Mathematical Prerequisites for Machine Learning Linear Algebra, Probability & Statistics, Calculus

Sarwan Ali

Department of Computer Science Georgia State University

🖬 The Mathematical Foundation of Intelligence 🖬

<ロト < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 0 へ () 1/27

- Linear Algebra: The Language of Data
- Probability & Statistics: Handling Uncertainty
- 3 Calculus: The Mathematics of Change
- Putting It All Together: ML Applications
- 5 Practical Tips & Next Steps

Why Does Math Matter for Machine Learning?

Think of ML as a Recipe

Traditional Cooking:

- Follow instructions step by step
- Measure ingredients
- Apply heat and time
- Get consistent results

Machine Learning:

- Linear Algebra: Organize ingredients (data)
- Probability: Handle uncertainty in cooking
- Calculus: Optimize the recipe automatically
- Get intelligent results!

Key Point

Math gives us the language to describe and solve learning problems systematically!

Linear Algebra: Why It's Everywhere in ML

Data is just numbers organized in arrays!

Real-world data as vectors/matrices:

- Images: Pixel intensity arrays
- Text: Word frequency vectors
- Audio: Wave amplitude sequences
- Customer data: Feature vectors

Example: House Price Prediction					
House =	[2000]				
	3				
	2				
	1995				
(sq ft, bedrooms, bathrooms, year built)					

Bottom Line

Linear algebra lets us manipulate ALL this data efficiently with simple operations!

Vectors: The Building Blocks

What is a Vector?

- A list of numbers
- Represents a point in space
- Or a direction and magnitude

Notation:

$$\mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

Example:

$$student = \begin{bmatrix} 85\\92\\78\\88 \end{bmatrix}$$

(Math, Science, English, History grades)



A D > A B > A B > A B >

Vector Operations:

- Addition: $\mathbf{a} + \mathbf{b}$
- Scalar multiplication: ca
- Dot product: $\mathbf{a} \cdot \mathbf{b}$

Essential Vector Operations

1. Dot Product (Inner Product):

$$\mathbf{a} \cdot \mathbf{b} = \sum_{i=1}^{n} a_i b_i = a_1 b_1 + a_2 b_2 + \dots + a_n b_n$$

- Geometric meaning: Measures similarity/correlation
- 2. Vector Norm (Length):

$$\|\mathbf{v}\| = \sqrt{\mathbf{v} \cdot \mathbf{v}} = \sqrt{v_1^2 + v_2^2 + \dots + v_n^2}$$

3. Unit Vector:

$$\hat{\mathbf{v}} = rac{\mathbf{v}}{\|\mathbf{v}\|}$$

ML Connection



イロト 不同 トイヨト イヨト

What is a Matrix?

- 2D array of numbers
- Collection of vectors
- Represents datasets or transformations Notation:

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

Size: $m \times n$ (rows × columns)

Dataset Example

Student grades matrix:

	[85	92	78
Tradas —	90	88	95
srades =	76	84	82
	88	91	89

- Rows = Students
- Columns = Subjects
- Each row is a student vector
- Each column is a subject vector

Matrix Operations That Power ML (Alternative)

1. Matrix Multiplication: AB = C

$$\mathbf{A}_{2\times 2} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \quad \mathbf{B}_{2\times 2} = \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} = \mathbf{C}_{2\times 2} = \begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix}$$

Where: $c_{ij} = \sum_k a_{ik} b_{kj}$

Matrix Transpose: A^T (flip rows and columns)
 Matrix Inverse: A⁻¹ (undoes the transformation)

ML Application

Linear regression: $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ (finds best fit line!)

The Big Idea: Some vectors don't change direction when transformed by a matrix!

 $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$

- v: eigenvector (special direction)
- λ : eigenvalue (scaling factor)

Why Care?

- Reveals the "natural" directions in data
- Used in dimensionality reduction (PCA)
- Helps understand data structure
- Critical for many ML algorithms



Same direction, different length!

PCA Example

Find directions of maximum variance in data to reduce dimensions while preserving information.

Probability: The Science of Uncertainty

Real world is messy and uncertain!

Why Probability in ML?

