⚆ OPEN ACCESS

# A Comprehensive Mathematical Exposition of Machine Learning Algorithms and Applications

Samuel Okon Essang[1*], Adie Benimpuye Emmanuel[2], Sylvia Adaobi Akpotuzor[1], Philia Agam Ayuk[3], Aigberemhon Ede Moses[4], Blessed Yahweh[5], Nko Samuel Bassey[6], Ekemini Anietie Johnson[7], Jude Alphonsus Inyangetoh[8], Anietie Emmanuel John[9], Jackson Efiong Ante[10]

[1]Department of Mathematics and Computer Science, Arthur Jarvis University, Akpabuyo
[2]Department of Mathematics and Computer Science Education, University of Calabar, Calabar
[3]Department of Statistics, University of Calabar, Calabar
[4]Electrical Electronics Department, University of Cross River State, Calabar
[5]Department of Research and Technological Development, The MindBook Group, Uyo,
[6]Department of Biological Sciences, Topfaith University, Mkpatak
[7]Department of Computer Science, Federal Polytechnic, Ukana
[8]Department of Statistics, Federal Polytechnic, Ukana
[9]Department of Computer Science, Ritman University, Ikot Ekpene
[10]Department of Mathematics, Topfaith University, Mkpatak

**\*Corresponding author:** Samuel Okon Essang
Department of Mathematics and Computer Science, Arthur Jarvis University, Akpabuyo

| Abstract | | Original Research Article |
|---|---|---|

A deep dive into the mathematical foundations of artificial intelligence that is both brief and comprehensive. The purpose of this article is to shed light on the fundamental rigor that lies behind every intelligent system by illuminating how the beautiful concepts of linear algebra, calculus, probability, and optimization drive machine learning's capacity to learn, adapt, and tackle the world's most challenging tasks.
**Keywords:** Machine Learning, Mathematical Principles, Optimization Algorithms, Data Transformation, Predictive Modeling.

## 1. INTRODUCTION

### Overview of Machine Learning (The Mathematical Perspetive)

Imagine a world where computers don't just follow instructions but actually learn from experience, adapt, and make smart decisions on their own. This isn't science fiction; it's the reality of machine learning (ML). At its heart, ML teaches algorithms to uncover hidden patterns and make predictions or choices directly from data, much like humans learn, but at an incredible scale and speed.

**This vast field breaks down into a few main ways machines learn:**

1. Supervised learning is like a diligent student with a teacher. It learns from examples where both the input and the correct answer are provided. It then figures out how to make predictions for new, unseen information. Think of it identifying spam emails or predicting house prices [3].

2. Unsupervised learning is the curious explorer. It dives into data without any labels or guidance, searching for hidden structures, relationships, or ways to simplify complex information [1]. This includes tasks like grouping similar data points together (clustering) or making complex data easier to understand (dimensionality reduction).

3. Reinforcement learning (RL) is a dynamic approach where an "agent" learns by interacting with its environment, much like a child learning to ride a bike through trial and error. It receives rewards or penalties for its actions, constantly refining its strategy to maximize its overall "score" over time [4].

The incredible power of ML isn't magic; it's built on deep mathematical foundations. Math isn't just a side tool for ML; it's the indispensable core for developing, analysing, and optimizing these systems. This mathematical bedrock provides the rigorous framework to precisely define algorithms, understand how they work, prove they will actually learn, quantify any uncertainties, and fine-tune their performance [2]. Without a solid grasp of these mathematical principles, it's impossible to truly understand, build, troubleshoot, or improve ML models for real-world impact.

### *Mathematical Foundations of Machine Learning*
### *Linear Algebra*

Linear algebra and calculus are foundational in machine learning. Linear algebra helps in organizing and manipulating data using structures like vectors, matrices, and tensors. Operations such as addition, multiplication, and transformations (e.g., scaling and rotation) are key for data manipulation and model computations. Eigenvalues and Singular Value Decomposition (SVD) help with dimensionality reduction, crucial for efficient data processing.

### *Calculus and Optimization Theory*

Calculus, especially derivatives and gradients, is essential for optimizing machine learning models. The gradient points to the steepest increase of a function, and algorithms like Gradient Descent use this to minimize errors by adjusting model parameters iteratively. Calculus also ensures models can learn effectively, especially with differentiable activation functions in neural networks, allowing backpropagation and optimization [7, 8].

### *Chain Rule and Backpropagation*:

Derivation and Significance in Neural Networks. The Chain Rule is a cornerstone of calculus, a fundamental rule for computing the derivative of a composite function. If a function f depends on g, which in turn depends on x, then the chain rule states $\frac{df}{dx} = \frac{df}{dg}\frac{dg}{dx}$. This elegant rule is absolutely central to Backpropagation, the key algorithm for efficiently training neural networks. Backpropagation efficiently computes the gradients of the loss function with respect to all network weights and biases. It involves a "forward pass" where input data propagates through the network to compute the final prediction, followed by a "backward pass" where gradients are computed by chaining derivatives layer-by-layer from the output back to the input, precisely using the Chain Rule [10, 12].

In machine learning, the Jacobian and Hessian matrices are used for more advanced optimization. The Jacobian tracks the first-order changes in multi-output functions, while the Hessian provides second-order derivatives, helping optimize functions faster using methods like Newton's method.
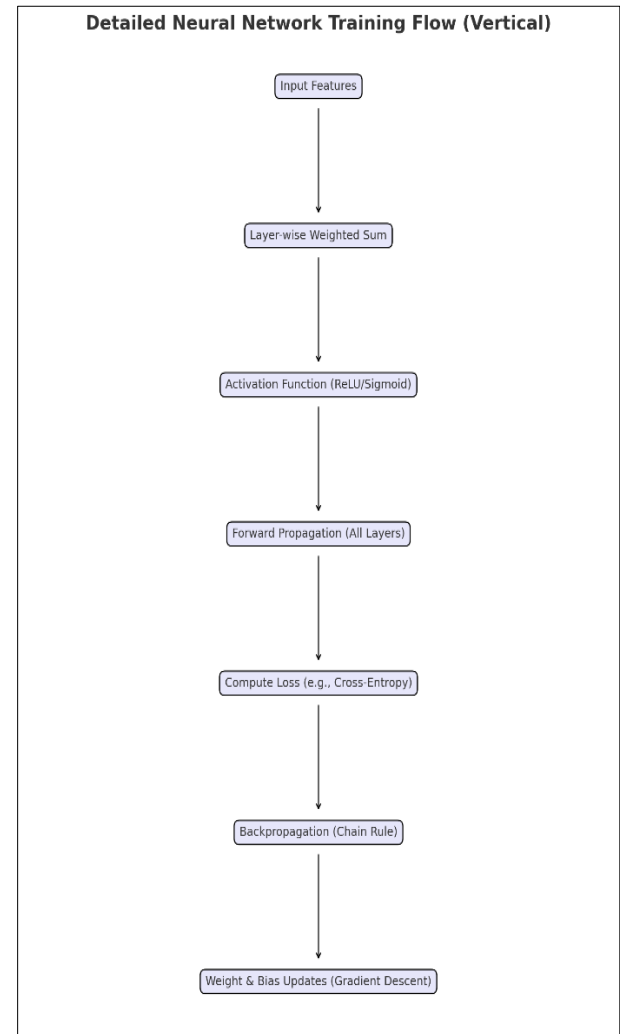


**Figure 2: Neural Network Training Flow Diagram, showing a top-down progression from input features through each training stage to parameter updates**

Probability and statistics are vital for modelling uncertainty in data. Random variables and distributions (Bernoulli, Binomial, Gaussian) help quantify uncertainty. Bayes' theorem is essential for updating beliefs, with applications like Naive Bayes classification. Techniques like Maximum Likelihood Estimation (MLE) and Maximum A Posteriori (MAP) estimate model parameters and reduce overfitting. Hypothesis testing and confidence intervals assess model reliability, and descriptive statistics summarize data characteristics [8,9].

