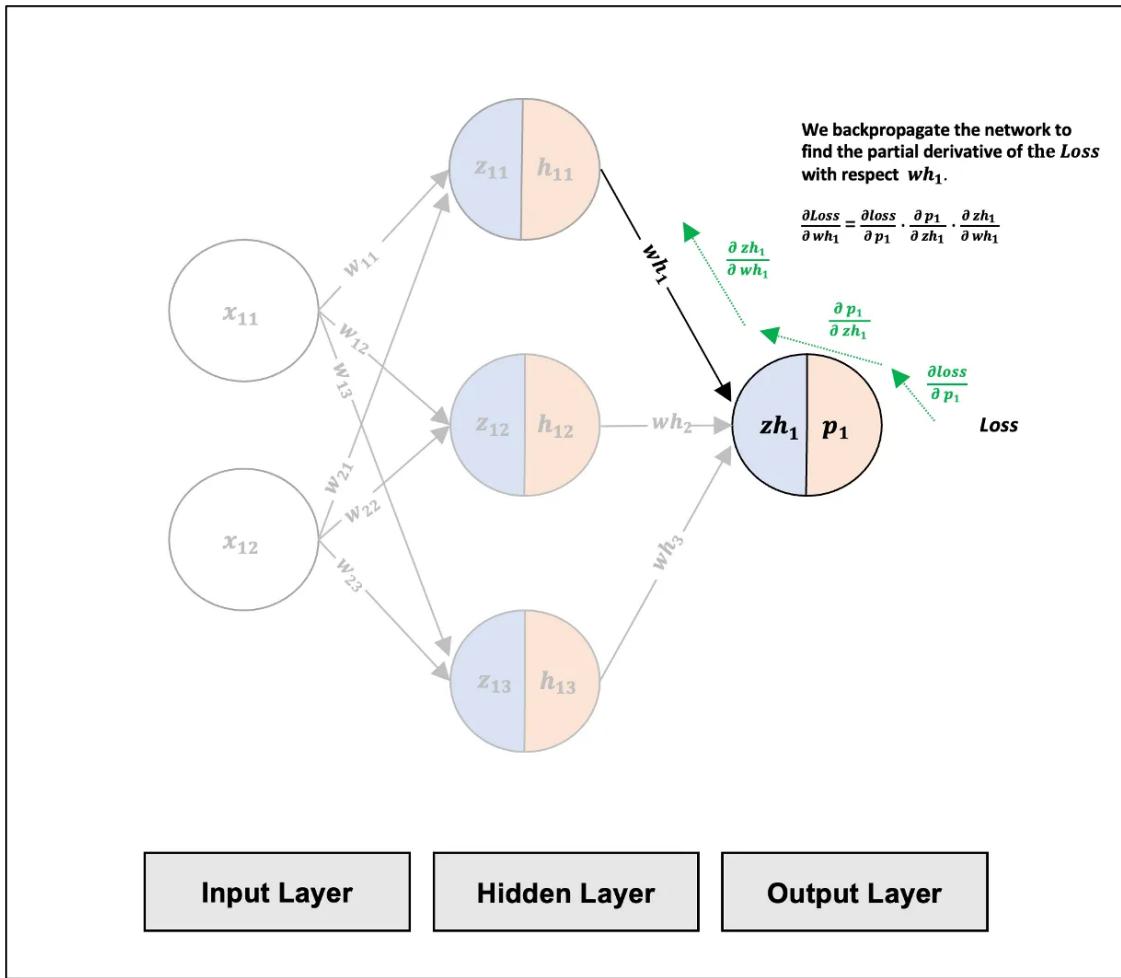


Figure 7.4: Image credit: <https://pub.towardsai.net/the-multilayer-perceptron-built-and-implemented-from-scratch-70d6b30f1964>

We have to update the weight vectors as

$$\begin{aligned}
 w^{(2)} &= w^{(2)} + \alpha \frac{\partial \mathcal{L}(y)}{\partial w^{(2)}} \\
 b^{(2)} &= b^{(2)} + \frac{\partial \mathcal{L}}{\partial b^{(2)}} \\
 w^{(2)} &= w^{(2)} - \alpha \frac{\partial \mathcal{L}(y)}{\partial w^{(2)}} \\
 b^{(2)} &= b^{(2)} + \frac{\partial \mathcal{L}}{\partial b^{(2)}}
 \end{aligned}$$

7.7.2 Computing the Gradient at the Output Layer

At the output layer L , the computation of $\frac{\partial \mathcal{L}}{\partial W^{(L)}}$ is similar to the gradient calculation used in the basic gradient descent algorithm.

