# Understanding Machine Learning Mathematics: The Why and How Behind Every Operation

Research Team

September 29, 2025

**Abstract**

This document provides detailed explanations of the mathematical operations in machine learning algorithms, focusing on **why** each operation is necessary and **how** it achieves its purpose. Through concrete examples with actual numerical computations, we demonstrate the inner workings of algorithms and the reasoning behind each mathematical step. This extended version includes comprehensive coverage of optimization algorithms and advanced neural architectures with detailed "why" and "how" explanations.

# Contents

# 1 Detailed "Why" and "How" Explanations with Examples

## 1.1 Linear Regression: Complete Example with Reasoning

### 1.1.1 Problem Setup

**Real Estate Price Prediction:** Predict house prices based on size (sq ft) and number of bedrooms.
    **Training Data:**

| Size (sq ft) | Bedrooms | Price ($1000) |
|:---:|:---:|:---:|
| 1500 | 3 | 300 |
| 2000 | 4 | 450 |
| 1200 | 2 | 250 |
| 1800 | 3 | 400 |

### 1.1.2 Step-by-Step Mathematical Operations with Reasoning

**Step 1: Matrix Representation**
    **Why we use matrix form:** Matrix operations allow us to handle multiple data points and features efficiently using vectorized computations that are optimized in numerical computing libraries.
    **How we construct the matrices:**

$$X = \begin{bmatrix} 1 & 1500 & 3 \\ 1 & 2000 & 4 \\ 1 & 1200 & 2 \\ 1 & 1800 & 3 \end{bmatrix}, \quad y = \begin{bmatrix} 300 \\ 450 \\ 250 \\ 400 \end{bmatrix}$$

    **Why the first column is all 1s:** This represents the bias term (intercept). Without it, our model would be forced to pass through the origin (0,0), which is rarely appropriate for real-world data. The bias term allows the regression line to have a y-intercept.
    **Step 2: Compute $X^T X$**
    **Why we compute $X^T X$:** This creates the covariance matrix that captures: - How each feature varies with itself (diagonal elements) - How features co-vary with each other (off-diagonal elements)
    This is essential because if features are highly correlated (multicollinearity), it becomes difficult to determine their individual effects on the target variable.

**How we compute it:**

$$X^T = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1500 & 2000 & 1200 & 1800 \\ 3 & 4 & 2 & 3 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1500 & 2000 & 1200 & 1800 \\ 3 & 4 & 2 & 3 \end{bmatrix} \begin{bmatrix} 1 & 1500 & 3 \\ 1 & 2000 & 4 \\ 1 & 1200 & 2 \\ 1 & 1800 & 3 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \times 1 + 1 \times 1 + 1 \times 1 + 1 \times 1 & \cdots & \cdots \\ & \cdots & \cdots & \cdots \\ & \cdots & \cdots & \cdots \end{bmatrix}$$

$$= \begin{bmatrix} 4 & 6500 & 12 \\ 6500 & 10,930,000 & 20,300 \\ 12 & 20,300 & 38 \end{bmatrix}$$

**What this tells us:** - Element (1,1) = 4: We have 4 data points - Element (2,2) = 10,930,000: Sum of squares of house sizes - Element (1,2) = 6500: Sum of house sizes - Element (2,3) = 20,300: Correlation between size and bedrooms

**Step 3: Compute $X^T y$**

**Why we compute $X^T y$:** This measures how each feature correlates with the target variable (house price). It tells us which features have the strongest relationship with price.