Optimization theory focuses on minimizing loss functions to improve machine learning models, making it essential for model performance and application.

***Supervised Learning***: The primary goal is to find an optimal mapping function $f(x_i, \theta)$ with parameters θ

that minimizes a loss function $L(y_i, f(x_i, \theta))$ over the training samples. Common loss functions include the squared Euclidean distance (Mean Squared Error), cross-entropy, and hinge loss. To mitigate overfitting, regularization terms (e.g., L2 norm) are frequently added to the objective function, resulting in a penalized loss: $min_\theta \frac{1}{N} \sum_{i=1}^{N} L(\mathbf{y}_i, f(\mathbf{x}_i, \boldsymbol{\theta})) + \lambda \|\boldsymbol{\theta}\|_2^2$ *Unsupervised*

*Learning*: Clustering (K-Means): Aims to partition samples into K clusters by minimizing the within-cluster sum of squares (WCSS), which measures the squared Euclidean distance between each point and its assigned cluster centroid: $min_s \sum_{k=1}^{K} \sum_{x \in S_k} \|\mathbf{x} - \boldsymbol{\mu}_k\|_2^2$ [2, 20]

**Dimensionality Reduction (PCA):** Seeks to retain as much original information as possible after projecting data into a low-dimensional space. This is formulated as minimizing the reconstruction error, which is equivalent to maximizing the variance of the projected data.

**Probabilistic Models:** For tasks like density estimation, the objective is to find an optimal probability density function p(x) that maximizes the logarithmic likelihood function (MLE) of the training samples: $\sum_{i=1}^{N} lnp(x_i; \theta)$

Reinforcement Learning: The objective is to find an optimal strategy function (policy π) whose output varies with the environment, with the goal of maximizing cumulative rewards. This is often expressed as maximizing the expected value function of a state s under policy π: $\max_\pi V_\pi(s) = E$ [21].

Optimization algorithms in machine learning are classified into three types: first-order, second-order, and heuristic methods.

1. **First-Order Methods:** These methods, like stochastic gradient descent, use gradient information to iteratively adjust model parameters. They are computationally efficient and scalable but may get stuck in local minima, especially in non-convex functions like those found in deep learning.

2. **Second-Order Methods:** An example is Newton's method, which uses second-order derivatives (Hessian matrix) to speed up convergence. These methods converge faster but are computationally expensive, especially in high-dimensional problems.

3. *Heuristic Methods*: These are derivative-free methods like coordinate descent, which are useful for problems where gradient information is not available or difficult to compute.



**Figure 1: Flow Diagram of Supervised Learning – This illustrates the sequential stages in a typical supervised learning pipeline. Comparison of Optimization Algorithms – This bar chart compares different optimization approaches used in machine learning based on relative speed and convergence quality**

The challenge of local minima arises in non-convex functions, where algorithms may converge to a local optimum rather than the global one. While first-order methods are often used in deep learning due to their scalability and lower computational cost, they may not always find the global minimum. In contrast, second-order methods offer better convergence but at the cost of higher computation [10-14].

Optimization algorithms are essential for machine learning, as they iteratively adjust model parameters to minimize the error, allowing models to learn and adapt from data.

**Table 1: Machine Learning Paradigm and Mathematical Objectives**

| ML Paradigm | General Objective | Example Algorithm | Mathematical Formula (Example) |
|---|---|---|---|
| Supervised Learning | Parameter Learning/Mapping | Linear Regression | $min_\theta \frac{1}{N}\sum_{i=1}^{N} L(\mathbf{y}_i, f(\mathbf{x}_i;\boldsymbol{\theta}))$ |
| Unsupervised Learning (Clustering) | Data grouping | K- Means | $min_s \sum_{k=1}^{K}\sum_{x\in S_k}\|\mathbf{x}-\boldsymbol{\mu}_k\|_2^2$ |
| Unsupervised Learning (Dimensionality Reduction) | Feature Transformation | PCA | $min\sum_{k=1}^{K}\|\mathbf{x_i}-\mathbf{x}_i^{'}\|_2^2$ |
| Unsupervised Learning (Probabilistic Models) | Density Estimation | Bayesian Network | $max_\theta \sum_{i=1}^{N}\log p(\mathbf{x}_i;\theta)$ |
| Reinforcement Learning | Policy Optimization | Q-Learning | $\pi: \max_\pi V_\pi(s) = E\left[\sum_{t=o}^{\infty}\gamma^t r_t\right]$ |

## 2. MATERIALS AND METHODS

This study looks at the math behind machine learning algorithms and how they are used in many industries, such as healthcare and finance. It brings together information from other sources to give a full picture of how arithmetic concepts are used to create and test machine learning models.

The methodology section talks about the main math tools used in machine learning, like linear algebra, calculus, and statistics and probability. It looks at how these fields function with machine learning algorithms such as support vector machines, decision trees, and neural networks. The study also talks about the trade-offs that come with model performance, like the bias-variance trade-off, and how different assessment measures help you choose the best model.

The paper's materials part talks about important uses in healthcare (predictive diagnostics, personalised medication) and finance (algorithmic trading, fraud detection). It shows how machine learning algorithms use strict maths to handle hard problems in the real world. It also stresses the need for openness and moral considerations, especially in fields with high stakes like healthcare.

The goal of this review is to give a methodical look at the mathematical foundations of machine learning while also showing how it may be used in the real world, the problems it faces, and its future potential.

### Core Machine Learning Algorithms: Mathematical Formulations and Applications

This section provides a detailed mathematical exposition of key machine learning algorithms, unveiling their objective functions, optimization strategies, and the elegant mathematical principles that govern their operation.

### Supervised Learning Algorithms

Supervised learning algorithms are the workhorses of predictive modelling, meticulously designed to infer a mapping function from input features (x) to output labels (y) using a dataset where both inputs and their corresponding correct outputs are provided. The learning process typically involves minimizing a predefined loss function that precisely quantifies prediction errors, guiding the model towards accuracy [20,21].

### Linear Regression

Linear regression stands as a foundational supervised learning algorithm, a cornerstone for modeling the straightforward, linear relationship between input features and a continuous output variable. It's often the first step in understanding predictive modeling.

*Mathematical Model and Objective Function (Mean Squared Error).*
The model posits a simple linear relationship, expressed as $\hat{y} = w^T x + b$ where $\hat{y}$ is the predicted output, x is the input feature vector, w represents the weight vector (determining the slope of the relationship), and b is the bias (or intercept). The most common objective function for linear regression is the Mean Squared Error (MSE), which quantifies the average of the squared differences between the predicted values and the actual observed values. Formally, the MSE is given by

$$J(w,b) = \frac{1}{N}\sum_{i=1}^{N}(y_i - (w^T x_i + b))^2$$

where $N$ is the number of samples, $y_i$ are the true labels, and $\hat{y} = \mathbf{W}^T\mathbf{X} + b$ are the predictions. Minimizing MSE is statistically equivalent to maximizing the likelihood of the observed data, under the assumption that the errors

are independently and identically normally distributed with zero mean [22,23,24].