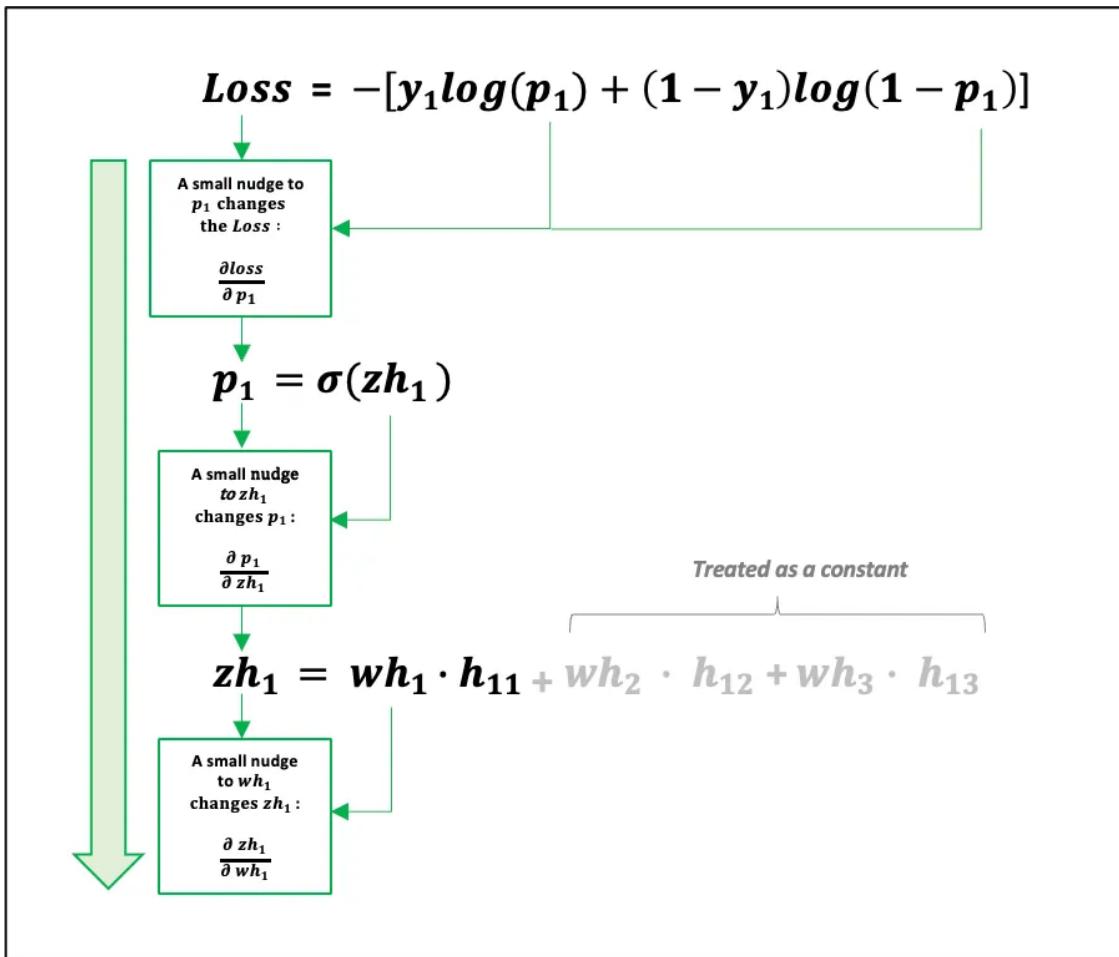


Figure 7.5: Image credit: <https://pub.towardsai.net/the-multilayer-perceptron-built-and-implemented-from-scratch-70d6b30f1964>

