

# Calculus Fundamentals for Data Science & Machine Learning

## Recap of Derivatives

### What is a Derivative?

A **derivative** measures the rate at which a function changes. In practical terms, it answers the question: "If I change the input by a tiny amount, how much does the output change?"

**Geometric interpretation:** The derivative is the slope of the tangent line to a curve at a specific point. If you zoom in very close to a point on a smooth curve, it looks like a straight line—that line's slope is the derivative.

**Notation:** The derivative of  $f(x)$  can be written as:

- $f'(x)$
- $df/dx$
- $dy/dx$  (if  $y = f(x)$ )

### The Limit Definition

The formal definition of a derivative is:

$$f'(x) = \lim_{h \rightarrow 0} [f(x+h) - f(x)] / h$$

This says: take two points on the curve very close together, compute the slope between them, and see what happens as the distance between points shrinks to zero.

### How to Compute Derivatives: Common Rules

Rather than using the limit definition every time, we use shortcuts:

**Power Rule:** If  $f(x) = x^n$ , then  $f'(x) = n \cdot x^{n-1}$

- Example:  $f(x) = x^3 \rightarrow f'(x) = 3x^2$

**Constant Multiple Rule:** If  $f(x) = c \cdot g(x)$ , then  $f'(x) = c \cdot g'(x)$

- Example:  $f(x) = 5x^2 \rightarrow f'(x) = 5 \cdot 2x = 10x$

**Sum Rule:** If  $f(x) = g(x) + h(x)$ , then  $f'(x) = g'(x) + h'(x)$

- Example:  $f(x) = x^3 + 2x \rightarrow f'(x) = 3x^2 + 2$

**Product Rule:** If  $f(x) = g(x) \cdot h(x)$ , then  $f'(x) = g'(x) \cdot h(x) + g(x) \cdot h'(x)$

- Example:  $f(x) = x^2 \cdot \sin(x) \rightarrow f'(x) = 2x \cdot \sin(x) + x^2 \cdot \cos(x)$

**Chain Rule:** If  $f(x) = g(h(x))$ , then  $f'(x) = g'(h(x)) \cdot h'(x)$

- Example:  $f(x) = (x^2 + 1)^3 \rightarrow f'(x) = 3(x^2 + 1)^2 \cdot 2x$

**Exponential:**  $d/dx[e^x] = e^x$

**Logarithm:**  $d/dx[\ln(x)] = 1/x$

**Sine/Cosine:**  $d/dx[\sin(x)] = \cos(x)$  and  $d/dx[\cos(x)] = -\sin(x)$

## What Derivatives Tell Us

A derivative value tells us about the behavior of the function:

- $f'(x) > 0$ : function is increasing at that point (slope is positive)
  - $f'(x) < 0$ : function is decreasing at that point (slope is negative)
  - $f'(x) = 0$ : function has a flat spot (potential minimum, maximum, or saddle point)
- 

## Recap of Integrals

### What is an Integral?

An **integral** is the inverse operation of a derivative. While derivatives break down change into infinitesimal pieces, integrals reassemble pieces into a whole.

**Geometric interpretation:** The definite integral represents the **area under the curve** between two points. If the curve is above the x-axis, the area is positive; if below, it's negative (contributing negative area).

**Notation:** The integral of  $f(x)$  from  $a$  to  $b$  is written as:

$$\int [a \text{ to } b] f(x) dx$$

Read as: "the integral from  $a$  to  $b$  of  $f(x)$  with respect to  $x$ "

## The Fundamental Theorem of Calculus

This theorem is the bridge connecting derivatives and integrals:

**If  $F'(x) = f(x)$ , then:**

$$\int [a \text{ to } b] f(x) dx = F(b) - F(a)$$

This means: to compute a definite integral, find an antiderivative  $F$  (a function whose derivative is  $f$ ), then evaluate it at the upper and lower bounds and subtract.

