# Algorithms and Machine Learning

**Comprehensive Study Notes**

Complete Coverage of All Units with Visual Diagrams

# UNIT I: Introduction to Algorithms and Machine Learning

## 1. Introduction to Algorithms

**What is an Algorithm?**
An algorithm is a step-by-step procedure or set of rules designed to solve a specific problem or perform a task. Think of it like a recipe for cooking - it tells you exactly what steps to follow to achieve a desired result.

**Key Properties of Algorithms:**
• **Input:** Algorithms take zero or more inputs
• **Output:** Algorithms produce at least one output
• **Definiteness:** Each step must be clearly defined
• **Finiteness:** Algorithm must terminate after a finite number of steps
• **Effectiveness:** Steps must be basic enough to be carried out

## 2. Tools to Analyze Algorithms

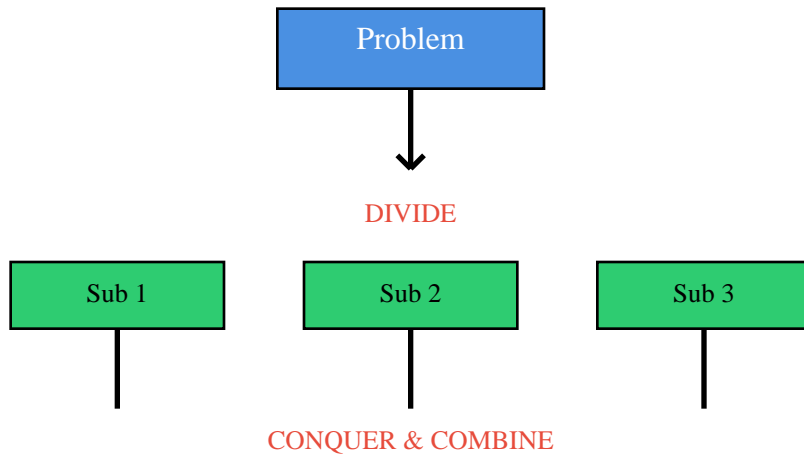**Time Complexity:** Measures how long an algorithm takes to run based on input size.
**Space Complexity:** Measures how much memory an algorithm uses.

**Big O Notation:** Used to describe the upper bound of algorithm performance
• $O(1)$ - Constant time: Same time regardless of input size
• $O(\log n)$ - Logarithmic: Time increases slowly with input
• $O(n)$ - Linear: Time increases proportionally with input
• $O(n^2)$ - Quadratic: Time increases with square of input
• $O(2^\blacksquare)$ - Exponential: Time doubles with each additional input

## 3. Divide and Conquer Technique

**Concept:** Break down a complex problem into smaller, more manageable sub-problems, solve each sub-problem independently, and then combine the solutions.

**Steps:**

1. **Divide:** Split the problem into smaller sub-problems
2. **Conquer:** Solve each sub-problem recursively
3. **Combine:** Merge solutions of sub-problems to get the final solution

**Examples:**
• **Merge Sort:** Divides array into halves, sorts each half, then merges them
• **Quick Sort:** Picks a pivot, partitions array, and recursively sorts partitions
• **Binary Search:** Divides search space in half at each step

## 4. Randomization

**What is Randomization?**
Using random choices in algorithms to improve performance or simplify design. Instead of following a fixed pattern, the algorithm makes random decisions.

**Benefits:**
• Often simpler to implement
• Can avoid worst-case scenarios
• Works well for many practical problems

**Example:** Randomized Quick Sort - picks a random pivot instead of fixed position, which helps avoid worst-case performance on already sorted data.