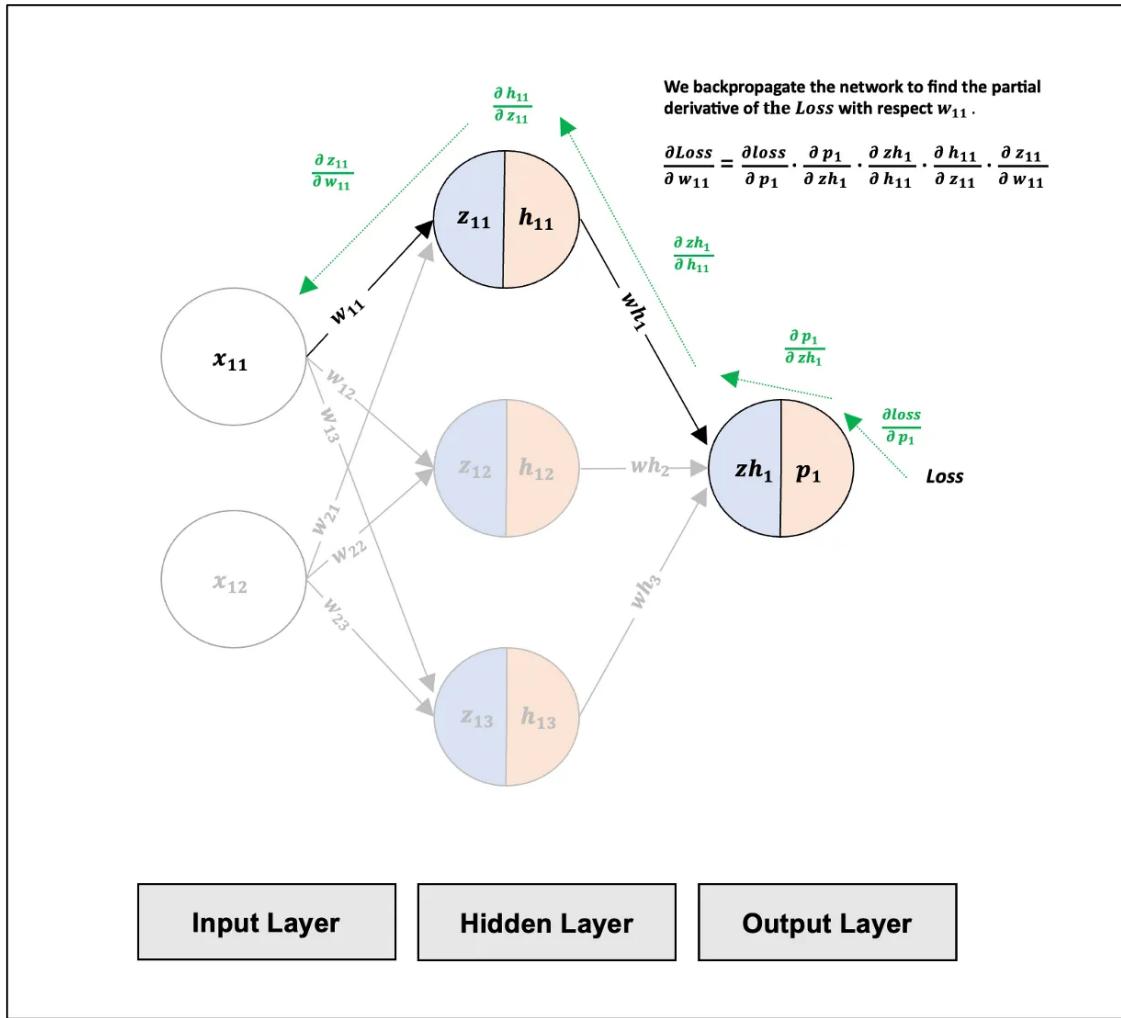


Figure 7.6: Image Credit: <https://pub.towardsai.net/the-multilayer-perceptron-built-and-implemented-from-scratch-70d6b30f1964>

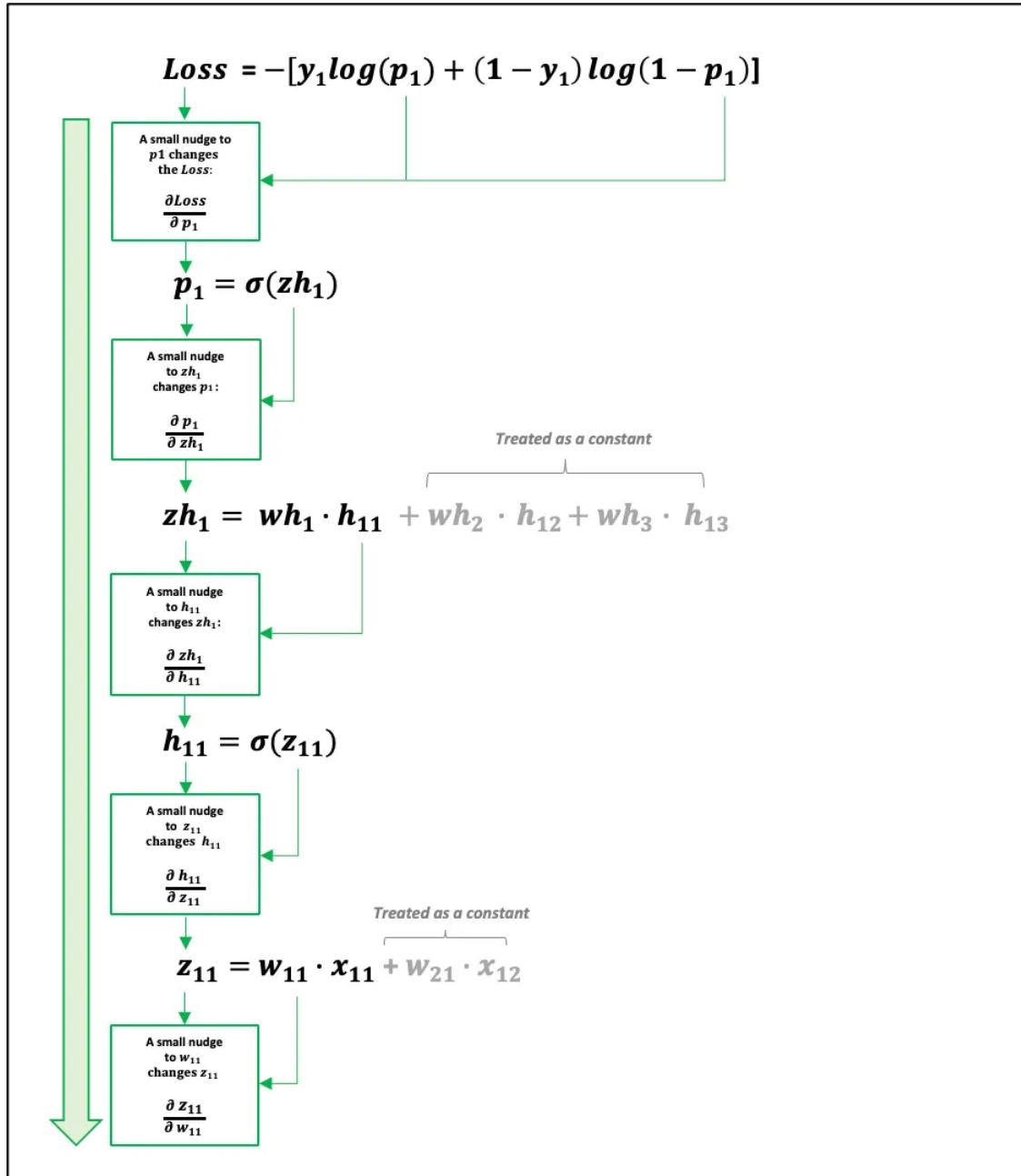


Figure 7.7: Image credit: <https://pub.towardsai.net/the-multilayer-perceptron-built-and-implemented-from-scratch-70d6b30f1964>

