

# *The Mathematical Trinity of Machine Learning*

# Gradient, Regularization, Eigenvalues & Eigenvectors

Data Science Expert

May 11, 2026

## Contents

<b>1</b>	<b>Part 1: The Gradient — The Compass of Machine Learning</b>	<b>4</b>
1.1	What is a Gradient? . . . . .	4
1.2	Mathematical Definition . . . . .	4
1.3	What the Gradient Tells Us . . . . .	4
1.4	The Three Critical Questions Gradient Answers . . . . .	4
1.4.1	Question 1: Which direction should I move? . . . . .	4
1.4.2	Question 2: How fast should I move? (Step size / Learning Rate) . . . . .	4
1.4.3	Question 3: Which variable matters most? . . . . .	5
1.5	What Happens WITHOUT Gradient? . . . . .	5
<b>2</b>	<b>Part 2: Regularization — The Guardrails of Learning</b>	<b>5</b>
2.1	The Big Picture Connection . . . . .	5
2.2	Mathematical Connection . . . . .	5
2.3	L1 vs L2 Regularization: The Sparsity Question . . . . .	6
2.3.1	Important Correction! . . . . .	6
2.3.2	Why L1 Pushes Weights Exactly to Zero . . . . .	6
2.3.3	Why L2 Only Shrinks (Never to Zero) . . . . .	6
2.3.4	Numerical Comparison Example . . . . .	7
<b>3</b>	<b>Part 3: Eigenvalues and Eigenvectors — The Hidden Structure</b>	<b>7</b>
3.1	What Are They? (Intuitive Explanation) . . . . .	7
3.2	Mathematical Definition . . . . .	7
3.3	Geometric Interpretation . . . . .	7
3.4	How to Compute Eigenvalues . . . . .	7
3.5	Example: $2 \times 2$ Matrix . . . . .	8
3.6	Why Eigenvalues and Eigenvectors Matter in Data Science . . . . .	8
3.6.1	1. Principal Component Analysis (PCA) . . . . .	8
3.6.2	2. Hessian Matrix and Optimization . . . . .	8
3.6.3	3. Convergence Rate of Gradient Descent . . . . .	8
3.6.4	4. PageRank Algorithm (Google) . . . . .	8
3.6.5	5. Spectral Clustering . . . . .	9
3.6.6	6. Stability of Dynamical Systems . . . . .	9
3.7	The Connection: Eigenvalues, Gradient, and Regularization . . . . .	9

3.7.1	How They Connect . . . . .	9
3.7.2	Regularization Improves Condition Number . . . . .	9
3.7.3	Why This Matters for Gradient Descent . . . . .	9
<b>4</b>	<b>Part 4: Complete Comparative Summary</b>	<b>9</b>
4.1	All Concepts at a Glance . . . . .	9
4.2	The Grand Unified Picture . . . . .	10
<b>5</b>	<b>Part 5: Practical Implications and Deep Concluding Thoughts</b>	<b>10</b>
5.1	What Happens Without Each Concept? . . . . .	10
5.2	Final Deep Concluding Thought . . . . .	10
5.3	Key Takeaways for Practice . . . . .	10
<b>A</b>	<b>Formula Reference Sheet</b>	<b>11</b>

# 1 Part 1: The Gradient — The Compass of Machine Learning

## 1.1 What is a Gradient?

Think of a gradient as a **direction + strength signal** that tells you:

*”If you want to improve something (reduce loss), move this way, and this much.”*

## 1.2 Mathematical Definition

For a function  $f(\mathbf{w})$  with multiple variables (weights), the gradient is a vector of partial derivatives:

$$\nabla f(\mathbf{w}) = \left[ \frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \dots, \frac{\partial f}{\partial w_n} \right]^T$$

## 1.3 What the Gradient Tells Us

## 1.4 The Three Critical Questions Gradient Answers

### 1.4.1 Question 1: Which direction should I move?

**Why it matters:** Without direction, you’re just wandering randomly.

In high dimensions (millions of parameters), random movement is useless. The gradient gives the **correct direction** for fastest improvement.

**Update rule (Gradient Descent):**

$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} - \eta \nabla L(\mathbf{w}_{\text{old}})$$

where  $\eta$  is the learning rate.