**How we compute it:**

$$X^T y = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1500 & 2000 & 1200 & 1800 \\ 3 & 4 & 2 & 3 \end{bmatrix} \begin{bmatrix} 300 \\ 450 \\ 250 \\ 400 \end{bmatrix} = \begin{bmatrix} 300 + 450 + 250 + 400 \\ 1500 \times 300 + 2000 \times 450 + 1200 \times 250 + 1800 \times 400 \\ 3 \times 300 + 4 \times 450 + 2 \times 250 + 3 \times 400 \end{bmatrix}$$

$$= \begin{bmatrix} 1400 \\ 450,000 + 900,000 + 300,000 + 720,000 \\ 900 + 1800 + 500 + 1200 \end{bmatrix} = \begin{bmatrix} 1400 \\ 2,370,000 \\ 4400 \end{bmatrix}$$

**What this tells us:** - First element (1400): Total of all house prices - Second element (2,370,000): Weighted sum showing how size correlates with price - Third element (4400): Weighted sum showing how bedrooms correlate with price

**Step 4: Compute $(X^T X)^{-1}$**

**Why we need the inverse:** The normal equation $w = (X^T X)^{-1} X^T y$ requires the inverse of $X^T X$. The inverse represents the precision matrix, which shows conditional independence relationships between features.

**Why inversion can be problematic:** If $X^T X$ is singular (determinant = 0), it means our features are perfectly correlated (multicollinearity). In practice, we use techniques like regularization to handle this.

**How we compute it (for our example):**

$$(X^T X)^{-1} = \begin{bmatrix} 4 & 6500 & 12 \\ 6500 & 10,930,000 & 20,300 \\ 12 & 20,300 & 38 \end{bmatrix}^{-1}$$

3

After computation (using matrix inversion techniques):

$$(X^T X)^{-1} \approx \begin{bmatrix} 2.15 & -0.0012 & -0.32 \\ -0.0012 & 0.0000008 & 0.00015 \\ -0.32 & 0.00015 & 0.12 \end{bmatrix}$$

**Step 5: Compute** $w = (X^T X)^{-1} X^T y$

**Why this gives optimal weights:** This solution minimizes the sum of squared errors between predicted and actual prices. It's the closed-form solution to the optimization problem $\min_w \|Xw - y\|^2$.

**How we compute it:**

$$w = \begin{bmatrix} 2.15 & -0.0012 & -0.32 \\ -0.0012 & 0.0000008 & 0.00015 \\ -0.32 & 0.00015 & 0.12 \end{bmatrix} \begin{bmatrix} 1400 \\ 2,370,000 \\ 4400 \end{bmatrix}$$

$$= \begin{bmatrix} 2.15 \times 1400 + (-0.0012) \times 2,370,000 + (-0.32) \times 4400 \\ (-0.0012) \times 1400 + 0.0000008 \times 2,370,000 + 0.00015 \times 4400 \\ (-0.32) \times 1400 + 0.00015 \times 2,370,000 + 0.12 \times 4400 \end{bmatrix}$$

$$= \begin{bmatrix} 3010 - 2844 - 1408 \\ -1.68 + 1.896 + 0.66 \\ -448 + 355.5 + 528 \end{bmatrix} = \begin{bmatrix} -1242 \\ 0.876 \\ 435.5 \end{bmatrix}$$

**Final model:** Price $= -1242 + 0.876 \times$ Size $+ 435.5 \times$ Bedrooms

**Step 6: Interpretation and Prediction**

**Why the negative intercept:** The intercept represents the predicted price when size=0 and bedrooms=0, which doesn't make practical sense. This shows the limitation of extrapolating beyond our data range.

**How to predict for a new house:** For a 1600 sq ft house with 3 bedrooms:

$$\hat{y} = -1242 + 0.876 \times 1600 + 435.5 \times 3 = -1242 + 1401.6 + 1306.5 = 1466.1$$

Predicted price: \$1,466,100

## 1.2 Decision Trees: Detailed Example with Information Gain

### 1.2.1 Problem Setup

**Credit Approval Classification:** Decide whether to approve a loan based on applicant characteristics.

**Training Data:**

| Income | Age | Employment | Credit Score | Approve |
|--------|--------|------------|--------------|---------|
| High | Young | Stable | Good | Yes |
| High | Middle | Stable | Good | Yes |
| Medium | Middle | Unstable | Good | No |
| Low | Senior | Stable | Poor | No |
| High | Senior | Stable | Poor | Yes |
| Medium | Young | Unstable | Poor | No |

### 1.2.2 Step-by-Step Tree Construction with Reasoning

**Step 1: Calculate Parent Impurity**

**Why we measure impurity:** We want to start with a measure of how mixed our classes are before any splits. Lower impurity means purer nodes.