The output layer has a single neuron. The linear equation to this neuron is given by:

$$Z_1^{(L)} = W^{(L)} \cdot a^{(L-1)} + b^{(L)} = w_1^{(L)} a_1^{(L-1)} + w_2^{(L)} a_2^{(L-1)} + \cdots + w_p^{(L)} a_p^{(L-1)} + b^{(L)} \quad (7.13)$$

Here, p is the number of neurons in the previous layer (layer $L - 1$), and each $a_i^{(L-1)}$ is the activation output from neuron i in that layer.

$$\frac{\partial \mathcal{L}}{\partial w_i^{(L)}} = \frac{\partial \mathcal{L}}{\partial a^L} \frac{\partial a^L}{\partial Z_1} \frac{\partial Z_1}{\partial w_i}, \text{ using chain rule. Recall, } \hat{y} = a^{(L)}(Z_1),$$

Recall that: $z_1^{(L)} = \sum_i w_{1i}^{(L)} a_i^{(1)} + b_j^{(2)}$, so,

$$\frac{\partial z_1^{(L)}}{\partial w_{1i}^{(L)}} = a_i^{(L-1)}.$$

Hence,

$$\frac{\partial \mathcal{L}}{\partial w_i^{(L)}} = \frac{\partial \mathcal{L}}{\partial a^L} \frac{\partial a^L}{\partial Z_1} a_i^{(L-1)}$$

and by using, $\frac{\partial Z_1}{\partial b_i^{(L)}} = 1$

$$\frac{\partial \mathcal{L}}{\partial b_i^{(L)}} = \frac{\partial \mathcal{L}}{\partial a^L} \frac{\partial a^L}{\partial Z_1} \cdot 1$$

Generally, this can be written as:

$$\frac{\partial \mathcal{L}}{\partial W^{(L)}} = \frac{\partial \mathcal{L}}{\partial a^L} \frac{\partial a^L}{\partial Z_1} \cdot a^{(L-1)}; \quad (7.14)$$

$$\frac{\partial \mathcal{L}}{\partial b^{(L)}} = \frac{\partial \mathcal{L}}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial Z_1}. \quad (7.15)$$

Finally, the update rules at the output layer, L can be written as:

$$W^{(L)} \leftarrow W^{(L)} - \alpha \frac{\partial \mathcal{L}}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial Z_1} \cdot (a^{(L-1)})^\top \quad (7.16)$$

$$b^{(L)} \leftarrow b^{(L)} - \alpha \frac{\partial \mathcal{L}}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial Z_1}. \quad (7.17)$$

Note: Here, The activation vector $a^{(L)}$ is transposed to turn it into a row vector so it can be multiplied with the error. This gives a full matrix of weight updates, matching the shape of $W^{(L)}$.

7.7.3 Understanding the Derivative of Activation with Respect to Pre-Activation

The derivative of the activation $a^{(i)}$ with respect to the input $Z^{(i)}$ at layer i is given by:

$$\frac{da^{(i)}}{dZ^{(i)}} = \phi'(Z^{(i)}) \quad (7.18)$$

Examples:

- **Sigmoid activation:**

$$\phi(z) = \frac{1}{1 + e^{-z}} \Rightarrow \phi'(z) = \phi(z)(1 - \phi(z))$$

So,

$$\frac{da^{(i)}}{dZ^{(i)}} = a^{(i)}(1 - a^{(i)}) \quad (7.19)$$

- **ReLU activation:**

$$\phi(z) = \max(0, z) \Rightarrow \phi'(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases}$$

So,

$$\frac{da^{(i)}}{dZ^{(i)}} = \begin{cases} 1 & \text{if } Z^{(i)} > 0 \\ 0 & \text{otherwise} \end{cases} \quad (7.20)$$

7.7.4 Understanding the Derivative of Loss with Respect to Output Activation

$\frac{\partial \mathcal{L}}{\partial a^{(L)}}$ represents the derivative of the loss function \mathcal{L} with respect to the activation output of the final layer — that is, how the loss changes as the network's prediction changes. The exact form depends on the loss function used. Let $\hat{y} = a^{(L)}$ be the activation at the output layer and y be the true label.

Mean Squared Error (MSE) Loss:

$$\mathcal{L} = \frac{1}{2}(\hat{y} - y)^2 \quad (7.21)$$

$$\frac{d\mathcal{L}}{da^{(L)}} = \hat{y} - y \quad (7.22)$$

Binary Cross-Entropy (BCE) Loss

$$\mathcal{L} = -(y \log \hat{y} + (1 - y) \log(1 - \hat{y})) \quad (7.23)$$

$$\frac{d\mathcal{L}}{da^{(L)}} = -\left(\frac{y}{\hat{y}} - \frac{1 - y}{1 - \hat{y}}\right) \quad (7.24)$$

7.7.5 Defining Error, $\frac{d\mathcal{L}}{dZ^{(l)}}$

Let us define,

$$\frac{\partial \mathcal{L}}{\partial Z^{(l)}} = \frac{\partial \mathcal{L}}{\partial a^{(l)}} \frac{\partial a^{(l)}}{\partial Z^{(l)}} = D^{(l)} \quad (7.25)$$

where, l is the layer.

The term $\frac{d\mathcal{L}}{dZ^{(l)}} = D_i^{(l)}$ represents how the loss \mathcal{L} changes with respect to the $Z_i^{(l)}$ of the i 'th neuron at layer l .