## 5. Applications

• **Search Engines:** Using efficient algorithms to search billions of web pages
• **GPS Navigation:** Finding shortest path between locations
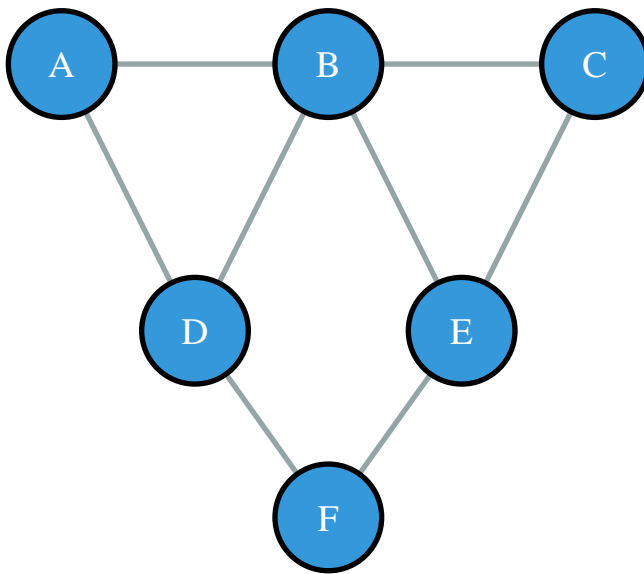• **Data Compression:** Reducing file sizes (ZIP, MP3)

- **Cryptography:** Securing online communications
- **Social Networks:** Recommending friends and content

# UNIT II: Algorithms - Graphs and Data Structures

## 1. Graphs and Maps

**What is a Graph?**
A graph is a collection of nodes (vertices) connected by edges. Think of it like a map where cities are nodes and roads are edges.



**Types of Graphs:**
• **Directed Graph:** Edges have direction (like one-way streets)
• **Undirected Graph:** Edges have no direction (like two-way roads)
• **Weighted Graph:** Edges have values (like distances or costs)
• **Unweighted Graph:** All edges are equal

## 2. Map Searching

**Graph Traversal Algorithms:**

**Breadth-First Search (BFS):**
• Explores all neighbors at current depth before moving deeper
• Uses a queue data structure
• Good for finding shortest path in unweighted graphs
• Example: Finding all friends within 2 connections on social media

**Depth-First Search (DFS):**
- Explores as far as possible along each branch before backtracking
- Uses a stack (or recursion)
- Good for maze solving and detecting cycles
- Example: Exploring all possible moves in a game

**Dijkstra's Algorithm:**
- Finds shortest path between nodes in weighted graph
- Always picks the closest unvisited node
- Used in GPS navigation and network routing

# 3. Stable Marriage Problem

**Problem:** Given n men and n women, where each person has ranked all members of the opposite sex in order of preference, find a stable matching where no two people would prefer each other over their current partners.

**Gale-Shapley Algorithm:**
1. Each unmatched man proposes to his highest-ranked woman who hasn't rejected him
2. Each woman tentatively accepts the best proposal and rejects others
3. Repeat until all are matched
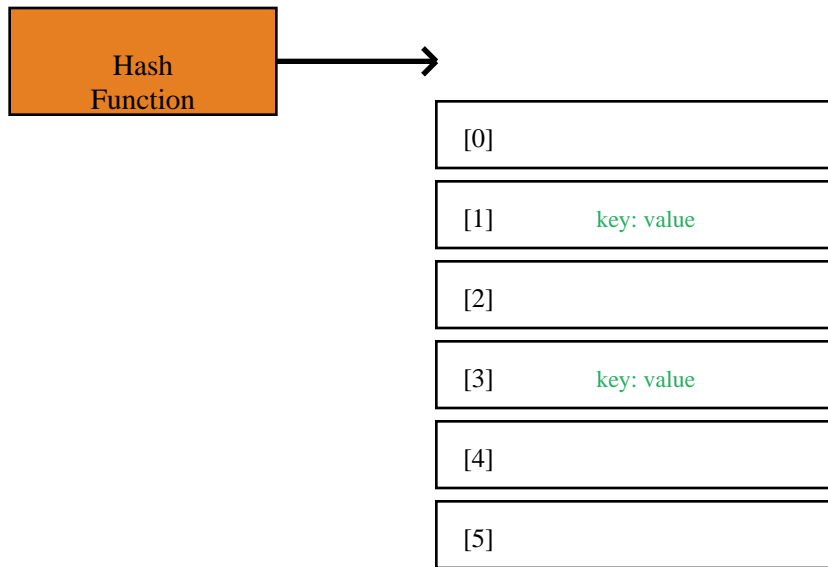
**Applications:**
- Medical residency matching
- College admissions
- Job recruitment