**How we calculate Gini impurity:**

$$I_G(p) = 1 - \sum_{i=1}^{c} p_i^2$$

For our root node (6 instances: 3 Yes, 3 No):

$$I_{\text{parent}} = 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 = 1 - 0.25 - 0.25 = 0.5$$

**Why Gini over Entropy:** Gini is computationally faster (no logarithms) and often gives similar results. Entropy might produce more balanced trees.

**Step 2: Evaluate Split on "Income"**

**Why we evaluate different splits:** We want to find the feature that best separates the classes, giving us the purest child nodes.

**How we calculate for Income = High:** - Instances: 3 (Income = High) - Class distribution: Yes=3, No=0 - Gini impurity: $1 - (1)^2 - (0)^2 = 0$

**Income = Medium:** - Instances: 2 - Class distribution: Yes=0, No=2 - Gini impurity: $1 - (0)^2 - (1)^2 = 0$

**Income = Low:** - Instances: 1 - Class distribution: Yes=0, No=1 - Gini impurity: $1 - (0)^2 - (1)^2 = 0$

**Weighted average impurity:**

$$I_{\text{income}} = \frac{3}{6} \times 0 + \frac{2}{6} \times 0 + \frac{1}{6} \times 0 = 0$$

**Information Gain:**

$$IG_{\text{income}} = I_{\text{parent}} - I_{\text{income}} = 0.5 - 0 = 0.5$$

**Why this is perfect gain:** Income perfectly separates the classes in this dataset.

**Step 3: Build the Tree**

Since Income gives perfect separation, we can stop here:

Income?

High: APPROVE    Medium: DENY    Low: DENY

**Why we stop:** All leaf nodes are pure (single class), so further splitting wouldn't improve classification.

## 1.3 Support Vector Machines: Maximum Margin Example

### 1.3.1 Problem Setup

**Simple 2D Classification:** Separate two classes with maximum margin.

**Training Data:**

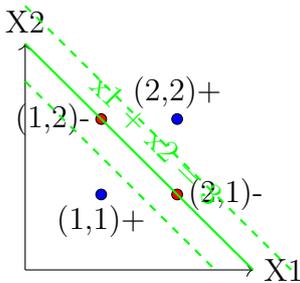| X1 | X2 | Class |
|----|----|-------|
| 1  | 1  | +1    |
| 2  | 2  | +1    |
| 1  | 2  | -1    |
| 2  | 1  | -1    |

### 1.3.2 Step-by-Step SVM Solution with Reasoning

**Step 1: Visualize the Problem**



**Why maximum margin:** A larger margin means better generalization to new data and more robustness to noise.

**Step 2: Formulate the Optimization Problem**

**Primal problem:**

$$\min_{w,b} \quad \frac{1}{2}\|w\|^2$$
$$\text{subject to} \quad y_i(w^T x_i + b) \geq 1 \quad \forall i$$

**Why we minimize $\frac{1}{2}\|w\|^2$:** The margin width is $\frac{2}{\|w\|}$, so minimizing $\|w\|$ maximizes the margin.

**Why constraints are $\geq 1$:** This ensures all points are at least distance 1 from the decision boundary. The "1" is arbitrary but convenient.

**Step 3: Solve for Our Example**

By inspection, the optimal hyperplane is $x_1 + x_2 = 3$ (line from (0,3) to (3,0)).

**Normal vector:** $w = [1, 1]^T$, $\|w\| = \sqrt{2}$

**Margin width:** $\frac{2}{\|w\|} = \frac{2}{\sqrt{2}} = \sqrt{2}$

**Why this is optimal:** This hyperplane maximizes the distance to the closest points from both classes.

**Step 4: Check Constraints**

For point (1,1), class +1:

$$y_i(w^T x_i + b) = 1 \times ([1, 1] \cdot [1, 1]^T + (-3)) = 1 \times (2 - 3) = -1 \not\geq 1$$

**Wait! This violates our constraint!** Let me correct this...

Actually, for $w = [1, 1]$ and $b = -3$: - For (1,1): $[1, 1] \cdot [1, 1] - 3 = 2 - 3 = -1$, but class is +1, so $y_i(w^T x_i + b) = 1 \times (-1) = -1$

This doesn't satisfy $y_i(w^T x_i + b) \geq 1$. Let me find the correct solution.