**Closed-Form Solution: Normal Equation Derivation.**
For linear regression, a remarkable property exists: the optimal parameters w and b can be determined directly, without any iterative process, by simply setting the gradient of the MSE objective function with respect to the parameters to zero. This elegant derivation leads to the Normal Equation: $w = \left(\mathbf{X}^T \mathbf{X}\right)^{-1} \mathbf{X}^T \mathbf{y}$ where X is the design matrix (augmented with a column of ones for the bias term) and y is the vector of true labels. This method offers a direct, one-shot calculation of the optimal parameters [25,26].

### *Iterative Solution: Gradient Descent Application.*

While a closed-form solution exists, Gradient Descent is widely applicable to linear regression and becomes an absolute necessity for larger datasets or more complex models where direct matrix inversion is computationally intractable or numerically unstable. In this iterative approach, the weights w and bias b are progressively adjusted by taking steps proportional to the negative of the gradient of the MSE. The update rules are: $\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_w J(\mathbf{w}, b)$ and $b \leftarrow b - \alpha \nabla_w J(\mathbf{w}, b)$ where α is the learning rate. The partial derivatives of the MSE with respect to **w** and b are

$$\nabla_w J = -\frac{2}{N} \sum_{i=1}^{N} (y_i - \widehat{y}_i)\mathbf{x}_i \qquad \text{and}$$

$$\nabla_b J = -\frac{2}{N} \sum_{i=1}^{N} (y_i - \widehat{y}_i)\mathbf{x}_i \quad \text{The MSE objective function}$$

for linear regression possesses the crucial mathematical property of convexity. This is a pivotal characteristic because, for convex functions, gradient descent, when applied with an appropriately chosen learning rate, is guaranteed to converge to the global minimum. This stands in stark contrast to more complex, non-convex models where local minima present a significant challenge. This property highlights how the mathematical nature of the objective function directly determines the convergence guarantees of the optimization algorithm [27,28]. [29,30].

### *3.1.2 Logistic Regression*

Logistic regression is a classification algorithm primarily used for binary classification tasks, modelling the probability of a binary outcome. Despite its name, it's a classifier, not a regressor, and it's a workhorse in many high-stakes domains due to its interpretability [31,32,33].

*Mathematical Model: Sigmoid Function and Probabilistic Interpretation.* The model posits that the log-odds of an event occurring are a linear combination of the input features:

$In\left(\dfrac{p}{1-p}\right) = \mathbf{w}^T \mathbf{x} + b$ to elegantly transform this linear combination into a probability ranging precisely between 0 and 1, the logistic (sigmoid) function is applied: $P(y = 1 \mid \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b) = \dfrac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$ This transformation allows the model's output to be directly interpreted as a probability, which is essential for classification tasks [3, 33].

*Objective Function: Cross-Entropy Loss Derivation.*
For logistic regression, the learning objective is to maximize the likelihood of correctly classifying the training samples, which is mathematically equivalent to minimizing the Log Loss or Binary Cross-Entropy Loss. For a single training example $(\mathbf{x}_i, y_i)$, where $y_i$ is the true binary label an $\widehat{y}_i$, is the predicted probability, the binary cross-entropy loss is $L(y_i, \widehat{y}_i) = -[y_i \log(\widehat{y}) + (1 - y_i)\log(1 - \widehat{y}_i)]$ The total objective function is the average cross-entropy over all $N$ samples:

$$J(\mathbf{w}, b) = -\frac{1}{N} \sum_{i=1}^{N} [y_i \log(\widehat{y}) + (1 - y_i)\log(1 - \widehat{y}_i)] \quad \text{Cross-}$$

entropy fundamentally measures the dissimilarity between the true probability distribution of the labels and the probability distribution predicted by the model, guiding the model to align its predictions with reality[34].

Logistic regression uses gradient descent to minimise cross-entropy loss by updating weights and bias iteratively. This update uses cross-entropy loss function gradients for gradient descent. Update parameters like linear regression but use the sigmoid activation function and cross-entropy loss for classification problems. Probabilistic models like logistic regression benefit from cross-entropy loss because it matches maximum likelihood estimation for categorical data, guiding probability predictions.

However, Support Vector Machines (SVMs) classify by finding the best separating hyperplane. Hard Margin SVMs maximise the "margin," the distance between the closest data points from each class and the decision boundary. This maximisation separates classes, making SVMs good for binary classification.

The hyperplane itself is mathematically defined by the equation $0 = \mathbf{w} \bullet \mathbf{x} + b$ The margin is the perpendicular distance between two parallel hyperplanes $1 = \mathbf{w} \bullet \mathbf{x} + b, -1 = \mathbf{w} \bullet \mathbf{x} + b$ which is given by $\dfrac{2}{\|\mathbf{w}\|}$. Therefore, maximizing this margin is mathematically equivalent to minimizing $\|\mathbf{w}\|^2 \ or \ \dfrac{1}{2}\|\mathbf{w}\|^2$

subject to the constraints $y_i(w\cdot xi+b)\geq 1$ for all training samples $(x_i,y_i)$ where $y_i \in \{-1,+1\}$ This formulation seeks the widest possible "street" between the classes [35,36].

*Soft Margin SVM:* Introduction of Slack Variables and Hinge Loss. The strict demands of Hard Margin SVM that data be perfectly linearly separable are rarely met in messy, real-world datasets. To address this limitation, Soft Margin SVM was introduced, gracefully allowing for some misclassifications or data points to fall within the margin. This relaxation is achieved through the ingenious introduction of non-negative slack variables, $\xi_i \geq 0$ which modify the constraint to $y_i(w\cdot xi+b)\geq 1-\xi_i$ A data point is considered misclassified if its corresponding $\xi_i > 0$ The objective function for Soft Margin SVM then becomes a minimization of a combination of the margin term and a penalty for constraint violations: $\min_{w,b,\xi} \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{N}\xi_i$ . Here, C is a regularization parameter that carefully balances the trade-off between maximizing the margin and minimizing the classification errors. The term $\xi_i$ can be expressed as $\max\left(0, 1-y_i(\mathbf{w}\cdot\mathbf{x}i+b)\right)$ , which is famously known as the Hinge Loss. Gradient Descent can be employed to optimize this hinge loss [36].

*Optimization: Lagrangian Duality and Karush-Kuhn-Tucker (KKT) Conditions.* SVM optimization problems are formulated as constrained convex optimization problems, which are particularly well-behaved. These problems are typically solved by transforming the constrained primal problem into an unconstrained dual problem using the powerful method of Lagrange Multipliers. For the hard margin SVM, the Lagrangian is $L(\mathbf{w},b,\boldsymbol{\lambda}) = \frac{1}{2}\|w\|^2 - C\sum_{i=1}^{N}\lambda_i\left(y_i(\mathbf{w}\cdot\mathbf{x}i+b)-1\right)$ The optimal solution must satisfy the Karush-Kuhn-Tucker (KKT) conditions, which are a set of necessary (and sufficient for convex problems) conditions for optimality in constrained optimization. These conditions involve the gradients of the Lagrangian with respect to both primal and dual variables, along with complementary slackness conditions, $\lambda_i\left(y_i(\mathbf{w}\cdot\mathbf{x}i+b)-1\right)=0$ . The KKT conditions are particularly significant as they identify the "support vectors" those crucial data points for which $\lambda_i>0$ and which lie precisely on the margin boundary, thus critically influencing the position of the hyperplane [37].