# 4. Dictionaries and Hashing

**Dictionary (Hash Table):**
A data structure that stores key-value pairs and allows fast lookup, insertion, and deletion.

Hash Table Structure

| | |
|---|---|
| Hash Function | → |

| | |
|---|---|
| [0] | |
| [1] | key: value |
| [2] | |
| [3] | key: value |
| [4] | |
| [5] | |

**Hash Function:** Converts keys into array indices
**Collision Handling:**
• **Chaining:** Store multiple items at same index using linked list
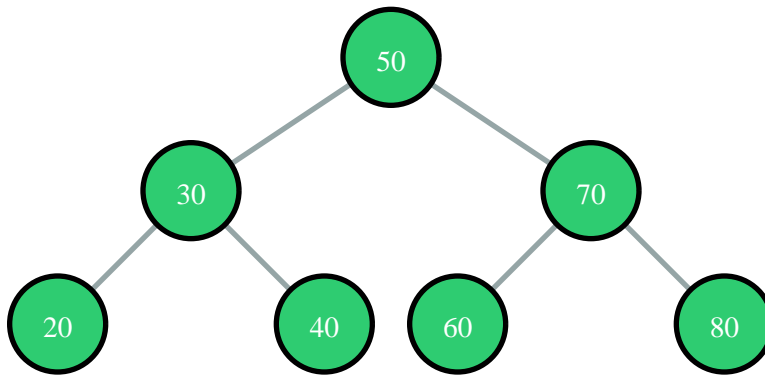• **Open Addressing:** Find another empty slot

**Time Complexity:**
• Average case: O(1) for search, insert, delete
• Worst case: O(n) when many collisions occur

**Applications:** Database indexing, caching, password verification

## 5. Search Trees

**Binary Search Tree (BST):**
A tree where each node has at most two children, and for every node: left child < parent < right child

**Operations:**
- Search: O(log n) average, O(n) worst
- Insert: O(log n) average, O(n) worst
- Delete: O(log n) average, O(n) worst

**Balanced Trees:** AVL and Red-Black trees maintain balance to ensure O(log n) operations

## 6. Dynamic Programming

**Concept:** Solving complex problems by breaking them into simpler sub-problems and storing solutions to avoid recomputing.

**Key Principles:**
1. **Optimal Substructure:** Optimal solution contains optimal solutions to sub-problems
2. **Overlapping Sub-problems:** Same sub-problems are solved multiple times

**Approaches:**
- **Top-Down (Memoization):** Start with main problem, recursively solve and store results
- **Bottom-Up (Tabulation):** Start with smallest sub-problems, build up to main problem

**Examples:**
- Fibonacci numbers
- Longest common subsequence
- Knapsack problem
- Shortest path problems

# UNIT III: Application to Personal Genomics

## 1. Linear Programming

**Definition:** A method to achieve the best outcome (maximum profit or minimum cost) in a mathematical model with linear relationships.

**Components:**
• **Objective Function:** What to maximize or minimize
• **Constraints:** Limitations or requirements
• **Variables:** Decision variables to determine

**Example:** A factory wants to maximize profit by producing products with limited resources.

**Applications:**
• Resource allocation
• Production planning
• Transportation optimization
• Portfolio optimization in finance

## 2. NP-Completeness

**P Problems:** Can be solved in polynomial time (efficiently)
**NP Problems:** Solutions can be verified in polynomial time
**NP-Complete:** Hardest problems in NP - if one is solved efficiently, all can be
**NP-Hard:** At least as hard as NP-Complete problems