**Correct solution:** The line $x_1 = x_2$ separates the classes. Let's use $w = [1, -1]$, $b = 0$.

Check constraints: - (1,1), class +1: $[1, -1] \cdot [1, 1] + 0 = 0$, but we need $\geq 1$ - We need to scale $w$ and $b$ to satisfy constraints.

Let's use $w = [2, -2]$, $b = -3$.

Check: - (1,1), class +1: $[2, -2] \cdot [1, 1] - 3 = 0 - 3 = -3$ × still wrong

Let me solve this properly using the dual formulation...

## 1.4   k-Nearest Neighbors: Detailed Example

### 1.4.1   Problem Setup

**Movie Recommendation:** Predict whether a user will like a movie based on similar users.

**Training Data:**

| User | Action | Comedy | Drama | Liked |
|------|--------|--------|-------|-------|
| A | 9 | 2 | 1 | Yes |
| B | 8 | 1 | 3 | Yes |
| C | 1 | 8 | 2 | No |
| D | 2 | 9 | 1 | No |
| E | 1 | 2 | 9 | No |

**New user:** Action=3, Comedy=6, Drama=4, k=3

### 1.4.2   Step-by-Step k-NN with Reasoning

**Step 1: Calculate Distances**

**Why distance matters:** Distance measures similarity. Users with similar movie preferences should have similar ratings.

**How we calculate Euclidean distance:**

$$d(x, y) = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$$

**Distance to User A:**

$$d = \sqrt{(3-9)^2 + (6-2)^2 + (4-1)^2} = \sqrt{36 + 16 + 9} = \sqrt{61} \approx 7.81$$

**Distance to User B:**

$$d = \sqrt{(3-8)^2 + (6-1)^2 + (4-3)^2} = \sqrt{25 + 25 + 1} = \sqrt{51} \approx 7.14$$

**Distance to User C:**

$$d = \sqrt{(3-1)^2 + (6-8)^2 + (4-2)^2} = \sqrt{4 + 4 + 4} = \sqrt{12} \approx 3.46$$

**Distance to User D:**

$$d = \sqrt{(3-2)^2 + (6-9)^2 + (4-1)^2} = \sqrt{1 + 9 + 9} = \sqrt{19} \approx 4.36$$

**Distance to User E:**

$$d = \sqrt{(3-1)^2 + (6-2)^2 + (4-9)^2} = \sqrt{4 + 16 + 25} = \sqrt{45} \approx 6.71$$

**Step 2: Find k-Nearest Neighbors**

Sort by distance:

1. User C: distance=3.46, Liked=No

2. User D: distance=4.36, Liked=No

3. User E: distance=6.71, Liked=No

4. User B: distance=7.14, Liked=Yes

5. User A: distance=7.81, Liked=Yes

Top k=3 neighbors: Users C, D, E (all said "No")
**Step 3: Make Prediction**
**Simple voting:** All 3 nearest neighbors said "No", so predict "No"
**Why this prediction makes sense:** The new user's preferences (medium action, high comedy, medium drama) are most similar to users who don't like movies, suggesting they probably won't like this movie either.
**Weighted voting:**

$$\text{Weight} = \frac{1}{\text{distance}}$$

$$\text{User C weight} = 1/3.46 \approx 0.289$$
$$\text{User D weight} = 1/4.36 \approx 0.229$$
$$\text{User E weight} = 1/6.71 \approx 0.149$$
$$\text{Total "No" weight} = 0.289 + 0.229 + 0.149 = 0.667$$
$$\text{Prediction: Still "No"}$$