- Data has noise and errors
- Predictions are uncertain
- Models make assumptions
- Need to quantify confidence

Examples:

- "70% chance of rain"
- "95% confidence interval"
- "Spam probability: 0.85"



Probability Distribution

Key Insight

ML algorithms learn probability distributions from data to make informed predictions!

Probability Basics

Probability Rules:

- $0 \le P(A) \le 1$ (probabilities are between 0 and 1)
- P(certain event) = 1
- *P*(impossible event) = 0

•
$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

Conditional Probability:

$$P(A|B) = rac{P(A ext{ and } B)}{P(B)}$$

"Probability of A given B has occurred" **Bayes' Theorem:**

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

Spam Detection

- P(spam|" free money") = ?
- Given email contains "free money", what's probability it's spam?
- Use Bayes' theorem with training data!

 $P(\text{spam}|\text{words}) = rac{P(\text{words}|\text{spam})P(\text{spam})}{P(\text{words})}$

Important Probability Distributions

1. Normal (Gaussian) Distribution:

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- Bell-shaped curve
- Characterized by mean μ and variance σ^2
- Central Limit Theorem
- Used everywhere in ML!

2. Bernoulli Distribution:

- Binary outcomes (0 or 1)
- Probability *p* of success
- Models coin flips, binary classification

Normal Distribution



ML Applications

• Normal: Linear regression errors, feature distributions

(a)

• Bernoulli: Logistic regression, binary classification

Statistics: Making Sense of Data

Descriptive Statistics:

- Mean: $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$
- Median: Middle value when sorted
- Mode: Most frequent value
- Variance: $\sigma^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i \bar{x})^2$
- Standard Deviation: $\sigma = \sqrt{\sigma^2}$

Why Important?

- Understand your data first!
- Detect outliers and anomalies
- Choose appropriate algorithms

Inferential Statistics:

- **Hypothesis Testing**: Is this result significant?
- **Confidence Intervals**: Range of plausible values
- p-values: Evidence against null hypothesis
- Correlation: Linear relationship strength

A/B Testing

- Test two versions of ML model
- Measure performance difference
- Use statistics to determine if difference is significant
- Make data-driven decisions!

How do machines learn?

They optimize! And optimization needs calculus.

The Learning Process:

- Start with random model
- Ø Measure how wrong it is (loss function)
- Find direction to improve
- Take a step in that direction
- Sepeat until optimal



Derivatives: The Rate of Change

What is a Derivative?

$$f'(x) = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

- Slope of the tangent line
- Rate of change at a point
- Direction of steepest increase **Common Derivatives:**

$$\frac{d}{dx}[x^n] = nx^{n-1} \tag{1}$$

$$\frac{d}{dx}[e^x] = e^x \tag{2}$$

$$\frac{d}{dx}[\ln x] = \frac{1}{x} \tag{3}$$

$$\frac{d}{dx}[\sin x] = \cos x \tag{4}$$



ML Connection

Gradient Descent:

- Calculate derivative of loss function
 - Move in opposite direction (negative gradient)
 - Minimize error automatically!

$$w_{new} = w_{old} - lpha rac{\partial Loss}{\partial w}$$

15 / 27

Partial Derivatives & Gradients: Multivariable Optimization

Partial Derivatives:

For function f(x, y):

$$\frac{\partial f}{\partial x} = \text{rate of change w.r.t. } x \qquad (5)$$
$$\frac{\partial f}{\partial y} = \text{rate of change w.r.t. } y \qquad (6)$$

Gradient Vector:

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}$$

- Points in direction of steepest increase
- Magnitude = rate of steepest increase
- Perpendicular to level curves



Gradient points away from minimum

Neural Networks

Backpropagation uses chain rule to compute gradients of loss w.r.t. all parameters efficiently!

(a)

Chain Rule: If y = f(g(x)), then:

 $\frac{dy}{dx} = \frac{dy}{dg} \cdot \frac{dg}{dx}$

Multiple Variables: If z = f(x, y) where x = g(t) and y = h(t):

$$\frac{dz}{dt} = \frac{\partial z}{\partial x}\frac{dx}{dt} + \frac{\partial z}{\partial y}\frac{dy}{dt}$$

Why Critical for ML?