**Famous NP-Complete Problems:**
• **Traveling Salesman:** Finding shortest route visiting all cities
• **Knapsack:** Selecting items to maximize value within weight limit
• **Graph Coloring:** Coloring graph nodes so no adjacent nodes share color
• **Boolean Satisfiability (SAT):** Finding variable assignments that satisfy formula

**Practical Approach:** Use approximation algorithms or heuristics for near-optimal solutions

## 3. Introduction to Personal Genomics

**What is Personal Genomics?**
The study of an individual's complete genetic information (genome) to understand health risks, ancestry, and personalized medicine.

**Key Concepts:**
• **DNA:** Contains genetic instructions (A, T, G, C bases)

- **Genome:** Complete set of DNA (~3 billion base pairs in humans)
- **Gene:** Segment of DNA that codes for a protein
- **SNP (Single Nucleotide Polymorphism):** Variation at single position in DNA

**Applications:**
- Disease risk prediction
- Personalized drug selection
- Ancestry tracing
- Pharmacogenomics (how genes affect drug response)

## 4. Massive Raw Data in Genomics

**Challenges:**
- Human genome = 3 billion base pairs
- Sequencing generates terabytes of data
- Need efficient storage and processing
- High computational requirements

**Data Processing Steps:**
1. **Sequencing:** Reading DNA sequences
2. **Assembly:** Piecing together short reads
3. **Alignment:** Comparing to reference genome
4. **Variant Calling:** Identifying differences
5. **Annotation:** Adding biological meaning

**Technologies:** Next-generation sequencing (NGS), high-performance computing, cloud storage

## 5. Data Science on Personal Genomes

**Techniques Used:**
- **Pattern Recognition:** Finding disease-associated genetic patterns
- **Statistical Analysis:** Determining significance of variants
- **Machine Learning:** Predicting disease risk from genetic data
- **Data Visualization:** Presenting complex genomic information

**Key Tasks:**
- Identifying disease-causing mutations
- Predicting drug responses
- Understanding genetic ancestry
- Discovering gene-environment interactions

## 6. Interconnectedness of Personal Genomes

**Concept:** Individual genomes don't exist in isolation - they're connected through:

• **Family Relationships:** Shared genetic variants with relatives
• **Population Structure:** Genetic similarities within ethnic groups
• **Evolution:** Common ancestors and historical migrations
• **Disease Networks:** Multiple genes contributing to same disease

**Privacy Implications:**
• Your genetic data reveals information about relatives
• Database matches can identify individuals
• Need for ethical guidelines and data protection

# 7. Case Studies

### Case 1: Predicting Diabetes Risk
Using genetic markers combined with lifestyle factors to predict Type 2 diabetes risk, enabling early intervention and prevention.

### Case 2: Pharmacogenomics in Cancer Treatment
Analyzing tumor genome to select most effective chemotherapy drugs and avoid treatments likely to cause severe side effects.

### Case 3: Ancestry and Migration
Using genomic data to trace human migration patterns and understand population history.

# UNIT IV: Machine Learning Fundamentals

## 1. Introduction to Machine Learning

**What is Machine Learning?**
A branch of AI where computers learn patterns from data without being explicitly programmed. Instead of following fixed rules, they improve through experience.

**Types of Machine Learning:**

**1. Supervised Learning:**
• Learn from labeled data (input-output pairs)
• Examples: Email spam detection, house price prediction
• Algorithms: Linear regression, decision trees, neural networks

**2. Unsupervised Learning:**
• Find patterns in unlabeled data
• Examples: Customer segmentation, anomaly detection
• Algorithms: K-means clustering, PCA

**3. Reinforcement Learning:**
• Learn by trial and error with rewards
• Examples: Game playing, robotics
• Algorithms: Q-learning, Deep Q-Networks

| Data Collection | → | Data Preprocessing | → | Model Training | → | Model Evaluation | → | Prediction |
|---|---|---|---|---|---|---|---|---|