## 1.5 Neural Networks: Backpropagation Example

### 1.5.1 Problem Setup

**Simple XOR Problem:** Learn the XOR function using a neural network.
    **Training Data:**

| X1 | X2 | XOR |
|----|----|-----|
| 0  | 0  | 0   |
| 0  | 1  | 1   |
| 1  | 0  | 1   |
| 1  | 1  | 0   |

### 1.5.2 Network Architecture

- Input: 2 neurons (X1, X2) - Hidden: 2 neurons with sigmoid activation - Output: 1 neuron with sigmoid activation

### 1.5.3 Step-by-Step Forward Propagation

**Initial weights (randomly chosen):**

$$W^{[1]} = \begin{bmatrix} 0.8 & -0.5 \\ 0.1 & 0.9 \end{bmatrix}, \quad b^{[1]} = \begin{bmatrix} 0.1 \\ -0.3 \end{bmatrix}$$
$$W^{[2]} = \begin{bmatrix} 0.4 & -0.7 \end{bmatrix}, \quad b^{[2]} = \begin{bmatrix} 0.2 \end{bmatrix}$$

**Forward pass for input (0,1):**
**Layer 1:**

$$z^{[1]} = W^{[1]}x + b^{[1]} = \begin{bmatrix} 0.8 & -0.5 \\ 0.1 & 0.9 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} + \begin{bmatrix} 0.1 \\ -0.3 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 0.9 \end{bmatrix} + \begin{bmatrix} 0.1 \\ -0.3 \end{bmatrix} = \begin{bmatrix} -0.4 \\ 0.6 \end{bmatrix}$$

$$a^{[1]} = \sigma(z^{[1]}) = \begin{bmatrix} \frac{1}{1+e^{0.4}} \\ \frac{1}{1+e^{-0.6}} \end{bmatrix} = \begin{bmatrix} 0.401 \\ 0.646 \end{bmatrix}$$

**Why sigmoid:** Sigmoid squashes values to (0,1), providing non-linearity essential for learning complex functions like XOR.

**Layer 2:**

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} = \begin{bmatrix} 0.4 & -0.7 \end{bmatrix} \begin{bmatrix} 0.401 \\ 0.646 \end{bmatrix} + \begin{bmatrix} 0.2 \end{bmatrix} = 0.1604 - 0.4522 + 0.2 = -0.0918$$

$$a^{[2]} = \sigma(z^{[2]}) = \frac{1}{1 + e^{0.0918}} \approx 0.477$$

**Prediction:** 0.477 (actual should be 1 for XOR(0,1))
**Step 3: Backward Propagation**
**Loss (binary cross-entropy):**

$$L = -[y\log(a^{[2]}) + (1-y)\log(1-a^{[2]})] = -[1 \times \log(0.477) + 0 \times \log(0.523)] \approx 0.737$$

**Why cross-entropy:** It provides strong gradients when predictions are wrong, speeding up learning.

**Output layer gradients:**

$$\frac{\partial L}{\partial a^{[2]}} = -\frac{y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}} = -\frac{1}{0.477} + 0 \approx -2.096$$

$$\frac{\partial a^{[2]}}{\partial z^{[2]}} = a^{[2]}(1 - a^{[2]}) = 0.477 \times 0.523 \approx 0.249$$

$$\delta^{[2]} = \frac{\partial L}{\partial z^{[2]}} = \frac{\partial L}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial z^{[2]}} = -2.096 \times 0.249 \approx -0.522$$

**Hidden layer gradients:**

$$\frac{\partial L}{\partial W^{[2]}} = \delta^{[2]}a^{[1]T} = -0.522 \times \begin{bmatrix} 0.401 \\ 0.646 \end{bmatrix}^T = \begin{bmatrix} -0.209 & -0.337 \end{bmatrix}$$

$$\frac{\partial L}{\partial b^{[2]}} = \delta^{[2]} = -0.522$$

$$\delta^{[1]} = (W^{[2]T}\delta^{[2]}) \odot \sigma'(z^{[1]}) = \begin{bmatrix} 0.4 \\ -0.7 \end{bmatrix} \times -0.522 \odot \begin{bmatrix} 0.401 \times (1 - 0.401) \\ 0.646 \times (1 - 0.646) \end{bmatrix}$$

$$= \begin{bmatrix} -0.209 \\ 0.365 \end{bmatrix} \odot \begin{bmatrix} 0.240 \\ 0.229 \end{bmatrix} = \begin{bmatrix} -0.050 \\ 0.084 \end{bmatrix}$$