The Error $D^{(l)}$ combines:

- The sensitivity of the loss to the neuron's activation: $\frac{\partial \mathcal{L}}{\partial a^l}$
- The sensitivity of the neuron's activation to its Z : $\frac{da}{dZ}$

7.7.6 Update Rule at the output layer

For the 2-layer feedforward neural network,

$$W^{(2)} \leftarrow W^{(2)} - \alpha D^{(2)} (a^{(L-1)})^\top \quad (7.26)$$

$$b^{(2)} \leftarrow b^{(2)} - \alpha D^{(2)} \quad (7.27)$$

7.7.7 Computing gradient at the Hidden Layer ($l = 1$)

Let $W^{(1)}$ be the weight matrix connecting the input layer to the hidden layer and $b^{(1)}$ the bias vector. Each weight $w_{ij}^{(1)}$ connects input neuron j to hidden neuron i , and let $a_j^{(0)} = x_j$ be the input feature at position j .

The linear equation Z_i at the hidden neuron i is written as:

$$Z_i^{(1)} = \sum_j w_{ij}^{(1)} a_j^{(0)} + b_i^{(1)}.$$

Generally, it can be written as:

$$Z^{(1)} = W^{(1)} \cdot x + b^{(1)}.$$

Using the chain rule of differentiation, the gradient of loss becomes

$$\frac{\partial \mathcal{L}}{\partial W^{(1)}} = \frac{\partial \mathcal{L}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial Z^{(2)}} \frac{\partial Z^{(2)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial Z^{(1)}} \frac{\partial Z^{(1)}}{\partial W^{(1)}}$$

and

$$\frac{\partial \mathcal{L}}{\partial b^{(1)}} = \frac{\partial \mathcal{L}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial Z^{(2)}} \frac{\partial Z^{(2)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial Z^{(1)}} \frac{\partial Z^{(1)}}{\partial b^{(1)}}.$$

Recall that $\frac{\partial \mathcal{L}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial Z^{(2)}} = D^2$ is already computed at the output layer. Now to compute $\frac{\partial Z^{(2)}}{\partial a^{(1)}}$, you may recall that the output layer contains only one neuron and the linear equation at that layer is written as:

$$Z_1^{(2)} = W^{(2)} \cdot a^{(1)} = \sum_j W_{1j}^{(2)} a_j^{(1)} + b_1^{(2)}$$

Now,

$$\frac{\partial Z_1^{(2)}}{\partial a_k^{(1)}} = W_{1k}^{(2)}. \quad (7.28)$$

This can be written more generally as:

$$\frac{\partial Z_p^{(2)}}{\partial a_k^{(1)}} = W_{pk}^{(2)} \quad (7.29)$$

$$\frac{\partial Z^{(2)}}{\partial a^{(1)}} = W^{(2)}. \quad (7.30)$$

The derivative of the activation $a^{(i)}$ with respect to the input $Z^{(i)}$ at layer i is given by:

$$\frac{da^{(i)}}{dZ^{(i)}} = \phi'(Z^{(i)}) \quad (7.31)$$

So, this will result to a column vector each containing the value of $\phi'(Z_i^1)$. Lastly,

$$\frac{\partial Z^{(l)}}{\partial W^{(l)}} = a^{(l-1)}. \quad (7.32)$$

For the hidden layer 1,

$$\frac{\partial Z^{(1)}}{\partial W^{(1)}} = a^{(0)} = x.$$

Putting it all together,

$$\frac{\partial \mathcal{L}}{\partial W^{(1)}} = (D^2 \cdot (W^{(2)})^\top) \frac{\partial a^{(1)}}{\partial Z^{(1)}} \odot (a^0)^\top \quad (7.33)$$

$$= (D^2 \cdot (W^{(2)})^\top) \frac{\partial a^{(1)}}{\partial Z^{(1)}} \odot x^\top \quad (7.34)$$

and,

$$\frac{\partial \mathcal{L}}{\partial b^{(1)}} = (D^2 \cdot (W^{(2)})^\top) \frac{\partial a^{(1)}}{\partial Z^{(1)}}. \quad (7.35)$$

Gradients of specific weight or bias entries at layer 1

Gradients for specific weight or bias entries allows us to compute how each individual weight and bias in the first layer should be adjusted to reduce the network's loss:

$$\frac{\partial \mathcal{L}}{\partial W_{ik}^{(1)}} = (D_1^{(2)} \cdot W_{1i}^{(2)}) \frac{\partial a_i^{(1)}}{\partial Z_i^{(1)}} a_k^{(0)} \quad (7.36)$$

$$= (D_1^{(2)} \cdot W_{1k}^{(2)}) \phi'(Z_i^{(1)}) x_k \quad (7.37)$$

and,

$$\frac{\partial \mathcal{L}}{\partial b_i^{(1)}} = (D_1^{(2)} \cdot W_{1i}^{(2)}) \phi'(Z_i^{(1)}) \quad (7.38)$$

Here:

- $W_{ik}^{(1)}$ is the weight connecting the k th input x_k to the i th neuron in the first hidden layer.
- $\frac{\partial \mathcal{L}}{\partial W_{ik}^{(1)}}$ shows how this weight affects the total loss \mathcal{L} .
- $D_1^{(2)} \cdot W_{1i}^{(2)}$ represents how much the i th neuron in the first layer influences the output layer's error.
- $\phi'(Z_i^{(1)})$ is the derivative of the activation function at neuron i in layer 1.
- x_k is the k th input value.