**Example:** Compute  $\int[0 \text{ to } 3] 2x \, dx$

- The antiderivative of  $2x$  is  $x^2$  (since  $d/dx[x^2] = 2x$ )
- So:  $\int[0 \text{ to } 3] 2x \, dx = [x^2] \text{ from } 0 \text{ to } 3 = 3^2 - 0^2 = 9$

## Common Antiderivatives

- $\int x^n \, dx = x^{(n+1)/(n+1)} + C$  (for  $n \neq -1$ )
- $\int e^x \, dx = e^x + C$
- $\int 1/x \, dx = \ln(x) + C$
- $\int \sin(x) \, dx = -\cos(x) + C$
- $\int \cos(x) \, dx = \sin(x) + C$
- $\int 1/(1+x^2) \, dx = \arctan(x) + C$

(The "+  $C$ " represents the constant of integration—when we compute definite integrals, it cancels out.)

## Why "Area Under the Curve"?

Imagine slicing the area under a curve into infinitely many thin vertical rectangles. Each rectangle has width  $dx$  (infinitesimally small) and height  $f(x)$ . The area of one rectangle is  $f(x) \cdot dx$ . The integral sums all these tiny areas:

$$\int[a \text{ to } b] f(x) \, dx = \text{sum of all } f(x) \cdot dx \text{ rectangles from } a \text{ to } b$$

## Applications of Calculus in Machine Learning (Overview)

### 1. Optimization and Parameter Tuning

Machine learning models have parameters (weights, biases) that need to be adjusted to make accurate predictions. Calculus enables optimization through derivatives.

**The Big Picture:** We define a loss function  $L$  that measures how wrong our model is. We want to find parameter values that minimize this loss. Derivatives tell us which direction to adjust parameters to reduce error. We repeatedly compute gradients (derivatives) and move parameters in the direction that decreases loss—this is how neural networks learn.

**Why it matters:** Without calculus, we'd have no systematic way to improve model parameters. Gradient-based optimization is what makes modern deep learning possible.

## 2. Understanding Model Behavior

Derivatives help us understand how sensitive a model is to changes in input or parameters.

**Sensitivity analysis:** If the derivative of output with respect to input is large, small changes in input cause large changes in output (the model is very sensitive). If the derivative is small, the model is stable to input variations.

**Example context:** In a classification model, knowing the gradient of the decision boundary tells us how easily the prediction can change—useful for understanding robustness to adversarial examples.

## 3. Probability and Statistical Modeling

Calculus is fundamental to probability theory, which underlies all statistical machine learning.

**Probability density functions (PDFs):** PDFs must sum (integrate) to 1 over their entire domain. We use integrals to ensure probabilities are valid and to compute the probability of events.

**Expected values:** The expected value (mean) of a random variable is computed as an integral:  $E[X] = \int x \cdot f(x) dx$ . This tells us the average outcome we expect.

**Maximum likelihood estimation:** When fitting models to data, we often maximize the likelihood of observing the data given model parameters. This optimization uses derivatives to find the best fit.

**Example context:** In logistic regression, we find parameters that maximize the probability of correct classifications. This involves computing derivatives of the log-likelihood function.

## 4. Loss Functions and Their Properties

Different ML tasks use different loss functions, and understanding their derivatives reveals their behavior.

**Smooth vs sharp loss functions:** A loss function with derivatives that change gradually is easier to optimize (we can use small steps reliably). A loss with sharp changes is harder to optimize.

**Convexity:** Some loss functions are convex (single global minimum), making optimization guaranteed to find the best solution. Others are non-convex (multiple local minima), requiring careful algorithm design. Understanding convexity comes from analyzing second derivatives (curvature).

## 5. Backpropagation in Neural Networks

The chain rule of calculus is the foundation of backpropagation, the algorithm for training neural networks.

**What happens:** A neural network is a chain of functions: inputs → layer 1 → layer 2 → ... → output → loss. To improve the network, we need gradients of loss with respect to each weight. The chain rule lets us compute these efficiently by working backward through the network, multiplying local gradients at each layer.

**Why it matters:** Without the chain rule, training deep networks would be computationally infeasible. The chain rule makes it possible to efficiently compute gradients for millions or billions of parameters.

## 6. Regularization and Constraints

Calculus helps enforce constraints on model complexity.