### 1.4.2 Question 2: How fast should I move? (Step size / Learning Rate)

#### The Learning Rate Trade-off

- **Too large ( $\eta$  big):** Overshoots minimum  $\rightarrow$  training unstable, diverges
- **Too small ( $\eta$  tiny):** Crawls slowly  $\rightarrow$  takes forever, may get stuck
- **Just right:** Stable convergence  $\rightarrow$  optimal efficiency

**Mathematical bound for stability:**

$$\text{If } \eta > \frac{2}{\lambda_{\max}}, \text{ gradient descent diverges}$$

where  $\lambda_{\max}$  is the largest eigenvalue of the Hessian matrix (more on eigenvalues later).

### 1.4.3 Question 3: Which variable matters most?

The gradient is a vector; each component  $\frac{\partial L}{\partial w_j}$  tells how sensitive the loss is to that specific weight.

**Why this matters:**

- Important features get large gradients  $\rightarrow$  large updates
- Irrelevant features get small gradients  $\rightarrow$  tiny updates
- This is how the model learns which weights are crucial

## 1.5 What Happens WITHOUT Gradient?

### Without Gradient

- **No direction:** You would guess random weight updates
- **No speed:** No way to know step size
- **No importance:** All weights updated equally
- **Result:** Training never converges; random search in high dimensions fails

## 2 Part 2: Regularization — The Guardrails of Learning

### 2.1 The Big Picture Connection

- **Gradient:** Tells you how to update parameters to reduce loss
- **Regularization:** Modifies the loss so the model doesn't overfit
- **Connection:** When you change the loss, you automatically change its gradient

### 2.2 Mathematical Connection

We don't minimize just the data loss anymore. We add a penalty:

$$L_{\text{total}}(\mathbf{w}) = L_{\text{data}}(\mathbf{w}) + \lambda R(\mathbf{w})$$

Taking the gradient:

$$\nabla L_{\text{total}}(\mathbf{w}) = \nabla L_{\text{data}}(\mathbf{w}) + \lambda \nabla R(\mathbf{w})$$

**The update rule becomes:**

$$\mathbf{w}_{\text{new}} = \mathbf{w}_{\text{old}} - \eta (\nabla L_{\text{data}}(\mathbf{w}_{\text{old}}) + \lambda \nabla R(\mathbf{w}_{\text{old}}))$$

## 2.3 L1 vs L2 Regularization: The Sparsity Question

### 2.3.1 Important Correction!

#### Correct Statement

- **L1 (Lasso):** Creates sparsity — pushes weights **exactly to zero**
- **L2 (Ridge):** Shrinks weights — makes them **small but non-zero**
- The statement "L2 increases sparsity" is **FALSE**

### 2.3.2 Why L1 Pushes Weights Exactly to Zero

**Mathematical derivation:** L1 penalty:  $R(w) = |w|$

The subgradient (for  $w \neq 0$ ):

$$\frac{\partial R}{\partial w} = \text{sign}(w) = \begin{cases} +1 & \text{if } w > 0 \\ -1 & \text{if } w < 0 \end{cases}$$

**The update rule:**

$$w_{\text{new}} = w_{\text{old}} - \eta \left( \frac{\partial L_{\text{data}}}{\partial w} + \lambda \cdot \text{sign}(w) \right)$$

**Why weights become exactly zero:**

- The regularization gradient  $\lambda \cdot \text{sign}(w)$  is **constant** ( $+\lambda$  or  $-\lambda$ )
- It keeps pushing with the same force even when  $w$  is tiny
- Once  $w$  crosses zero, it can stay at zero due to subgradient

### 2.3.3 Why L2 Only Shrinks (Never to Zero)

**Mathematical derivation:** L2 penalty:  $R(w) = w^2$

The gradient:

$$\frac{\partial R}{\partial w} = 2w$$

**The update rule:**

$$w_{\text{new}} = w_{\text{old}} - \eta \left( \frac{\partial L_{\text{data}}}{\partial w} + 2\lambda w_{\text{old}} \right)$$

**Why weights never become exactly zero:**

- The regularization gradient  $2\lambda w$  is **proportional to the weight**
- As  $w$  gets smaller, the push gets weaker
- It asymptotically approaches zero but never reaches it

### 2.3.4 Numerical Comparison Example

#### Example: L1 vs L2 Update

Suppose  $w = 0.5$ ,  $\eta = 0.1$ ,  $\lambda = 0.1$ , and assume data gradient  $= 0$ .

**L1:**

$$w_{\text{new}} = 0.5 - 0.1 \times 0.1 \times (+1) = 0.5 - 0.01 = 0.49$$

Each step subtracts exactly 0.01. After 50 steps,  $w = 0$ .