Back Propagation of Loss

Notice that we computed the error term $\frac{\partial \mathcal{L}}{\partial Z^l}$ at the hidden layer ($l = 1$) using the error term and weight matrix from the output layer ($l = 2$):

$$D^{(1)} = D^{(2)} \cdot (W^{(2)})^\top \quad (7.39)$$

This shows that, for any number of layers, the error at a given layer can be computed from the error and weights of the next layer. This process of sending the error backward through the network is called backpropagation.

$$D^{(l)} = D^{(l+1)} \cdot (W^{(l+1)})^\top \cdot \frac{\partial a^{(l)}}{\partial z^{(l)}} \quad (7.40)$$

Computing the error contributions of a specific neuron

The error contribution by the k 'th neuron at the l 'th layer can be computed as

$$D_k^{(l)} = \left(\sum_p D_p^{(l+1)} W_{pk}^{(l+1)} \right) \phi'_k \quad (7.41)$$

Here,

- The term $\left(\sum_p D_p^{(l+1)} W_{pk}^{(l+1)} \right)$ is the backpropagation term. It represents the total contribution of the k th neuron in layer l to the errors in all neurons p in the next layer $(l + 1)$, weighted by the corresponding weights $W_{pk}^{(l+1)}$.
- $\phi' = \frac{\partial a_k^{(l)}}{\partial z_k^{(l)}}$ is the derivative of the activation function at the k 'th neuron of layer- l .

Generalized Rules for Gradient Computation

$$\frac{\partial \mathcal{L}}{\partial w^{(l)}} = \left((W^{(l+1)})^T \delta^{(l+1)} \right) \odot \frac{\partial a^{(l)}}{\partial z^{(l)}} a^{(l-1)} \quad (7.42)$$

$$\frac{\partial \mathcal{L}}{\partial b^{(l)}} = \left((W^{(l+1)})^T \delta^{(l+1)} \right) \odot \frac{\partial a^{(l)}}{\partial z^{(l)}} \quad (7.43)$$

Generalized Rules for Gradient Computation of Specific Entries

$$\frac{\partial \mathcal{L}}{\partial w_{ik}^{(l)}} = D_i^{(l)} a_k^{(l-1)} \quad (7.44)$$

$$\frac{\partial \mathcal{L}}{\partial b_{(l)}} = D^l \quad (7.45)$$

Generalized Update Rules for Weights and Bias

$$W_{ik}^{(l)} \leftarrow W_{ik}^{(l)} - \alpha D_i^{(l)} a_k^{(l-1)}, \quad (7.46)$$

and,

$$b_i^{(l)} \leftarrow b_i^{(l)} - \alpha D_i^{(l)} a_k^{(l-1)} \quad (7.47)$$

7.7.8 Simplified Example: Backpropagation in 2-layer Feed-forward Neural Network

Inputs:

$$x = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad y = 1, \quad \alpha = 0.1 \quad (7.48)$$

Initial Weights and Biases:

$$W^{(1)} = \begin{bmatrix} 0.2 & -0.3 \\ -0.1 & 0.4 \end{bmatrix}, \quad b^{(1)} = \begin{bmatrix} 0.1 \\ 0.05 \end{bmatrix}$$

$$W^{(2)} = [0.3 \quad -0.2], \quad b^{(2)} = 0.2$$

Step 1: Forward Pass

Hidden layer:

$$z^{(1)} = W^{(1)}x + b^{(1)} = \begin{bmatrix} 0.2 \cdot 1 + (-0.3) \cdot 2 + 0.1 \\ -0.1 \cdot 1 + 0.4 \cdot 2 + 0.05 \end{bmatrix} = \begin{bmatrix} -0.3 \\ 0.75 \end{bmatrix} \quad (7.49)$$

$$a^{(1)} = \sigma(z^{(1)}) = \begin{bmatrix} \frac{1}{1+e^{0.3}} \\ \frac{1}{1+e^{-0.75}} \end{bmatrix} \approx \begin{bmatrix} 0.4256 \\ 0.6792 \end{bmatrix} \quad (7.50)$$

Step 2: Output Layer Backward Pass

$$\frac{d\mathcal{L}}{da^{(2)}} = \hat{y} - y = 0.5319 - 1 = -0.4681 \quad (7.51)$$

$$\frac{da^{(2)}}{dZ^{(2)}} = \hat{y}(1 - \hat{y}) \approx 0.5319(1 - 0.5319) \approx 0.2490 \quad (7.52)$$

$$D^{(2)} = \frac{d\mathcal{L}}{dZ^{(2)}} = -0.4681 \cdot 0.2490 \approx -0.1166 \quad (7.53)$$

Step 3: Backpropagate to Hidden Neurons

Hidden Layer 1:

Error term for neuron 1:

$$D_1^{(1)} = D_1^{(2)} \cdot w_{11}^{(2)} \cdot a_1^{(1)}(1 - a_1^{(1)}) = -0.1166 \cdot 0.3 \cdot 0.4256(1 - 0.4256) \approx -0.0085 \quad (7.54)$$

Error term for neuron 2:

$$D_2^{(1)} = D_1^{(2)} \cdot w_{12}^{(2)} \cdot a_2^{(1)}(1 - a_2^{(1)}) = -0.1166 \cdot (-0.2) \cdot 0.6792(1 - 0.6792) \approx 0.0051 \quad (7.55)$$