- Neural networks are compositions $f_n(f_{n-1}(...f_1(x)))$
- Need gradients w.r.t. all parameters
- Chain rule propagates error backwards

Backpropagation
$$\frac{\partial Loss}{\partial w_1} = \frac{\partial Loss}{\partial y} \cdot \frac{\partial y}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_1}$$

Optimization Problem:

 $\min_w f(w)$

- Where f(w) is the loss/cost function **Critical Points:**
 - $\nabla f(w) = 0$ (gradient is zero)
 - Could be minimum, maximum, or saddle point

Second Derivative Test:

- f''(w) > 0: Local minimum
- f''(w) < 0: Local maximum
- f''(w) = 0: Inconclusive

Gradient Descent Algorithm:

Initialize parameters w₀

2 For
$$t = 0, 1, 2, ...$$
:

- Compute gradient: $g_t = \nabla f(w_t)$
- Update: $w_{t+1} = w_t \alpha g_t$
- Stop when converged

Learning Rate

- α too large: Overshooting
- α too small: Slow convergence
- Need to tune carefully!

Problem: Predict house prices from features **Linear Algebra:**

- Data matrix **X** (houses \times features)
- Target vector **y** (prices)
- Parameter vector **w** (weights)
- Prediction: $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$

Statistics:

- ullet Assume errors \sim Normal distribution
- Least squares estimation
- R-squared for model evaluation

Calculus:

- Loss function: $L(\mathbf{w}) = \|\mathbf{X}\mathbf{w} \mathbf{y}\|^2$
- Take derivative: $\frac{\partial L}{\partial \mathbf{w}} = 2\mathbf{X}^T (\mathbf{X}\mathbf{w} \mathbf{y})$

• Set to zero:
$$\mathbf{X}^{\mathcal{T}}(\mathbf{X}\mathbf{w} - \mathbf{y}) = 0$$

• Solve:
$$\mathbf{w} = (\mathbf{X}^{\mathsf{T}} \mathbf{X})^{-1} \mathbf{X}^{\mathsf{T}} \mathbf{y}$$

All Three Together!

Linear algebra for data representation, statistics for assumptions, calculus for optimization.

Example 2: Logistic Regression

Problem: Classify emails as spam or not spam

The Model:

$$P(\mathsf{spam}|\mathbf{x}) = rac{1}{1 + e^{-\mathbf{w}^{ au}\mathbf{x}}}$$

Linear Algebra:

- Feature vectors **x** (word counts, etc.)
- Weight vector **w**
- Linear combination $\mathbf{w}^T \mathbf{x}$

Probability:

- Sigmoid function maps to [0,1]
- Bernoulli distribution for binary outcome
- Maximum likelihood estimation

Calculus Optimization:

Loss function (negative log-likelihood):

$$L(\mathbf{w}) = -\sum_{i=1}^{n} [y_i \log p_i + (1 - y_i) \log(1 - p_i)]$$

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{X}^{\mathsf{T}}(\mathbf{p} - \mathbf{y})$$

Update rule:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \mathbf{X}^T (\mathbf{p} - \mathbf{y})$$



Example 3: Principal Component Analysis (PCA)

Problem: Reduce data dimensions while preserving information Mathematical Foundation:

The Algorithm:

- **O** Center the data: $\mathbf{X}_c = \mathbf{X} \bar{\mathbf{X}}$
- 2 Compute covariance: $\mathbf{C} = \frac{1}{\pi} \mathbf{X}_c^T \mathbf{X}_c$
- Sind eigenvalues/eigenvectors of С
- Sort by eigenvalue magnitude
- **5** Keep top k eigenvectors as principal components
- Project data: $\mathbf{Y} = \mathbf{X}_{c} \mathbf{W}_{k}$

Linear Algebra:

- Eigendecomposition: $\mathbf{C}\mathbf{v} = \lambda\mathbf{v}$
- Matrix multiplication for projection
- Orthogonal transformation

Statistics:

- Variance measures information content
- Covariance shows feature relationships
- Principal components capture maximum variance

Applications

- Image compression, Data visualization
- Feature extraction. Noise reduction

Example 4: Neural Networks - The Full Package

Neural Networks Use ALL Mathematical Concepts!