**L2:**

$$w_{\text{new}} = 0.5 - 0.1 \times 0.1 \times (2 \times 0.5) = 0.5 - 0.01 = 0.49$$

Next step:  $w = 0.49$ ,  $2w = 0.98$ , subtract  $0.1 \times 0.1 \times 0.98 = 0.0098$ , so  $w = 0.4802$ .  
Step size decreases as  $w$  decreases  $\rightarrow$  never reaches zero.

## 3 Part 3: Eigenvalues and Eigenvectors — The Hidden Structure

### 3.1 What Are They? (Intuitive Explanation)

Imagine you have a rubber sheet (a transformation). Some directions stretch or compress the sheet. An **eigenvector** is a special direction that keeps its orientation after transformation. The **eigenvalue** tells you how much it gets stretched (or compressed).

*"Eigenvectors are the natural axes of a transformation. Eigenvalues tell you the scaling factor along those axes."*

### 3.2 Mathematical Definition

For a square matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , if:

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

then:

- $\mathbf{v}$  is an **eigenvector** (non-zero vector)
- $\lambda$  is the corresponding **eigenvalue** (scalar)

### 3.3 Geometric Interpretation

### 3.4 How to Compute Eigenvalues

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

$$\mathbf{A}\mathbf{v} - \lambda\mathbf{v} = 0 \implies (\mathbf{A} - \lambda\mathbf{I})\mathbf{v} = 0$$

For non-zero  $\mathbf{v}$ , the matrix  $(\mathbf{A} - \lambda\mathbf{I})$  must be singular:

$$\det(\mathbf{A} - \lambda\mathbf{I}) = 0$$

This gives the **characteristic polynomial** of degree  $n$ .

### 3.5 Example: 2×2 Matrix

Let  $\mathbf{A} = \begin{bmatrix} 4 & 1 \\ 2 & 3 \end{bmatrix}$ .

Characteristic equation:

$$\det \begin{bmatrix} 4 - \lambda & 1 \\ 2 & 3 - \lambda \end{bmatrix} = (4 - \lambda)(3 - \lambda) - 2 = 0$$

$$12 - 4\lambda - 3\lambda + \lambda^2 - 2 = 0 \implies \lambda^2 - 7\lambda + 10 = 0$$

$$(\lambda - 2)(\lambda - 5) = 0 \implies \lambda = 2 \text{ or } \lambda = 5$$

**Eigenvectors:** For  $\lambda = 5$ :  $\begin{bmatrix} 4 - 5 & 1 \\ 2 & 3 - 5 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} -1 & 1 \\ 2 & -2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 0$

$$-v_1 + v_2 = 0 \implies v_1 = v_2 \implies \mathbf{v} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

### 3.6 Why Eigenvalues and Eigenvectors Matter in Data Science

#### 3.6.1 1. Principal Component Analysis (PCA)

PCA finds the eigenvectors of the covariance matrix. The eigenvector with the largest eigenvalue points in the direction of maximum variance.

$$\text{Covariance matrix } \Sigma = \frac{1}{n} \mathbf{X}^T \mathbf{X}$$

$$\Sigma \mathbf{v}_i = \lambda_i \mathbf{v}_i$$

where  $\lambda_i$  = variance explained by component  $i$ .

#### 3.6.2 2. Hessian Matrix and Optimization

The Hessian  $\mathbf{H}$  (matrix of second derivatives) tells us about curvature:

- All eigenvalues  $> 0 \rightarrow$  local minimum
- All eigenvalues  $< 0 \rightarrow$  local maximum
- Mixed signs  $\rightarrow$  saddle point

#### 3.6.3 3. Convergence Rate of Gradient Descent

For a quadratic objective, gradient descent converges at a rate determined by the **condition number**:

$$\kappa = \frac{\lambda_{\max}}{\lambda_{\min}}$$

Larger  $\kappa$  means slower convergence.

#### 3.6.4 4. PageRank Algorithm (Google)

PageRank computes the dominant eigenvector of the web's link matrix:

$$\mathbf{r} = \alpha \mathbf{M} \mathbf{r} + (1 - \alpha) \mathbf{1}$$

The eigenvector gives page importance scores.

### 3.6.5 5. Spectral Clustering

Uses eigenvectors of the graph Laplacian to cluster data points.