Other machine learning algorithms for classification and regression include decision trees. They recursively split feature space into subgroups using simple rules. Gini Impurity and Entropy measure how mixed or pure the subsets are, determining each split's effectiveness. Gini Impurity measures subset misclassification, helping the tree make more accurate predictions.Its formula is $GI = 1 - \sum_{j=1}^{c}(p_j)^2$ , where $p_j$ is the proportion of instances belonging to class j in the given set. For binary classification, Gini Impurity ranges from 0 (perfect purity) to 0.5 (maximum impurity); for multi-class classification, its range is $\left[0, 1-\frac{1}{C}\right]$. A significant computational advantage of Gini Impurity is that it does not involve logarithmic functions, making its calculation faster than entropy [38].

Entropy: In information theory, entropy quantifies the amount of uncertainty, randomness, or "information" associated with a randomly chosen variable's potential outcome. Its formula is $H = -\sum_{j=1}^{C}p_j \log_2(p_j)$ , where $p_j$ is the proportion of instances of class $j$ in the set. Entropy ranges from 0 (perfect purity) to $log_2C$ (maximum impurity), where $C$ is the number of classes; for binary classification, the range is [0,1]. Entropy tends to favour splits that result in a higher reduction of uncertainty, seeking to maximize information gain [40].

In decision trees, the use of Gini Impurity and Entropy ties the tree-building process to key principles from information theory, where the goal is to maximize information gain or minimize uncertainty with each split. Each split aims to create purer subsets, reducing the randomness or impurity of the data within the nodes. The choice between Gini Impurity and Entropy often involves a trade-off: Gini Impurity is faster to compute, making it suitable for large datasets, while Entropy, though more theoretically pure, comes at the cost of higher computational expense.

For regression trees, the objective is to minimize Mean Squared Error (MSE), which ensures that the resulting child nodes contain more homogeneous output values, leading to better predictions.

In unsupervised learning, K-Means clustering is a popular algorithm that groups data into clusters based on the proximity to the nearest centroid. The core goal of K-Means is to minimize the Within-Cluster Sum of Squares (WCSS), a measure of how tightly grouped the data points are within each cluster, thus aiming to create more cohesive and distinct groups.

Formally, the objective is to find the partition $S = \{S_1, S_2, ..., S_k\}$ that minimizes $\sum \sum_{x \in S_k} \|x - \mu_k\|_2^2$ where $\mu k$ represents the mean (centroid) of the points within cluster $S_k$. This formulation is equivalent to minimizing the pairwise squared deviations of points within the same cluster [41,42].

*Algorithm Steps*: Mathematical Basis for Assignment and Update Steps. The K-Means algorithm proceeds iteratively through two main steps, a dance between assignment and refinement:

1. *Initialization:* The process begins by randomly generating K initial cluster centroids within the data domain.
2. *Assignment Step (E-Step):* In this phase, each observation is assigned to the cluster whose centroid exhibits the least squared Euclidean distance to that observation. This effectively partitions the observations according to the Voronoi diagram generated by the current set of centroids. Mathematically, for each data point x(i), its cluster assignment c(i) is determined by

$$c^{(i)} = \arg \min_k \left\| \mathbf{x}^{(i)} - \mathbf{\mu}_k \right\|_2^2$$

3. *Update Step (M-Step)***:** Following the assignment, the centroids for each cluster are recalculated. Each new centroid $\mu_k$ is determined as the mean (geometric mean) of all the data points that have been assigned to that specific cluster;

$$\mu_k = \frac{1}{|S_k|} \sum_{x \in S_k} \mathbf{x}$$

If a centroid ends up with no points assigned to it, it is typically re-initialized randomly [33,34,37].

Moreover, the algorithm's reliance on squared Euclidean distance ($L_2$ norm) to measure cluster similarity implies that it tends to favour spherical clusters. If the true data structure requires a different distance metric (e.g., Manhattan distance or cosine similarity), K-Means may not perform optimally, highlighting the importance of choosing an appropriate distance metric for the data at hand.This process enables clearer insights from high-dimensional data while reducing computational complexity.:

$$\max_{\|\mathbf{w}\|=1} \{\|\mathbf{Xw}\|^2\} = \max_{\|\mathbf{w}\|=1} \{\mathbf{w}^T \mathbf{X}^T \mathbf{Xw}\}$$

This formulation seeks to identify the direction along which the data exhibits the greatest spread [38].

*Derivation***:** *Covariance Matrix, Eigenvectors, and Eigenvalues.* The derivation of PCA involves several key mathematical steps:

1. *Centering Data:* Prior to performing PCA, the data is typically centered by subtracting the mean of each variable from all its observations. This ensures that the principal components pass through the centroid of the data, simplifying the subsequent mathematical formulation.

2. *Covariance Matrix:* The principal components are mathematically proven to be the eigenvectors of the data's *covariance matrix*. The sample covariance matrix $\mathbf{C}$ is computed from the mean-subtracted data matrix $\mathbf{B}$ as $C = \frac{1}{N-1} \mathbf{B}^{\mathbf{T}} \mathbf{B}$ (or $\mathbf{X}^{\mathbf{T}} \mathbf{X}$) for mean-centered data). This matrix is inherently symmetric and positive semidefinite.

3. *Eigenvectors and Eigenvalues***:** The problem of maximizing the variance of the projected data is mathematically equivalent to finding the eigenvectors of the covariance matrix. The *eigenvalues* ($\lambda_k$) associated with these eigenvectors ($v_k$) quantify the amount of variance captured along their respective principal components. Larger eigenvalues correspond to more significant principal components, indicating directions of greater data spread. A notable property is that the principal components are uncorrelated with each other.

In dimensionality reduction, Principal Component Analysis (PCA) projects high-dimensional data onto a lower-dimensional subspace by selecting eigenvectors with the largest eigenvalues. This process effectively retains most of the data's variance, which is often used as a proxy for information. While PCA is useful for reducing noise and facilitating visualization, it has a crucial assumption: variance indicates the most "important" information in the data. However, this can result in the loss of information if it's encoded in directions with low variance. Despite this, PCA remains an invaluable tool, particularly because it combines linear algebra (eigenvectors) and statistics (covariance) to reduce dimensionality while preserving key data characteristics.For a single neuron, the net input to the activation function is $z = \sum_i w_i x_i + b = \mathbf{w}T\mathbf{x} + b$

When considering an entire layer, these operations can be efficiently represented using matrix multiplication: $\mathbf{Z} = \mathbf{XW} + \mathbf{B}$,where X is the input matrix, W is the weight matrix for the layer, and B is the bias vector [15, 24, 43].

***Activation Functions***

Activation functions are critical components within neural networks, primarily responsible for introducing non-linearity into the system. This non-linearity is what enables networks to learn and model complex, non-linear mappings between inputs and outputs. These functions dictate the output of a neuron given its aggregated input.

***Mathematical Definitions and Properties (Sigmoid, Tanh, ReLU, Softmax).***

*Linear Activation Function:* Defined as *f(x)=x*, this function simply returns the input as the output. While simple, it restricts the network to modeling only linear relationships, rendering the "depth" of the network irrelevant for learning complexity. It is primarily used in the output layer for regression problems where a numerical value is predicted. Sigmoid *(Logistic) Function*: Mathematically expressed as $\sigma(x) = \frac{1}{1+e^{-x}}$ ,the sigmoid function squashes any real-valued input into the range (0,1). Its output can be interpreted as a probability, making it suitable for binary classification. Historically significant, it suffers from the *vanishing gradient problem* for very large or very small inputs, which can impede learning in deep networks.