Step 4: Update Hidden Layer Weights**Hidden Layer 1:****Neuron 1:**

$$w_{11}^{(1)} \leftarrow w_{11} - \alpha D_1^{(1)} x_1 \quad (7.56)$$

$$= 0.2 - 0.1 \cdot (-0.0085) \cdot 1 \quad (7.57)$$

$$= 0.20085 \quad (7.58)$$

$$w_{12}^{(1)} \leftarrow w_{12} - \alpha D_1^{(1)} x_2 \quad (7.59)$$

$$= -0.3 - 0.1 \cdot (-0.0085) \cdot 2 \quad (7.60)$$

$$= -0.2983 \quad (7.61)$$

$$b_1^{(1)} \leftarrow b_1^1 - \alpha D_1^{(1)} \quad (7.62)$$

$$= 0.1 - 0.1 \cdot (-0.0085) \quad (7.63)$$

$$= 0.10085 \quad (7.64)$$

Neuron 2:

$$w_{21}^{(1)} \leftarrow w_{21} - \alpha D_2^{(1)} x_1 \quad (7.65)$$

$$= -0.1 - 0.1 \cdot 0.0051 \cdot 1 = -0.10051 \quad (7.66)$$

$$w_{22}^{(1)} \leftarrow w_{22} - \alpha D_2^{(1)} x_2 \quad (7.67)$$

$$= 0.4 - 0.1 \cdot 0.0051 \cdot 2 = 0.39898 \quad (7.68)$$

$$b_2^{(1)} \leftarrow b_2 - \alpha D_2^{(1)} \quad (7.69)$$

$$= 0.05 - 0.1 \cdot 0.0051 = 0.04949 \quad (7.70)$$

7.7.9 Example: Backpropagation on a 3-Layer Neural Network**Network Structure**

- Input Layer: $x_1 = 1.0, x_2 = 2.0$
- Hidden Layer 1 (3 neurons): sigmoid activation
- Hidden Layer 2 (2 neurons): sigmoid activation
- Output Layer (1 neuron): sigmoid activation

Randomly Initialized Weights and Biases**Layer 1 (Input to Hidden 1):**

$$W^{(1)} = \begin{bmatrix} 0.1 & -0.2 \\ 0.4 & 0.3 \\ -0.5 & 0.2 \end{bmatrix}, \quad b^{(1)} = \begin{bmatrix} 0.0 \\ 0.1 \\ -0.1 \end{bmatrix}$$

Layer 2 (Hidden 1 to Hidden 2):

$$W^{(2)} = \begin{bmatrix} 0.2 & -0.1 & 0.3 \\ -0.4 & 0.2 & 0.1 \end{bmatrix}, \quad b^{(2)} = \begin{bmatrix} 0.05 \\ -0.05 \end{bmatrix}$$

Layer 3 (Hidden 2 to Output):

$$W^{(3)} = [0.3 \quad -0.2], \quad b^{(3)} = 0.1$$

Step 1: Forward Pass

Compute pre-activations and activations:

$$\begin{aligned} z^{(1)} &= W^{(1)} \cdot x + b^{(1)} = \begin{bmatrix} 0.1(1) + (-0.2)(2) \\ 0.4(1) + 0.3(2) \\ -0.5(1) + 0.2(2) \end{bmatrix} + \begin{bmatrix} 0 \\ 0.1 \\ -0.1 \end{bmatrix} = \begin{bmatrix} -0.3 \\ 1.1 \\ -0.2 \end{bmatrix} \\ a^{(1)} &= \sigma(z^{(1)}) = \begin{bmatrix} \sigma(-0.3) \\ \sigma(1.1) \\ \sigma(-0.2) \end{bmatrix} = \begin{bmatrix} 0.4256 \\ 0.7503 \\ 0.4502 \end{bmatrix} \\ z^{(2)} &= W^{(2)}a^{(1)} + b^{(2)} = \begin{bmatrix} 0.2 & -0.1 & 0.3 \\ -0.4 & 0.2 & 0.1 \end{bmatrix} \begin{bmatrix} 0.4256 \\ 0.7503 \\ 0.4502 \end{bmatrix} + \begin{bmatrix} 0.05 \\ -0.05 \end{bmatrix} = \begin{bmatrix} 0.1734 \\ -0.1800 \end{bmatrix} \\ a^{(2)} &= \sigma(z^{(2)}) = \begin{bmatrix} 0.5432 \\ 0.4551 \end{bmatrix} \\ z^{(3)} &= W^{(3)}a^{(2)} + b^{(3)} = [0.3, -0.2] \begin{bmatrix} 0.5432 \\ 0.4551 \end{bmatrix} + 0.1 = 0.1629 \\ \hat{y} &= a^{(3)} = \sigma(0.1629) = 0.5406 \end{aligned}$$

Step 2: Backward Pass

Output layer delta:

$$D^{(3)} = \frac{\partial \mathcal{L}}{\partial a^{(3)}} \cdot \sigma'(z^{(3)}) = (\hat{y} - y) \cdot \hat{y}(1 - \hat{y}) = (0.5406 - 1)(0.5406)(0.4594) = -0.1139$$

Layer 2 delta using:

$$D_k^{(2)} = \left(\sum_p D_p^{(3)} W_{pk}^{(3)} \right) \cdot \sigma'(z_k^{(2)})$$

