

Mechanistic Interpretability 101

Episode 3

Calculus



$$A = \pi r^2$$

$$C = 2\pi r$$

	30°	45°	60°
sin	$\frac{1}{2}$	$\frac{\sqrt{2}}{2}$	$\frac{\sqrt{3}}{2}$
cos	$\frac{\sqrt{3}}{2}$	$\frac{\sqrt{2}}{2}$	$\frac{1}{2}$
tan	$\frac{\sqrt{3}}{3}$	1	$\sqrt{3}$



$$\int \sin x dx = -\cos x + C$$

$$\int \frac{dx}{\cos^2 x} = \tan x + C$$

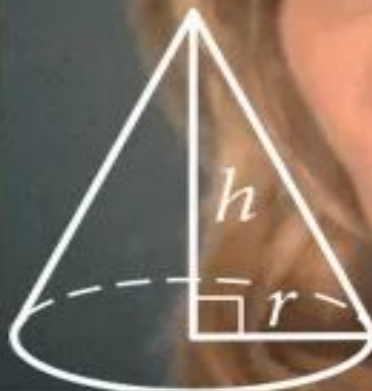
$$\int \tan x dx = -\ln|\cos x| + C$$

$$\int \frac{dx}{\sin x} = \ln\left|\tan \frac{x}{2}\right| + C$$

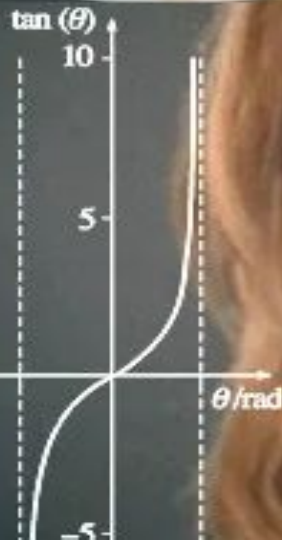
$$\int \frac{dx}{a^2 + x^2} = \frac{1}{a} \arctan \frac{x}{a} + C$$

$$\int \frac{dx}{x} = \ln|x| + C$$

$$V = \frac{1}{3} \pi r^2 h$$



$$V = \pi r^2 h$$



$$ax^2 + bx + c = 0$$

$$a\left(x^2 + \frac{b}{a}x + \frac{c}{a}\right) = 0$$

$$x^2 + 2\frac{b}{2a}x + \left(\frac{b}{2a}\right)^2 - \left(\frac{b}{2a}\right)^2 + \frac{c}{a} = 0$$

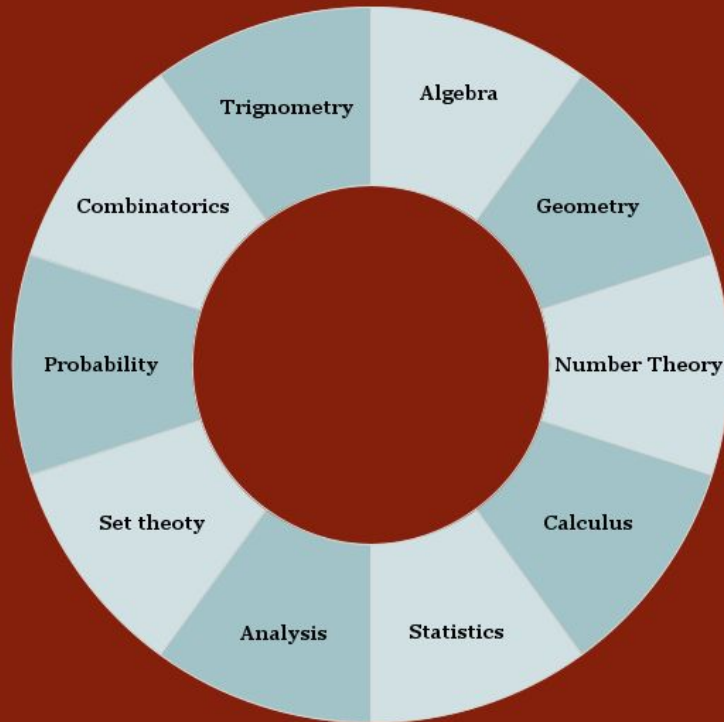
$$\left(x + \frac{b}{2a}\right)^2 - \frac{b^2 - 4ac}{4a^2} = 0$$



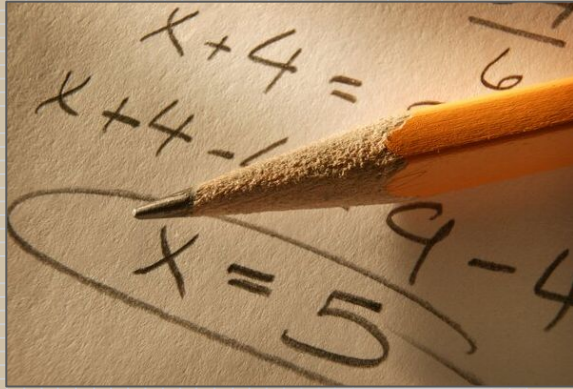
**CALCULUS, I SEE
YOU!**

AND I HAVE GOT YOU!