*Tanh (Hyperbolic Tangent) Function*: Defined as $f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$, the tanh function outputs values in the range [−1,1]. It is essentially a scaled version of the sigmoid function, but its zero-centered output and generally stronger gradient (steeper derivatives) can lead to faster convergence in certain problems compared to sigmoid.

*ReLU (Rectified Linear Unit):* The ReLU function is defined as $f(x) = \max(0, x)$, It outputs the input directly if it is positive, and zero otherwise. ReLU is computationally efficient and helps mitigate the vanishing gradient problem for positive inputs. However, it can suffer from the "dying ReLU" problem where neurons can become inactive if their input is consistently negative [27].

*Softmax Function:* Primarily used in the output layer for multi-class classification problems, the softmax function transforms a vector of raw scores (logits) into a probability distribution where all elements are non-negative and sum to 1. For an input vector x with elements $x_1,...,x_C$,the softmax for element $x_i$ is $f(x_i) = \frac{e^{x_i}}{\sum_{j=1}^{C} e^{x_i}}$ .The primary role of activation functions in neural networks is to introduce non-linearity, enabling the network to learn and represent complex, non-linear relationships in data. Without non-linearity, even a deep multi-layered network would behave like a single linear model, unable to capture the intricate patterns in real-world data. By introducing non-linear behaviour, activation functions significantly increase the flexibility and power of neural networks, allowing them to model more complex data.

In deep learning, loss functions are crucial for training as they measure the discrepancy between the model's predictions and the actual labels. These functions guide the optimization process, helping the model minimize errors and learn effectively. Common loss functions include:

1.  Mean Squared Error (MSE): Measures the average of the squared differences between predicted and actual values. It is commonly used in regression tasks.
2.  Binary Cross-Entropy: Used for binary classification, it quantifies the difference between predicted probabilities and true binary labels.
3.  Categorical Cross-Entropy: Used for multi-class classification, it measures the difference between predicted probabilities and the true categorical labels.

Each of these loss functions plays a vital role in determining how the model's parameters should be updated to improve its predictions.

Widely used for regression tasks where the output is continuous. It calculates the average of the squared differences between predicted values $(\hat{y})$

And actual values $y_i : J = \frac{1}{N} \sum_{i=1}^{N} (y_i - \hat{y}_i)^2$ , due to the squaring of errors, MSE penalizes larger errors more heavily and is sensitive to outliers [40].

Binary Cross-Entropy Loss (Log Loss): Employed for binary classification problems where the model's output is a probability value between 0 and 1. It measures the performance of a classification model by penalizing false classifications, particularly when the predicted probability deviates significantly from the true label. The formula for N samples is $J = -\frac{1}{N} \sum_{i=1}^{N} \left[ y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right]$, where $y_i$ is the true binary label and $\hat{y}_i$ is the predicted probability.

Categorical Cross-Entropy Loss: Used for multi-class classification problems, typically in conjunction with a Softmax activation function in the output layer. It calculates the error for predicted class probabilities across multiple categories. For N samples and M classes, the formula is $J = -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{M} y_{ij} \left[ \log(p_{ij}) \right]$ where $y_{ij}$ is 1 if sample *i* belongs to class *j* and 0 otherwise, and $p_{ij}$ is the predicted probability for class *j*.

Deep learning uses specialized loss functions tailored to specific tasks, like Huber Loss (less sensitive to outliers), Hinge Loss (for SVMs), Dice Loss (in image segmentation), and Kullback-Leibler Divergence (for measuring distribution differences). The choice of loss function depends on the task (e.g., regression, classification) and the type of output, ensuring the model is optimized for its goal, whether it's minimizing errors or maximizing accuracy.

Backpropagation is the key algorithm for training neural networks, efficiently computing gradients using the chain rule. It works in two phases:

1. Forward Pass: Data flows from the input to the output, with each neuron processing inputs through weighted sums and activation functions.
2. Backward Pass: Gradients are calculated from the output back through the network, adjusting the weights and biases of earlier layers. This helps the model improve by iteratively reducing errors.

For a composite function $f(g(x))$, the chain rules states $\dfrac{\partial f}{\partial x} = \dfrac{\partial f}{\partial g}\dfrac{\partial g}{\partial x}$, In a neural network, this translates to calculating how changes in a weight affect the loss by multiplying the sensitivity of the loss to the neuron's output by the sensitivity of the neuron's output to the weight. These computed gradients indicate the sensitivity of the overall loss to changes in individual parameters, guiding their adjustment [44].

Backpropagation implicitly treats the neural network as a differentiable computational graph, where nodes represent mathematical operations and edges represent the flow of data. The chain rule then enables the efficient flow of gradients backward through this graph. This conceptualization is fundamental to modern deep learning frameworks (e.g., TensorFlow, PyTorch), which automate the complex process of gradient computation, allowing researchers and practitioners to focus on model architecture and data rather than tedious manual gradient derivation [33].

### Gradient Descent Variants

Gradient Descent, as the core optimization algorithm, has several variants that differ in how much data they use to compute the gradient at each step. These differences profoundly impact training speed, stability, and convergence characteristics, making the choice of variant a crucial decision for practitioners.

### Batch, Stochastic, Mini-batch Gradient Descent (Mathematical Update Rules).

Batch Gradient Descent (BGD): In BGD, the gradient of the loss function is computed using the *entire* training dataset before each parameter update. The update rule is $\theta \leftarrow \theta - \eta \nabla J(\theta; \mathbf{X}, \mathbf{y})$ where *X, y* represents the full dataset. While BGD is guaranteed to converge to the global minimum for convex functions and provides stable updates, it can be extremely slow for large datasets due to the computational cost of processing all data points in each iteration [37,46].

*Stochastic Gradient Descent (SGD):* In stark contrast to BGD, SGD computes the gradient and updates parameters using only *one* randomly chosen training example at a time for each update. The update rule is $\theta \leftarrow \theta - \eta \nabla J(\theta; \mathbf{x}_i, \mathbf{y}_i)$ for a single sample *i*. SGD offers significantly faster updates, particularly beneficial for very large datasets or online learning scenarios, and its noisy updates can help escape shallow local minima in non-convex landscapes. However, the high variance in updates can lead to oscillations around the minimum, making convergence less smooth.

*Mini-batch Gradient Descent*: This variant strikes a practical balance between BGD and SGD. It computes the gradient using a small *subset* (a "mini-batch") of the training data, typically comprising 32 or 64 examples. The update rule is $\theta \leftarrow \theta - \eta \nabla J(\theta; \mathbf{x}_{batch}, \mathbf{y}_{batch})$ Mini-batch GD balances speed and stability, providing more stable convergence than SGD while leveraging vectorized operations for computational efficiency. The optimal mini-batch size often requires careful tuning to achieve the best performance [9,42].

### The Importance of the Learning Rate

Think of training a machine learning model as trying to find the lowest point in a hilly landscape. The learning rate (η) is like the size of your steps as you walk down. It's an incredibly important setting (what we call a "hyperparameter") in all variations of gradient descent, which is the mathematical engine that helps models learn.

Mathematically, the learning rate precisely controls how big each step is when the model adjusts its internal parameters based on the "slope" of the error landscape. If your steps are too big, you might overshoot the lowest point entirely or even bounce wildly out of control – this is called divergence. On the other hand, if your steps are too small, it will take an incredibly long time, possibly forever, to reach the bottom – this is called slow convergence. For simpler, "bowl-shaped" problems (what mathematicians call convex functions), we can only guarantee that we'll reach the absolute lowest point if we choose the learning rate just right. This shows that the learning rate isn't just something you tweak; it's a mathematically critical factor that directly impacts how stable and fast the learning process is, often posing a significant challenge in practical deep learning [7, 40].