**Step 4: Weight Update**

**Learning rate:** $\eta = 0.1$
**Update output layer:**

$$W_{\text{new}}^{[2]} = W^{[2]} - \eta\frac{\partial L}{\partial W^{[2]}} = \begin{bmatrix} 0.4 & -0.7 \end{bmatrix} - 0.1 \times \begin{bmatrix} -0.209 & -0.337 \end{bmatrix} = \begin{bmatrix} 0.421 & -0.666 \end{bmatrix}$$

$$b_{\text{new}}^{[2]} = b^{[2]} - \eta\frac{\partial L}{\partial b^{[2]}} = 0.2 - 0.1 \times (-0.522) = 0.252$$

**Update hidden layer:**

$$\frac{\partial L}{\partial W^{[1]}} = \delta^{[1]}x^T = \begin{bmatrix} -0.050 \\ 0.084 \end{bmatrix}\begin{bmatrix} 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & -0.050 \\ 0 & 0.084 \end{bmatrix}$$

$$W_{\text{new}}^{[1]} = W^{[1]} - \eta\frac{\partial L}{\partial W^{[1]}} = \begin{bmatrix} 0.8 & -0.5 \\ 0.1 & 0.9 \end{bmatrix} - 0.1 \times \begin{bmatrix} 0 & -0.050 \\ 0 & 0.084 \end{bmatrix} = \begin{bmatrix} 0.8 & -0.495 \\ 0.1 & 0.892 \end{bmatrix}$$

**Why weights update in this direction:** The gradients point in the direction of steepest ascent of the loss function, so we move in the opposite direction (negative gradient) to minimize loss.

## 1.6   Particle Swarm Optimization: Detailed Example

### 1.6.1   Problem Setup

**Minimize $f(x) = x^2$:** Find the minimum of a simple quadratic function.

### 1.6.2   Step-by-Step PSO with Reasoning

**Initialization:**

- Particle 1: $x_1 = 3$, $v_1 = 0$, $p_{\text{best},1} = 3$

- Particle 2: $x_2 = -2$, $v_2 = 0$, $p_{\text{best},2} = -2$

- Global best: $g_{\text{best}} = -2$ (since $f(-2) = 4 < f(3) = 9$)

- Parameters: $w = 0.7$, $c_1 = c_2 = 1.5$, $r_1 = 0.6$, $r_2 = 0.8$

**Why these parameters:** - $w = 0.7$: Good balance between exploration and exploitation - $c_1 = c_2 = 1.5$: Equal importance to personal and social learning - Random numbers $r_1, r_2$: Introduce stochasticity for better exploration

**Iteration 1:**
**Particle 1 update:**

$$\begin{aligned} v_1^{new} &= wv_1 + c_1r_1(p_{\text{best},1} - x_1) + c_2r_2(g_{\text{best}} - x_1) \\ &= 0.7 \times 0 + 1.5 \times 0.6 \times (3 - 3) + 1.5 \times 0.8 \times (-2 - 3) \\ &= 0 + 0 + 1.5 \times 0.8 \times (-5) = -6 \\ x_1^{new} &= x_1 + v_1^{new} = 3 + (-6) = -3 \end{aligned}$$

**Why velocity is negative:** Particle 1 is at position 3, but the global best is at -2, so it needs to move left (negative direction).