**Regularization:** We often add penalty terms to loss functions (like L2 regularization) to encourage simpler models. The derivatives of these penalty terms guide the optimization to find models that generalize better.

## 7. Uncertainty Quantification

In probabilistic ML (Bayesian methods), integrals appear when we need to compute probabilities by integrating over uncertain variables.

**Example context:** Computing posterior distributions often requires integrating over parameter space, which can be done numerically using Monte Carlo methods.

---

# Working with Functions of Multiple Variables

Most real ML problems don't involve single-variable functions. A neural network might have millions of parameters. Understanding how to work with multivariable functions is essential.

## Partial Derivatives

When a function depends on multiple variables, we can ask: "How does the output change if I vary just one input, keeping others fixed?" This is a **partial derivative**.

**Notation and computation:** For a function  $f(x, y, z)$ , we write:

- $\partial f / \partial x$ : partial derivative with respect to  $x$  (treat  $y$  and  $z$  as constants)
- $\partial f / \partial y$ : partial derivative with respect to  $y$  (treat  $x$  and  $z$  as constants)
- $\partial f / \partial z$ : partial derivative with respect to  $z$  (treat  $x$  and  $y$  as constants)

**Example:**  $f(x, y) = x^2 + 3xy + y^3$

$\partial f / \partial x = 2x + 3y$  (differentiate with respect to  $x$ , treat  $y$  as constant)

$\partial f / \partial y = 3x + 3y^2$  (differentiate with respect to  $y$ , treat  $x$  as constant)

We compute partial derivatives using the same rules as single-variable derivatives—just treat other variables as constants.

## The Gradient Vector

The **gradient** is a vector containing all partial derivatives:

$$\nabla f(x, y, z) = [\partial f / \partial x, \partial f / \partial y, \partial f / \partial z]$$

**Geometric meaning:** The gradient points in the direction of steepest ascent. At any point on a surface, if you move in the direction the gradient points, you go uphill as steeply as possible.

**In ML:** We move opposite the gradient (negative gradient) because we want to descend (minimize loss). The gradient tells us which parameter adjustments help most.

## Higher-Order Partial Derivatives

Just as we can take the second derivative of a single-variable function, we can take second partial derivatives of multivariable functions:

**Pure second partials:**  $\partial^2 f / \partial x^2$  or  $\partial^2 f / \partial y^2$  (take the same partial twice)

**Mixed partials:**  $\partial^2 f / \partial x \partial y$  (take partial with respect to x, then with respect to y)

**The Hessian matrix:** For a function  $f(x, y)$ , the Hessian is:

$$H = \begin{bmatrix} \partial^2 f / \partial x^2 & \partial^2 f / \partial x \partial y \\ \partial^2 f / \partial y \partial x & \partial^2 f / \partial y^2 \end{bmatrix}$$

The Hessian captures all curvature information. Its eigenvalues tell us if we're near a minimum (all positive eigenvalues), maximum (all negative), or saddle point (mixed signs).

**In ML:** The Hessian reveals the landscape around our current parameters. Where the Hessian is positive definite, we're near a local minimum and can trust our optimization is working well.

## Level Sets and Contours

A **level set** (or contour) is a curve where the function has a constant value.

**Example:** For  $f(x, y) = x^2 + y^2$ , the level sets are circles centered at the origin. Each circle represents points where  $x^2 + y^2$  equals some constant value.

**Key insight:** The gradient is always perpendicular to level sets. If you're standing on a hillside (the 3D surface of a function), contour lines show paths where elevation is constant. The direction of steepest climb is perpendicular to these contour lines—which is exactly where the gradient points.

**In ML visualization:** When optimizing a function of two variables, we often plot contours and show how gradient descent traces a path perpendicular to contours, spiraling toward the minimum.

## Directional Derivatives

Sometimes we want to know how fast the function changes in a specific direction (not necessarily the steepest direction).

The **directional derivative** in the direction of a unit vector  $\mathbf{u}$  is:

$$D_u f = \nabla f \cdot u$$

(The dot product of the gradient with the unit direction vector)

This is useful for understanding how the function behaves along any direction we choose.