Forward Pass (Linear Algebra):

$$\mathbf{z}^{(1)} = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$$
 (7)
 $\mathbf{a}^{(1)} = \sigma(\mathbf{z}^{(1)})$ (8)

$$\mathbf{z}^{(2)} = \mathbf{W}^{(2)}\mathbf{a}^{(1)} + \mathbf{b}^{(2)}$$
 (9)

$$\mathbf{y} = \operatorname{softmax}(\mathbf{z}^{(2)})$$
 (10)

Probability:

- Softmax outputs probabilities
- Cross-entropy loss
- Dropout for regularization
- Bayesian neural networks

The Beauty of Neural Networks

They seamlessly combine linear algebra (efficient computation), probability (uncertainty handling), and calculus (automatic optimization) into one powerful framework! 22/27

Backward Pass (Calculus): $\frac{\partial L}{\partial \mathbf{W}^{(2)}} = \frac{\partial L}{\partial \mathbf{z}^{(2)}} (\mathbf{a}^{(1)})^{T} \qquad (11)$ $\frac{\partial L}{\partial \mathbf{a}^{(1)}} = (\mathbf{W}^{(2)})^{T} \frac{\partial L}{\partial \mathbf{z}^{(2)}} \qquad (12)$ $\frac{\partial L}{\partial \mathbf{W}^{(1)}} = \frac{\partial L}{\partial \mathbf{z}^{(1)}} \mathbf{x}^{T} \qquad (13)$

Chain Rule Everywhere:

 $\frac{\partial L}{\partial \mathbf{W}^{(1)}} = \frac{\partial L}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{z}^{(2)}} \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{a}^{(1)}} \frac{\partial \mathbf{a}^{(1)}}{\partial \mathbf{z}^{(1)}} \frac{\partial \mathbf{z}^{(1)}}{\partial \mathbf{W}^{(1)}}$

How to Master These Mathematical Concepts

Study Strategy:

- Start with intuition Understand the "why"
- Practice calculations Build mechanical skills
- Ode implementations Solidify understanding
- Apply to ML problems See connections
- **5** Teach others Test your knowledge

Don't Fear the Math!

- You don't need to prove theorems
- Focus on understanding and application
- Use computational tools when appropriate

Recommended Resources: Books:

- "Mathematics for Machine Learning" -Deisenroth et al.
- "Linear Algebra Done Right" Axler
- "Introduction to Statistical Learning" -James et al.

Online:

- Khan Academy (basics)
- 3Blue1Brown (visual explanations)
- MIT OpenCourseWare
- Coursera/edX math courses

Practice:

Common Misconceptions & How to Avoid Them

Misconception 1: "I need to memorize all formulas"

Reality: Understanding concepts *i* memorizing formulas. Focus on intuition and derive when needed.

Misconception 2: "Linear algebra is just matrix arithmetic"

Reality: It's about understanding transformations, spaces, and geometric relationships.

Misconception 3: "Statistics is just calculating means and variances"

Reality: It's about reasoning under uncertainty and making valid inferences from data.

Misconception 4: "Calculus is just taking derivatives"

Reality: It's about understanding rates of change and optimization in high dimensions.

Remember: Math is a tool to solve problems, not an end in itself!

Your Mathematical Journey in Machine Learning



Timeline Estimate:

- Foundations: 2-3 months
- Basic ML: 2-3 months
- Advanced: 6-12 months
- Research: Ongoing!

Key Success Factors:

- Consistent daily practice
- Balance theory with implementation
- Join study groups/communities
- Work on real projects

Summary: The Mathematical Toolkit for ML

You now have the roadmap to mathematical mastery in ML!

Ħ

Linear Algebra

- Data representation
- Efficient computation
- Transformations
- Dimensionality reduction

Probability & Statistics

- Handle uncertainty
- Model assumptions
- Inference and testing
- Confidence quantification



- Optimization
- Gradient descent
- Backpropagation
- Parameter tuning

Final Message

Mathematics is not a barrier to machine learning - it's your superpower! Start with one concept, practice consistently, and watch as the entire ML landscape opens up before you.

Questions? Let's dive deeper into the mathematics of intelligence!

Thank You!

Mathematical Prerequisites for Machine Learning

"The only way to learn mathematics is to do mathematics." - Paul Halmos

Ready to start your mathematical journey in ML? => = onc