## 4. Advanced Concepts and Model Evaluation Metrics

This section ventures into crucial concepts that enhance model performance and interpretability, alongside the rigorous mathematical metrics used to objectively evaluate machine learning models. It's about refining models and truly understanding their capabilities.

### Regularization Techniques

Regularization techniques are essential mathematical strategies employed to combat the pervasive problem of overfitting in machine learning

models. They achieve this by adding a penalty term to the loss function, thereby subtly constraining model complexity. This constraint encourages the learning of simpler, more robust models that generalize better to unseen data, preventing them from memorizing noise.

### L1 (Lasso) and L2 (Ridge) Regularization: Mathematical Penalty Terms and their Effect on Model Complexity.

*L2 (Ridge) Regularization*: This technique adds a penalty proportional to the sum of the squared magnitudes of the model's coefficients (the L2 norm) to the original loss function. The objective function becomes

$$J_{Ridge}(\theta) = J(\theta) + \lambda \sum_{j=1}^{m} w_j^2 = J(\theta) + \lambda \|\theta\|_2^2,$$

where J(θ) is the original loss, θ represents the model parameters (weights $w_j$), and λ is the regularization parameter. The effect of L2 regularization is to shrink the coefficients towards zero, but it rarely forces them to be exactly zero. This helps in reducing the model's variance, making it less sensitive to training data fluctuations.

L1 (Lasso) Regularization: In contrast, L1 regularization adds a penalty proportional to the sum of the absolute magnitudes of the coefficients (the L1 norm). The objective function is

$$J_{Lasso}(\mathbf{\theta}) = J(\theta) + \lambda \sum_{j=1}^{m} |w_j|$$ A key characteristic of L1 regularization is its remarkable ability to shrink some coefficients exactly to zero, effectively performing automatic feature selection by eliminating less important predictors from the model [30,44].

λ (Regularization Parameter): The parameter λ is a crucial tuning parameter that controls the strength of the regularization. A value of λ=0 implies no regularization, reverting to the original loss function. Increasing λ increases the penalty, leading to simpler models.

Elastic Net. Elastic Net Regression ingeniously combines both L1 and L2 regularization, incorporating both the absolute and squared measures of the weights into the penalty term. Its objective function is

$$J_{ElasticNet}(\mathbf{\theta}) = J(\mathbf{\theta}) + \lambda \left( \alpha \sum_{j=1}^{m} |w_j| + (1-\alpha) \sum_{j=1}^{m} w_j^2 \right)$$ Here, α is a mixing parameter, ranging from 0 to 1, that balances the contributions of L1 and L2 regularization, offering a flexible approach to model complexity [45,46].

Regularization techniques directly address the bias-variance trade-off by imposing a penalty on model complexity. This mathematical intervention aims to improve the model's generalization performance on unseen data. L1 regularization, in particular, promotes sparsity in the model by driving some coefficients to zero, which simplifies the model and can effectively reduce variance. This demonstrates a mathematically defined strategy for navigating the bias-variance dilemma, leading to models that generalize more effectively.

The impact of regularization can be understood through a compelling geometric interpretation of the penalty terms. The L1 penalty term, $\left( \sum_{j=1}^{m} |w_j| \right)$ corresponds to a diamond-shaped constraint region in the parameter space (for a 2D example, $\left( |w_1| + |w_2| \le s \right)$. The L2 penalty term $\sum_{j=1}^{m} w_j^2$ corresponds to a circular (or spherical) constraint region $\left( w_1^2 + w_2^2 \le s \right)$ The minimization of the loss function, subject to these constraints, often leads to solutions at the "corners" of the $L_1$ diamond, where some coefficients are exactly zero. In contrast, the smooth, rounded nature of the $L_2$ circle typically results in coefficients that are shrunk towards zero but rarely become exactly zero. This geometric difference directly explains why L1 regularization leads to sparse solutions (feature selection), while $L_2$ regularization primarily shrinks coefficients without eliminating them [11,45,46].

Bias-Variance Trade-off: In machine learning, the bias-variance trade-off refers to the challenge of minimizing two sources of error: bias and variance, which both affect a model's ability to generalize beyond its training set.

1. Bias is the error from overly simplistic assumptions made by the model, often leading to underfitting. A high-bias model doesn't capture the underlying patterns well, resulting in poor performance on both training and new data.
2. Variance is the error caused by the model's sensitivity to fluctuations or noise in the training data. High variance leads to overfitting, where the model learns the noise in the training data, making it perform well on the training set but poorly on new, unseen data.
3. Irreducible Error represents noise in the data itself, which no model can reduce.

The goal is to strike a balance between bias and variance, finding a model that is complex enough to capture important patterns but simple enough to avoid overfitting to the noise. This balance is often addressed using techniques like regularization, early stopping, pruning, and ensemble methods.

*Underfitting and Overfitting*:
1. Underfitting occurs when a model is too simple, failing to capture the true patterns in the data. It typically shows high bias and low variance.

2. Overfitting happens when a model becomes too complex and learns even the noise in the training data, leading to low bias but high variance.

The bias-variance trade-off involves adjusting the model's complexity to avoid both underfitting and overfitting, optimizing its ability to generalize to new data.

Evaluating Models: To measure a model's performance, we use rigorous evaluation metrics that quantify its accuracy and ability to generalize. In classification tasks (e.g., determining spam or disease), the chosen metric depends on the problem, ensuring that the model is assessed objectively for its predictive quality and generalization ability.

***Accuracy:***
*Definition:* The proportion of all correct classifications (both true positives and true negatives) relative to the total number of input samples.

Formula for accuracy is given as;
$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN},$$

*Interpretation***:** Provides a quick overview of overall correctness. However, it can be misleading for datasets with imbalanced class distributions, where a model might achieve high accuracy by simply predicting the majority class.

*Precision:*
*Definition***:** The proportion of all positive predictions made by the model that are actually correct (True Positives). Formula: $Accuracy = \frac{TP}{TP + FP}$

*Interpretation:* Measures the exactness or quality of positive predictions. High precision indicates a low rate of false positives, crucial when false alarms are costly (e.g., flagging legitimate transactions as fraudulent).
*Recall (Sensitivity, True Positive Rate - TPR):*

*Definition:* The proportion of all actual positive instances that were correctly identified by the model.
$$Accuracy = \frac{TP}{TP + FN}$$

*Interpretation:* Measures the completeness or coverage of positive predictions. High recall indicates a low rate of false negatives and is crucial when false negatives carry high costs (e.g., missing a disease diagnosis) [44].

*F1-Score:*

Definition: The harmonic means of precision and recall, providing a single metric that balances both.

Formula: $F_1 = 2 \Box \frac{\text{Pr}ecision \times \text{Re}call}{\text{Pr}ecision + \text{Re}call} = \frac{2TP}{2TP + FP + FN}$

*Precision and Recall*: These metrics are especially useful for class-imbalanced datasets, as they account for both false positives and false negatives. A value of 1.0 represents perfect precision and recall.

Confusion Matrix: This N×N matrix displays the performance of a classification model by showing the counts of True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN). It provides a detailed breakdown of the correct and incorrect classifications, serving as the foundation for calculating other metrics.

AUC-ROC: The ROC curve plots the True Positive Rate (Recall) against the False Positive Rate (FPR) at various threshold settings. The AUC quantifies the classifier's overall performance, ranging from 0 to 1. An AUC of 1.0 means a perfect classifier, while 0.5 indicates random guessing. AUC is valuable for comparing models, but for imbalanced datasets, Precision-Recall curves can provide a more informative evaluation.