### Multivariable Chain Rule

When we have a composite function of multiple variables, the chain rule generalizes:

If  $z = f(x, y)$  where  $x$  and  $y$  depend on  $t$  (time, iteration number, etc.):

$$dz/dt = (\partial f / \partial x) \cdot (dx/dt) + (\partial f / \partial y) \cdot (dy/dt)$$

We sum over all paths: each partial derivative times the rate of change in that variable.

**In neural networks:** When computing  $\partial \text{Loss} / \partial w$ , we apply this multivariable chain rule across all layers. Each layer contributes to the total gradient through this summation principle.

---

## Visualizing Derivatives and Integrals with Desmos

[Optional: Interactive visualizations can deepen understanding]

### Visualizing Derivatives

You can use **Desmos** ([desmos.com/calculator](https://desmos.com/calculator)) to visualize derivatives:

**Setup:** Graph a function like  $f(x) = x^2 - 2x + 1$

#### Visualize the derivative as slope:

1. Add a slider for point  $a$ :  $a = 0$  (or any value)
2. Plot the point  $(a, f(a))$
3. Plot the tangent line with slope  $f'(a)$ :  $y = f(a) \cdot (x - a) + f(a)$
4. Drag the slider and watch the tangent line pivot around the curve
5. The tangent line's slope at any point is the derivative at that point

#### Visualize the derivative function:

1. Graph  $f(x) = x^2 - 2x + 1$
2. Graph  $f(x) = 2x - 2$  (the derivative function)
3. Notice: where  $f'(x) = 0$ , the original function has a flat spot (minimum or maximum)

4. Where  $f(x) > 0$ ,  $f$  is increasing; where  $f(x) < 0$ ,  $f$  is decreasing

## Visualizing Integrals

**Setup:** Graph a function like  $f(x) = \sin(x)$

### Visualize the area under the curve:

1. Add sliders for bounds:  $a = 0, b = \pi$
2. Use Desmos's built-in tools or shade the region from  $x = a$  to  $x = b$
3. The shaded area is the definite integral  $\int[a \text{ to } b] \sin(x) dx$
4. Change bounds with sliders and watch the area change
5. Compute:  $\int[0 \text{ to } \pi] \sin(x) dx = [-\cos(x)] \text{ from } 0 \text{ to } \pi = 2$

### Visualize the antiderivative:

1. Graph  $f(x) = \sin(x)$
2. Graph  $F(x) = -\cos(x)$  (an antiderivative)
3. Note:  $F'(x) = \sin(x) = f(x)$
4. The antiderivative is the "reverse" operation of taking the derivative

## Visualizing Multivariable Functions

**3D surfaces:** Desmos doesn't natively display 3D surfaces, but you can use:

- **GeoGebra 3D calculator** ([geogebra.org](http://geogebra.org))
- **Wolfram Alpha** ([wolframalpha.com](http://wolframalpha.com))

**Setup:** Graph  $f(x, y) = x^2 + y^2$

### Understand the gradient:

1. Visualize level sets (contours) of the function
2. The gradient vector is perpendicular to these contours
3. The gradient points toward higher values (uphill)
4. Imagine water flowing downhill—it follows paths opposite the gradient

### Watch optimization happen:

1. Start at a random point on the surface
2. Compute the gradient

3. Move a small step opposite the gradient direction
  4. Repeat—you trace a path downhill toward the minimum
  5. This is gradient descent in action
- 

## Summary

**Derivatives** measure rates of change and reveal function behavior. They answer: "If I change the input, how much does the output change?"

**Integrals** accumulate quantities and compute areas. They're the inverse of derivatives and fundamental to probability theory.

**Multivariable calculus** extends these concepts to functions with many inputs. Partial derivatives and gradients guide optimization in high-dimensional spaces.

In **machine learning**, calculus isn't just theoretical—it's the engine. Derivatives power optimization algorithms that train models, integrals underpin probability theory, and understanding multivariable functions is essential for working with complex models.

**Visualization** through tools like Desmos builds intuition and transforms abstract mathematics into tangible understanding.