## 2. Classification

**Definition:** Predicting which category or class an input belongs to.
Example: Is this email spam or not spam? (Binary classification)
Example: What digit is in this image? (Multi-class classification)

**Common Classification Algorithms:**

**1. Logistic Regression:**
• Despite name, used for classification
• Outputs probability of belonging to a class
• Simple and interpretable

**2. Decision Trees:**
• Makes decisions based on asking series of questions
• Easy to understand and visualize
• Can handle both numerical and categorical data

**3. K-Nearest Neighbors (KNN):**
• Classifies based on majority vote of k nearest neighbors
• Simple but can be slow for large datasets

**4. Support Vector Machines (SVM):**
• Finds optimal boundary between classes
• Works well for high-dimensional data

**5. Neural Networks:**
• Inspired by brain structure
• Can learn complex patterns
• Requires more data and computation

# 3. Linear Classification

**Concept:** Separating classes using a straight line (or hyperplane in higher dimensions).

**Perceptron Algorithm:**
• Simplest linear classifier
• Finds a line that separates two classes
• Updates weights based on misclassified points

**Linear Discriminant Analysis (LDA):**
• Finds linear combination of features that best separates classes
• Assumes normal distribution of data
• Also used for dimensionality reduction

**Limitations:**
• Only works when data is linearly separable
• Can't handle complex, non-linear patterns
• Solution: Use kernel methods or non-linear classifiers

# 4. Ensemble Classifiers

**Idea:** Combine multiple models to get better predictions than any single model. Like asking multiple experts and combining their opinions.

**Types of Ensemble Methods:**

**1. Bagging (Bootstrap Aggregating):**
• Train multiple models on different random subsets of data
• Average predictions (regression) or vote (classification)
• **Random Forest:** Popular bagging method using decision trees
• Reduces variance and prevents overfitting

**2. Boosting:**
• Train models sequentially, each correcting previous errors
• **AdaBoost:** Adjusts weights of misclassified examples
• **Gradient Boosting:** Builds trees to minimize loss function
• **XGBoost:** Efficient implementation, widely used in competitions

**3. Stacking:**
• Train multiple different models
• Use another model to combine their predictions
• More complex but can achieve best performance

**Benefits:**
• Improved accuracy
• Reduced overfitting
• More robust predictions

# 5. Model Selection

**Challenge:** Choosing the best model and parameters for your problem.

**Considerations:**
• **Problem Type:** Classification, regression, clustering?
• **Data Size:** Some models need more data than others
• **Feature Types:** Numerical, categorical, text, images?
• **Interpretability:** Do you need to explain predictions?
• **Speed:** Training time vs prediction time
• **Accuracy:** How precise do predictions need to be?

**Bias-Variance Tradeoff:**
• **High Bias:** Model too simple, underfits data

- **High Variance:** Model too complex, overfits training data
- **Goal:** Find balance for best generalization

**Hyperparameter Tuning:**
- **Grid Search:** Try all combinations of parameters
- **Random Search:** Try random combinations
- **Bayesian Optimization:** Smart search based on previous results

# 6. Cross-Validation

**Purpose:** Reliably estimate how well model will perform on new, unseen data.

**K-Fold Cross-Validation:**
1. Split data into k equal parts (folds)
2. Train on k-1 folds, test on remaining fold
3. Repeat k times, each fold used as test once
4. Average the k performance scores

**Common Choices:**
- k=5 or k=10 most common
- k=n (Leave-One-Out) for small datasets

**Stratified K-Fold:**
- Maintains class distribution in each fold
- Important for imbalanced datasets

**Benefits:**
- Uses all data for both training and testing
- Reduces variance in performance estimates
- Helps detect overfitting

# 7. Holdout Method

**Approach:** Split data into separate sets for different purposes.

**Typical Split:**
- **Training Set (60-80%):** Used to train the model
- **Validation Set (10-20%):** Used to tune hyperparameters
- **Test Set (10-20%):** Used for final evaluation

**Important Rules:**
- Never train on test data

• Test set touched only once at the end
• Validation set used for model selection

**When to Use:**
• Large datasets where cross-validation is too slow
• When you need a completely independent test set

**Limitations:**
• Smaller effective training set
• Performance depends on split
• Can be unreliable for small datasets

# UNIT V: Machine Learning Applications

## 1. Probabilistic Modeling

**Concept:** Using probability theory to model uncertainty in data and predictions. Instead of giving single answer, gives probability distribution over possible answers.