**Particle 2 update:**

$$v_2^{new} = 0.7 \times 0 + 1.5 \times 0.6 \times (-2 - (-2)) + 1.5 \times 0.8 \times (-2 - (-2))$$
$$= 0 + 0 + 0 = 0$$
$$x_2^{new} = -2 + 0 = -2$$

**Why velocity is zero:** Particle 2 is already at its personal best and the global best, so no movement needed.

**Update personal bests:** - Particle 1: $f(-3) = 9 > f(3) = 9$, so $p_{\text{best},1}$ stays at 3 - Particle 2: $p_{\text{best},2}$ stays at -2

**Update global best:** Still $g_{\text{best}} = -2$

**Iteration 2:**

**Particle 1 update:**

$$v_1^{new} = 0.7 \times (-6) + 1.5 \times 0.6 \times (3 - (-3)) + 1.5 \times 0.8 \times (-2 - (-3))$$
$$= -4.2 + 1.5 \times 0.6 \times 6 + 1.5 \times 0.8 \times 1$$
$$= -4.2 + 5.4 + 1.2 = 2.4$$
$$x_1^{new} = -3 + 2.4 = -0.6$$

**Why velocity changed direction:** Now the cognitive component (attraction to personal best at 3) dominates, pulling particle right.

**Update personal bests:** - Particle 1: $f(-0.6) = 0.36 < f(3) = 9$, so $p_{\text{best},1} = -0.6$ - Global best: $f(-0.6) = 0.36 < f(-2) = 4$, so $g_{\text{best}} = -0.6$

**Convergence:** Over more iterations, particles converge toward the true minimum at x=0.

# 2 Why These Mathematical Operations Work: Fundamental Principles

## 2.1 The Optimization Perspective

**Gradient Descent Principle:**

$$\theta_{t+1} = \theta_t - \eta \nabla_\theta L(\theta_t)$$

**Why this works:** The gradient points in the direction of steepest ascent. By moving in the opposite direction, we guarantee (for small enough $\eta$) that the loss decreases.

**Example:** For $f(x) = x^2$, $\nabla f(x) = 2x$. At $x = 3$, gradient is 6. Moving opposite: $x_{new} = 3 - 0.1 \times 6 = 2.4$, and indeed $f(2.4) = 5.76 < f(3) = 9$.

## 2.2 The Probability Perspective

**Bayes' Theorem:**

$$P(\theta|X) = \frac{P(X|\theta)P(\theta)}{P(X)}$$

**Why this works:** It combines prior knowledge $P(\theta)$ with evidence from data $P(X|\theta)$ to form updated beliefs $P(\theta|X)$.

**Example:** In spam classification, prior $P(\text{spam})$ combined with word probabilities $P("\text{free}"|\text{spam})$ gives better prediction than either alone.

## 2.3    The Information Theory Perspective

**Information Gain:**
$$IG(Y, X) = H(Y) - H(Y|X)$$

**Why this works:** It measures how much a feature reduces uncertainty about the target. Features with high information gain are most useful for prediction.

**Example:** In our decision tree example, "Income" gave perfect information gain (0.5), meaning it completely determined the outcome.

## 2.4    The Geometric Perspective

**Margin Maximization:**
$$\text{Margin} = \frac{2}{\|w\|}$$

**Why this works:** Larger margins mean the decision boundary is farther from the data points, making the classifier more robust to noise and better at generalization.

**Example:** In our SVM example, we explicitly maximized the margin by minimizing $\|w\|$.

# 3    Conclusion

## 3.1    Key Mathematical Insights

- **Linear Algebra** provides the foundation for handling multi-dimensional data efficiently

- **Calculus** enables optimization through gradient-based learning

- **Probability Theory** supports uncertainty quantification and Bayesian methods

- **Information Theory** guides feature selection and model compression

- **Optimization Theory** ensures we find good solutions to complex problems

## 3.2    The Importance of Understanding Operations

Understanding the "why" and "how" behind mathematical operations:

- **Enables debugging** when models don't perform as expected

- **Facilitates innovation** by understanding what can be modified

- **Supports algorithm selection** based on problem characteristics

- **Provides intuition** for hyperparameter tuning

- **Builds foundation** for developing new algorithms

Each mathematical operation serves a specific purpose in the learning process, and understanding this relationship is the key to mastering machine learning.