For regression models, Mean Squared Error (MSE) is commonly used. It measures the average squared difference between predicted and actual values, quantifying the magnitude of prediction errors.

Formula: $MSE = \frac{1}{N} \sum_{i=1}^{N} (y_i - \widehat{y}_i)^2$

Interpretation: MSE penalizes larger errors more heavily due to the squaring operation, making it sensitive to outliers. It's a good choice when large errors are particularly undesirable.

Root Mean Squared Error (RMSE):
Definition: The square root of the Mean Squared Error.

Formula: $RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (y_i - \widehat{y}_i)^2}$

Interpretation: RMSE has the same units as the target variable, which often makes it more interpretable than MSE, as it represents the typical magnitude of error. Like MSE, it is also sensitive to outliers. RMSE is frequently preferred over Mean Absolute Error (MAE) when the errors are assumed to be normally distributed.

*R-squared ($R^2$):*

Definition: Measures the proportion of the variance in the dependent (target) variable that is predictable from the independent (feature) variables. It quantifies how well the model explains the variability in the data.

*Formula*: $R^2 = 1 - \dfrac{Variance(error)}{Variance(y)} = 1 - \dfrac{MSE}{Var(y)}$.

Interpretation: $R^2$ values typically range from 0 to 1, with higher values indicating a better fit of the model to the data. A value of 1 implies that the model perfectly explains all the variance in the target variable.

Evaluating a model's performance goes beyond measuring raw errors. It involves choosing the right error measurement technique based on the data's error distribution. For instance, MSE and RMSE are ideal for errors following a normal distribution, where larger errors are more impactful, while MAE is better for distributions with long tails, as it's more forgiving of outliers. Additionally, R-squared (R2) helps assess how well the model explains the data variation and its ability to generalize to new data, providing insights into its predictive power beyond just error metrics.

## 5. Applications of Machine Learning
### Real-World Impact in Healthcare:
Machine learning is transforming healthcare by enabling personalized treatments based on an individual's genetics, lifestyle, and environment. It's especially powerful in predictive diagnostics, such as detecting diseases early, and in accelerating drug discovery, which reduces development time. However, the complexity of personalized medicine and the high stakes of healthcare require ML models to be accurate, trustworthy, and transparent. Challenges like the "black box" problem and the bias-variance trade-off are critical, as they influence the safety, fairness, and effectiveness of ML in clinical settings.

### Finance:
In algorithmic trading, ML models analyse vast amounts of financial data to make quick, informed decisions, minimizing human error and optimizing returns. In fraud detection, ML efficiently sifts through transaction data to spot suspicious activity, improving security and reducing financial losses. Credit risk assessment benefits from ML's ability to evaluate a wider range of data, resulting in faster, more inclusive lending. Portfolio management also benefits from ML, which interprets market data to enhance returns and minimize risks. Reinforcement learning is especially valuable here, as it allows models to adaptively make decisions that maximize rewards in dynamic market conditions. The mathematical efficiency of ML in processing fast-moving financial data is key to its transformative role in the finance sector.

### Natural Language Processing (NLP):
NLP enables machines to understand and generate human language. Early systems were rule-based, but with the shift to statistical models and deep learning, NLP capabilities—like sentiment analysis, machine translation, and text generation—have vastly improved. However, converting language data into high-dimensional numerical representations presents challenges like the curse of dimensionality. Techniques such as Latent Dirichlet Allocation (LDA) help manage large datasets by identifying underlying topics, making the data more digestible and interpretable.

### COMPUTER VISION:
Computer vision empowers machines to interpret visual data. Key tasks like image recognition, object detection, and image segmentation rely on Convolutional Neural Networks (CNNs). CNNs are effective because they learn hierarchical features, starting from simple elements like edges and textures and progressing to complex representations such as objects. This layered learning process, enabled by deep learning, is crucial for applications in autonomous vehicles and healthcare diagnostics. The success of computer vision models relies on the precise application and continuous refinement of mathematical principles, ensuring high reliability and performance in real-world scenarios.

## 6. CONCLUSIONS
At its core, machine learning (ML) isn't just supported by math; it *is* math. Think of mathematics as the fundamental language and operating system for all of AI. Every time an ML algorithm "learns"—whether it's predicting outcomes, finding hidden patterns, or figuring things out through trial and error—what it's really doing is solving an optimization problem. This means it's constantly trying to make something as good as it can possibly be.

Different branches of math each play a vital role. Linear algebra is like the backbone, giving ML models the structure they need to handle and manipulate vast amounts of data. Then there's calculus, especially its gradient-based methods and that clever backpropagation algorithm; that's the engine that lets models fine-tune themselves and learn with remarkable precision. And finally, probability and statistics are essential for dealing with all the real-world uncertainty, making sure the models can draw solid conclusions from noisy data and perform reliably even on information they've never seen before.

These mathematical ideas aren't just separate tools; they weave together in every ML algorithm. For instance, whether it's the straightforward approach of linear regression, the probabilistic nature of logistic regression, the geometric elegance of Support Vector Machines, the smart decision-making of decision trees, the grouping power of K-Means, the pattern-finding ability of PCA, or the complex, multi-layered learning in deep neural networks—it all comes back to these underlying mathematical principles.

Why does all this math matter so much? Well, optimization is universal across all of ML, giving

seemingly different algorithms a common mathematical goal. A deep understanding of these principles isn't just for academics; it's critical for building trust and ensuring that complex AI systems work reliably, especially when they're used in sensitive areas. Math provides the inherent structure for data, explains the fundamental "how" and "why" of learning, and helps us ensure our models are robust and generalize well. It's also about finding the right balance between how efficient an algorithm is and how confident we can be it's found the best possible solution.

In the real world, from revolutionizing healthcare to optimizing finance, understanding human language, or empowering computer vision, ML's mathematical foundations are what allow it to tackle problems characterized by immense data, high complexity, and dynamic environments. The ongoing evolution of ML itself is often driven by the need to overcome tough mathematical challenges, leading to even more efficient, robust, and powerful models. Ultimately, the extraordinary power, versatility, and transformative potential of machine learning are deeply rooted in its profound mathematical foundations. A strong grasp of these principles isn't just about theoretical advancement; it's absolutely essential for the practical development, responsible deployment, and ethical application of intelligent systems in our increasingly data-driven world.