**Key Probability Concepts:**
• **Prior Probability:** What we believe before seeing data
• **Likelihood:** Probability of data given model
• **Posterior Probability:** Updated beliefs after seeing data
• **Bayes' Theorem:** $P(A|B) = P(B|A) \times P(A) / P(B)$

**Naive Bayes Classifier:**
• Assumes features are independent (naive assumption)
• Despite simplification, works surprisingly well
• Fast and efficient
• Common in text classification (spam detection)

**Bayesian Networks:**
• Graph representing probabilistic relationships
• Nodes = variables, edges = dependencies
• Can reason about uncertain situations
• Used in medical diagnosis, risk assessment

## 2. Topic Modeling

**Purpose:** Automatically discover abstract topics in a collection of documents.

**Latent Dirichlet Allocation (LDA):**
• Most popular topic modeling algorithm
• Assumes each document is mixture of topics
• Each topic is mixture of words

**How it Works:**
1. Decide number of topics (k)
2. Randomly assign words to topics
3. Iteratively improve assignments
4. Words frequently together form topics

**Example:**
Topic 1: car, engine, fuel, drive → Transportation

Topic 2: computer, software, code, program → Technology
Topic 3: health, patient, disease, treatment → Medicine

**Applications:**
• Document organization
• Recommendation systems
• Understanding large text collections
• Content discovery

# 3. Probabilistic Inference

**Definition:** Computing probability distributions over variables of interest given observed data.

**Types of Inference:**

**1. Exact Inference:**
• Calculate exact probabilities
• Variable elimination
• Works for small networks

**2. Approximate Inference:**
• Estimate probabilities when exact is too slow
• **Sampling Methods:**
- Monte Carlo sampling
- Markov Chain Monte Carlo (MCMC)
- Gibbs sampling
• **Variational Methods:**
- Approximate complex distributions with simpler ones

**Hidden Markov Models (HMM):**
• Model sequences with hidden states
• Applications: Speech recognition, gene finding
• Inference: Find most likely sequence of states

# 4. Application: Prediction of Preterm Birth

**Problem:** Predicting which pregnancies are at risk of preterm birth (before 37 weeks) to enable early intervention.

**Data Sources:**
• Medical history (previous pregnancies, conditions)
• Clinical measurements (blood pressure, weight gain)

- Laboratory tests (hormone levels, biomarkers)
- Ultrasound measurements
- Demographics (age, ethnicity, socioeconomic factors)

**Machine Learning Approach:**
1. **Data Collection:** Gather comprehensive medical records
2. **Feature Engineering:** Select and create relevant predictors
3. **Model Selection:** Compare algorithms (logistic regression, random forest, neural networks)
4. **Training:** Learn patterns from historical data
5. **Validation:** Test on separate patient group
6. **Clinical Integration:** Deploy in healthcare setting

**Challenges:**
- Imbalanced data (preterm births are minority)
- Missing data in medical records
- Need for interpretability for clinical use
- Ethical considerations in healthcare AI

**Impact:**
- Early identification of high-risk pregnancies
- Targeted interventions
- Reduced healthcare costs
- Improved outcomes for mothers and babies

# 5. Data Description and Preparation

**Importance:** Data quality determines model quality. "Garbage in, garbage out."
Typically 60-80% of ML project time spent on data preparation!