### 3.6.6 6. Stability of Dynamical Systems

In recurrent neural networks (RNNs), eigenvalues of the weight matrix determine stability:

- $|\lambda| < 1 \rightarrow$  stable (gradients vanish)
- $|\lambda| > 1 \rightarrow$  unstable (gradients explode)

## 3.7 The Connection: Eigenvalues, Gradient, and Regularization

### 3.7.1 How They Connect

#### The Trinity Connection

1. **Gradient** gives the direction of movement
2. **Hessian** (matrix of second derivatives) has eigenvalues that tell us about curvature
3. **Regularization** modifies the Hessian's eigenvalues to improve conditioning

### 3.7.2 Regularization Improves Condition Number

For linear regression, the Hessian of the loss is  $\mathbf{X}^T \mathbf{X}$  (without regularization) or  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$  (with L2 regularization).

**Without regularization:** The smallest eigenvalue might be tiny  $\rightarrow$  ill-conditioned  $\rightarrow$  slow convergence.

**With L2 regularization:** Add  $\lambda$  to every eigenvalue:

Eigenvalues become  $\lambda_i + \lambda$

This makes the smallest eigenvalue larger, improving the condition number:

$$\kappa_{\text{reg}} = \frac{\lambda_{\max} + \lambda}{\lambda_{\min} + \lambda} < \frac{\lambda_{\max}}{\lambda_{\min}} = \kappa$$

### 3.7.3 Why This Matters for Gradient Descent

Better condition number = faster convergence. This is one reason why regularization helps beyond just preventing overfitting — it also **speeds up optimization!**

## 4 Part 4: Complete Comparative Summary

### 4.1 All Concepts at a Glance

---

Concept	Mathematical Form	What It Tells Us	Why Important
Gradient	$\nabla f = [\partial f / \partial x_i]$	Direction and rate of steepest ascent	Guides optimization
L1 Regularization	$\lambda \ w\ _1$	Constant pull toward zero	Creates sparsity
L2 Regularization	$\lambda \ w\ _2^2$	Pull proportional to $w$	Shrinks weights
Eigenvalue	$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$	Scaling along invariant direction	Compression, stability
Eigenvector	Same equation	Direction that doesn't change orientation	Natural axes of system
Hessian	$\mathbf{H}_{ij} = \partial^2 f / \partial x_i \partial x_j$	Curvature matrix	Convexity, local minima
Condition Number	$\kappa = \lambda_{\max} / \lambda_{\min}$	How ill-conditioned the problem is	Convergence speed

## 4.2 The Grand Unified Picture

### The Data Science Trinity

**Gradient** → tells you **where to go**  
**Regularization** → tells you **how to stay safe**  
**Eigenvalues** → tell you **the landscape's shape**

Without gradient: you're lost (random guessing).

Without regularization: you're reckless (overfitting).

Without eigenvalues: you're blind to curvature (can't diagnose convergence).

## 5 Part 5: Practical Implications and Deep Concluding Thoughts

### 5.1 What Happens Without Each Concept?

### 5.2 Final Deep Concluding Thought

*"Gradient gives you the direction. Regularization gives you restraint. Eigenvalues reveal the geometry. Together, they form the mathematical foundation of learning from data."*

### 5.3 Key Takeaways for Practice

1. Always compute gradients (automatic in PyTorch/TensorFlow)
2. Always add some regularization — even small  $\lambda$  helps

3. **Monitor eigenvalues** when possible — for small problems, check Hessian eigenvalues to verify convexity
4. **Use the condition number as a diagnostic** — if  $\kappa$  is huge, consider better scaling or more regularization
5. **Remember L1 vs L2:** L1 for sparsity (feature selection), L2 for shrinkage (all features contribute)
6. **The gradient connects everything** — regularization modifies it; eigenvalues describe the landscape it moves through

## A Formula Reference Sheet

- Gradient:  $\nabla f = \left[ \frac{\partial f}{\partial x_i} \right]$
- Gradient descent:  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla L$
- L1:  $R(w) = \lambda|w|$ , gradient =  $\lambda \cdot \text{sign}(w)$
- L2:  $R(w) = \lambda w^2$ , gradient =  $2\lambda w$
- Eigenvalue equation:  $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$
- Characteristic polynomial:  $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$
- Hessian:  $\mathbf{H}_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$
- Condition number:  $\kappa = \frac{\lambda_{\max}}{\lambda_{\min}}$
- Regularized Hessian:  $\mathbf{H}_{\text{reg}} = \mathbf{H} + \lambda\mathbf{I}$
- PCA:  $\Sigma\mathbf{v} = \lambda\mathbf{v}$
- PageRank:  $\mathbf{r} = \alpha\mathbf{M}\mathbf{r} + (1 - \alpha)\mathbf{1}$