Since $D^{(3)} = -0.1139$, $W^{(3)} = [0.3, -0.2]$

$$D_1^{(2)} = (-0.1139)(0.3) \cdot \sigma'(0.1734) = -0.03417 \cdot (0.5432)(0.4568) = -0.0085$$

$$D_2^{(2)} = (-0.1139)(-0.2) \cdot \sigma'(-0.1800) = 0.02278 \cdot (0.4551)(0.5449) = 0.0057$$

Layer 1 delta:

$$D_k^{(1)} = \left(\sum_p D_p^{(2)} W_{pk}^{(2)} \right) \cdot \sigma'(z_k^{(1)})$$

$$\begin{aligned} D_1^{(1)} &= (-0.0085 \cdot 0.2 + 0.0057 \cdot (-0.4)) \cdot \sigma'(-0.3) \\ &= (-0.0017 - 0.00228) \cdot 0.2445 \\ &= -0.00097 \end{aligned}$$

$$\begin{aligned} D_2^{(1)} &= (-0.0085 \cdot (-0.1) + 0.0057 \cdot 0.2) \cdot \sigma'(1.1) \\ &= (0.00085 + 0.00114) \cdot 0.1879 \\ &= 0.00037 \end{aligned}$$

$$\begin{aligned} D_3^{(1)} &= (-0.0085 \cdot 0.3 + 0.0057 \cdot 0.1) \cdot \sigma'(-0.2) \\ &= (-0.00255 + 0.00057) \cdot 0.2475 \\ &= -0.00049 \end{aligned}$$

Step 3: Gradient Update for Weights and Biases

Using:

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^{(l)}} = -\alpha D_i^{(l)} \cdot a_j^{(l-1)} \quad \frac{\partial \mathcal{L}}{\partial b_i^{(l)}} = -\alpha D_i^{(l)}$$

Assume learning rate $\alpha = 0.1$, then for example:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W_{1,1}^{(1)}} &= -0.1 \cdot (-0.00097) \cdot 1.0 = 0.000097 \\ \frac{\partial \mathcal{L}}{\partial W_{1,2}^{(1)}} &= -0.1 \cdot (-0.00097) \cdot 2.0 = 0.000194 \\ \frac{\partial \mathcal{L}}{\partial b_1^{(1)}} &= -0.1 \cdot (-0.00097) = 0.000097 \end{aligned}$$

Similarly, compute updates for all weights and biases.

What Happens After We Update the Weights?

Once the weights and biases are updated using backpropagation, they are used in the next iteration of training. Here's what happens:

- In the **next forward pass**, the updated weights and biases are used to make a new prediction.

- If the prediction is still incorrect, **new gradients are computed** based on the updated prediction, and weights are further updated.
- This cycle repeats over many training examples and epochs until the network **converges to a solution** that minimizes the loss.
- Once training is complete, the final set of weights is **frozen and used for inference** on unseen data.

Backpropagation trains a neural network by making small improvements to its weights with each example it sees. These updates accumulate over time, gradually improving the network's predictions.

Training a Neural Network: Step-by-Step Procedure

Training a neural network involves iteratively adjusting its parameters to minimize prediction error on a dataset. The process uses feedforward computation, loss evaluation, backpropagation, and weight updates. Below are the key steps:

Step 1: Initialize Parameters

- Randomly initialize the weights $W^{(l)}$ and biases $b^{(l)}$ for each layer l .
- Choose a learning rate η , number of epochs, and batch size (if applicable).

Step 2: Forward Pass

- For each input x , compute the outputs of each layer using:

$$z^{(l)} = W^{(l)}a^{(l-1)} + b^{(l)}, \quad a^{(l)} = \sigma(z^{(l)})$$

- The final activation $a^{(L)}$ is the predicted output \hat{y} .

Step 3: Compute Loss

- Compare the predicted output \hat{y} with the true label y using a suitable loss function, e.g., binary cross-entropy for classification:

$$\mathcal{L} = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})]$$

Step 4: Backpropagation

- Compute gradients of the loss w.r.t. output layer and recursively propagate errors backward:

$$\begin{aligned}\delta^{(L)} &= (\hat{y} - y) \odot \sigma'(z^{(L)}) \\ \delta^{(l)} &= (W^{(l+1)})^T \delta^{(l+1)} \odot \sigma'(z^{(l)})\end{aligned}$$

- Calculate gradients of weights and biases:

$$\frac{\partial \mathcal{L}}{\partial W^{(l)}} = \delta^{(l)} (a^{(l-1)})^T, \quad \frac{\partial \mathcal{L}}{\partial b^{(l)}} = \delta^{(l)}$$

Step 5: Gradient Descent at Each Layer

- Use the computed gradients to update the weights and biases layer by layer:

For each layer $l = 1, 2, \dots, L$:

$$W^{(l)} \leftarrow W^{(l)} - \eta \cdot \frac{\partial \mathcal{L}}{\partial W^{(l)}}, \quad b^{(l)} \leftarrow b^{(l)} - \eta \cdot \frac{\partial \mathcal{L}}{\partial b^{(l)}}$$

- This ensures that each layer's parameters are adjusted to reduce the error based on their contribution.

Step 6: Repeat for All Examples and Epochs

- Repeat steps 2 to 5 for each training example or mini-batch.
- After each pass through the entire dataset, increment the epoch counter.
- Continue until the loss converges or a stopping condition is met.