# REFERENCES

1. Arora, S. (2018). Mathematics of machine learning: an introduction. In Proceedings of the International Congress of Mathematicians: Rio de Janeiro 2018 (pp. 377-390).
2. Boryshchak, Y. (2020). Mathematical Perspective of Machine Learning. arXiv preprint arXiv:2007.01503.
3. Deisenroth, M. P., Faisal, A. A., & Ong, C. S. (2020). Mathematics for machine learning. Cambridge University Press.
4. El Khatib, O., & Alkhatib, N. A Comprehensive Overview of Kernels in Machine Learning: Mathematical Foundations and Applications. International Journal of Computer (IJC), 53(1), 150-172.
5. Gabriel, F., Signolet, J., &Westwell, M. (2018). A machine learning approach to investigating the effects of mathematics dispositions on mathematical literacy. International Journal of Research & Method in Education, 41(3), 306-327.
6. Knights, V., &Prchkovska, M. (2024). From equations to predictions: Understanding the mathematics and machine learning of multiple linear regression. J. Math. Comput. Appl, 3(2), 1-8.
7. Kumar, G., Banerjee, R., Kr Singh, D., & Choubey, N. (2020). Mathematics for machine learning. Journal of Mathematical Sciences & Computational Mathematics, 1(2), 229-238.
8. Essang, S. O., Ante, J. E., Otobi, A. O., Okeke, S., I., Akpan, U. D., Francis, R. E., Auta, J. T., Essien, D. E., Fadugba, S. E., Kolawole, O. M., Asanga, E. E., & Ita, B., I. (2025). Optimizing Neural Networks with Convex Hybrid Activations for Improved Gradient Flow. Abjournals, 5(1), 10–27. https://doi.org/10.52589/ajste-uobyfv1b
9. Lamba, S., Saini, P., Kukreja, V., & Sharma, B. (2021, April). Role of mathematics in machine learning. In Proceedings of the International Conference on Innovative Computing & Communication (ICICC).
10. Rajendra, P., Ravi, P. V., & Meenakshi, K. (2024, August). Machine learning from a mathematical perspective. In AIP Conference Proceedings (Vol. 3149, No. 1). AIP Publishing.
11. Shi, B., & Iyengar, S. S. (2020). Mathematical theories of machine learning-Theory and applications. Springer International Publishing.
12. Wang, J., & Zhang, W. (2020). Fuzzy mathematics and machine learning algorithms application in educational quality evaluation model. Journal of Intelligent & Fuzzy Systems, 39(4), 5583-5593.
13. Wilmott, P. (2022). Machine learning: an applied mathematics introduction. Machine Learning and the City: Applications in Architecture and Urban Design, 217-248.
14. Okon, E. S., Michael, K. O., Friday, R. E., Efiong, A. J., Obi, O. M., Timothy, A. J., Edet, O. P., Dominic, E. R., & Jimmy, U. A. (2024). Application of AI algorithms for the prediction of the likelihood of sickle cell crises. Scholars Journal of Engineering and Technology, 12(12), 394–403. https://doi.org/10.36347/sjet.2024.v12i12.008
15. Weinan, E. (2022). A mathematical perspective of machine learning. Plenary LectureS, 2, 914-954.
16. Wilmott, P. (2022). Machine learning: an applied mathematics introduction. Machine Learning and the City: Applications in Architecture and Urban Design, 217-248.
17. Wojtowytsch, S. (2024). Mathematics of Machine Learning: An Introduction. Pittsburgh Interdisciplinary Mathematics Review, 2, 1-25.
18. Shi, B., Iyengar, S. S. (2019). Mathematical Theories of Machine Learning - Theory and Applications. Germany: Springer International Publishing.
19. Deisenroth, M. P., Faisal, A. A., Ong, C. S. (2020). Mathematics for Machine Learning. United Kingdom: Cambridge University Press.
20. Kneusel, R. T. (2021). Math for Deep Learning: What You Need to Know to Understand Neural Networks. United States: No Starch Press.
21. Chaudhury, K. (2024). Math and Architectures of Deep Learning. United States: Manning.
22. Dawani, J. (2020). Hands-On Mathematics for Deep Learning: Build a Solid Mathematical Foundation for Training Efficient Deep Neural Networks. Germany: Packt Publishing.
23. Essang, S. O., Ante, J. E., Fadugba, S. E., Auta, J. T., Ezeorah, J. N., Francis, R. E., &Otobi, A. O. (2025). Optimizing neural networks with linearly

combined activation functions: A novel approach to enhance gradient flow and learning dynamics. International Journal of Mathematical Sciences and Optimization: Theory and Applications, 11(2), 29–44.

24. Claster, W. (2020). Mathematics and Programming for Machine Learning with R: From the Ground Up. United Kingdom: CRC Press.

25. Kroese, D. P., Botev, Z. I., Taimre, T., Vaisman, R. (2019). Data Science and Machine Learning: Mathematical and Statistical Methods. United States: Chapman & Hall/CRC.

26. Hack, S. (2019). Machine Learning Mathematics: Study Deep Learning Through Data Science. How to Build Artificial Intelligence Through Concepts of Statistics, Algorithms, Analysis and Data Mining. United Kingdom: Amazon Digital Services LLC - KDP Print US.

27. Suzuki, J. (2021). Statistical Learning with Math and Python: 100 Exercises for Building Logic. South Korea: Springer Nature Singapore.

28. Aggarwal, C. C. (2020). Linear Algebra and Optimization for Machine Learning: A Textbook. Germany: Springer International Publishing.

29. Calin, O. (2020). Deep Learning Architectures: A Mathematical Approach. Germany: Springer International Publishing.

30. Goodfellow, I., Bengio, Y., Courville, A. (2016). Deep Learning. United Kingdom: MIT Press.

31. Shalev-Shwartz, S., Ben-David, S. (2014). Understanding Machine Learning: From Theory to Algorithms. United Kingdom: Cambridge University Press.

32. Dua, R., Ghotra, M. S., Pentreath, N. (2017). Machine Learning with Spark - Second Edition. United Kingdom: Packt Publishing.

33. Essang, S. O., Okeke, S. I., Effiong, J. A., Francis, R., Fadugba, S. E., Otobi, A. O., Auta, J. T., Chukwuka, C. F., Ogar-Abang, M. O., & Moses, A. (2025). Adaptive hybrid optimization for backpropagation neural networks in image classification. Proceedings of the Nigerian Society of Physical Sciences, 150. https://doi.org/10.61298/pnspsc.2025.2.150

34. Howard, J., Gugger, S. (2020). Deep Learning for Coders with Fastai and PyTorch. Japan: O'Reilly Media.

35. The Mathematics of Data. (2018). United States: American Mathematical Society.

36. Simovici, D. A. (2018). Mathematical Analysis for Machine Learning and Data Mining. Japan: World Scientific Publishing Company.

37. Chen, L. M., Su, Z., Jiang, B. (2015). Mathematical Problems in Data Science: Theoretical and Practical Methods. Germany: Springer International Publishing.

38. Knox, S. W. (2018). Machine Learning: A Concise Introduction. United Kingdom: Wiley.

39. Blum, A., Hopcroft, J., Kannan, R. (2020). Foundations of Data Science. India: Cambridge University Press.

40. Kepner, J., Jananthan, H. (2018). Mathematics of Big Data: Spreadsheets, Databases, Matrices, and Graphs. (n.p.): MIT Press.

41. Conway, D., White, J. M. (2012). Machine Learning for Hackers. United States: O'Reilly Media.

42. Halmos, P. R. (1995). Linear Algebra Problem Book. United Kingdom: Mathematical Association of America.

43. Bohn, B., Garcke, J., Griebel, M. (2024). Algorithmic Mathematics in Machine Learning. United States: Society for Industrial and Applied Mathematics.

44. Ni, H., Dong, X., Yu, G., Zheng, J. (2021). An Introduction to Machine Learning in Quantitative Finance. Singapore: World Scientific.

45. ] Essang, S. O., Ante, J. E., Francis, R. E., Otobi, A. O., Ita, B. I., Kolawole, O. M., Aigberemhon, E. M., Oluwagbemi, J. T., Fadugba, E. S., Essien, D. E., Ogar-Abang, M. O., Auta, J. T., & Asanga, E. E. (2025). Digital literacy awareness to enhance the learning of mathematics. Computing and Applied Sciences Impact, 2(2), 7–15. https://www.researchgate.net/publication/391572279

46. Dudek, G. Computational Intelligence and Machine Learning: Advances in Models and Applications. Electronics 2025, 14, 1530. https://doi.org/10.3390/electronics14081530