**Data Description:**
- **Exploratory Data Analysis (EDA):**
- Understand data distribution
- Identify patterns and anomalies
- Visualize relationships
- **Statistical Summaries:**
- Mean, median, mode
- Standard deviation, quartiles
- Correlation between features

**Data Cleaning:**
- **Handle Missing Values:**
- Remove rows/columns

- Impute with mean/median/mode
- Use predictive models to fill
• **Remove Duplicates:** Identify and eliminate redundant records
• **Fix Errors:** Correct typos, inconsistencies, outliers

**Data Transformation:**
• **Scaling:** Normalize features to similar ranges
- Min-Max scaling: Scale to [0,1]
- Standardization: Zero mean, unit variance
• **Encoding:** Convert categorical to numerical
- Label encoding: Assign numbers to categories
- One-hot encoding: Binary columns for each category
• **Feature Engineering:** Create new meaningful features
- Combinations of existing features
- Domain-specific transformations

**Feature Selection:**
• Remove irrelevant or redundant features
• Methods: correlation analysis, feature importance, PCA
• Benefits: Faster training, better generalization, easier interpretation

# 6. Relationship Between Machine Learning and Statistics

**Machine Learning vs Statistics:**
Both fields overlap significantly but have different emphasis and terminology.

**Similarities:**
• Both work with data to make inferences
• Use similar mathematical foundations
• Many algorithms originated in statistics
• Both concerned with prediction and understanding

**Key Differences:**

| Aspect | Statistics | Machine Learning |
|---|---|---|
| Primary Goal | Inference & interpretation | Prediction & automation |
| Focus | Understanding relationships | Performance on new data |
| Model Complexity | Often simpler, interpretable | Can be very complex |
| Assumptions | Explicit assumptions about data | Often fewer assumptions |
| Sample Size | Works with smaller samples | Benefits from large data |

| Validation | Hypothesis testing, p-values | Cross-validation, test sets |
|---|---|---|
| Example | Linear regression analysis | Deep neural networks |

**Modern Convergence:**

• Fields increasingly borrowing from each other

• Statistical methods used for ML model interpretation

• ML techniques handling complex statistical models

• "Statistical learning" bridges both fields

**When to Use Which:**

• **Use Statistics when:** Need to understand causation, test hypotheses, explain relationships

• **Use ML when:** Focus on prediction accuracy, have large data, don't need full interpretation

• **Use Both:** Most modern data science projects benefit from both approaches

# Summary and Key Takeaways

**Unit I - Algorithms Fundamentals:**
Algorithms are step-by-step procedures for solving problems. We analyze them using time and space complexity. Key techniques include divide-and-conquer (breaking problems into smaller parts) and randomization (using random choices to improve performance).

**Unit II - Data Structures & Advanced Algorithms:**
Graphs represent connections between objects and can be searched efficiently using BFS and DFS. Hash tables provide fast lookups, search trees maintain sorted data, and dynamic programming solves complex problems by storing sub-problem solutions.

**Unit III - Genomics Applications:**
Personal genomics uses algorithms to analyze massive DNA data. Linear programming optimizes resources, while NP-completeness identifies hard problems. Genomic data reveals health risks and requires careful handling due to privacy and interconnectedness with relatives.

**Unit IV - Machine Learning Basics:**
Machine learning lets computers learn from data without explicit programming. Classification predicts categories using various algorithms. Ensemble methods combine multiple models for better accuracy. Cross-validation and holdout methods ensure models generalize well to new data.

**Unit V - Advanced ML Applications:**
Probabilistic modeling handles uncertainty using probability theory. Topic modeling discovers themes in documents. Real applications like preterm birth prediction show ML's healthcare potential. Data preparation is crucial and time-consuming. ML and statistics overlap but emphasize different goals.

## Study Tips

• Understand core concepts rather than memorizing
• Practice implementing algorithms in code
• Draw diagrams to visualize data structures and processes
• Work through examples step-by-step
• Connect theoretical concepts to real-world applications
• Review regularly and test yourself
• Form study groups to discuss complex topics