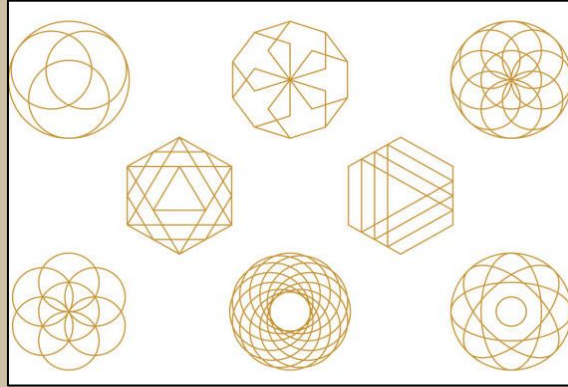
Fundamental branches of Math



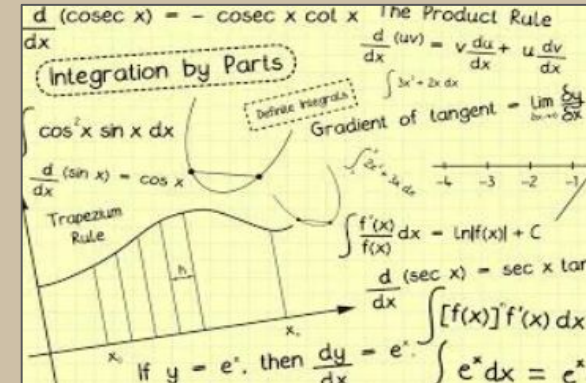
Fundamental branches of Math



Algebra



Geometry



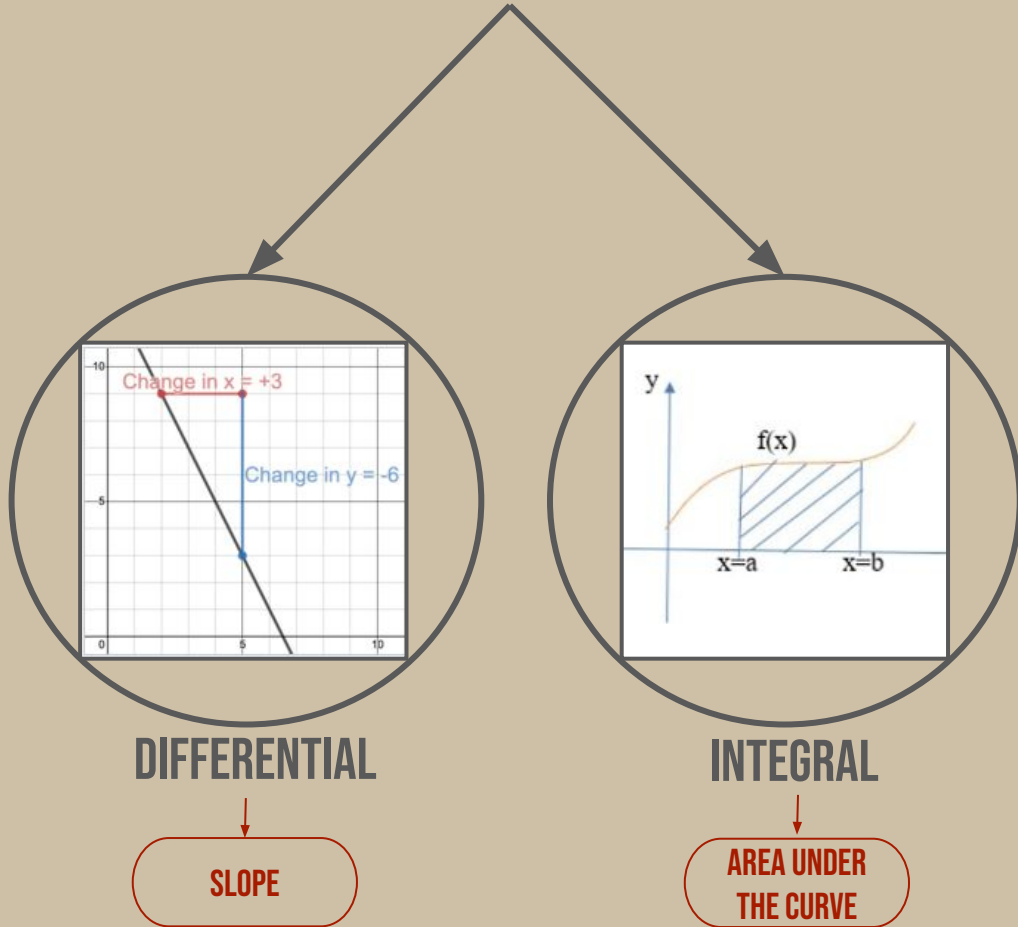
Calculus



Let us
'Change'
this
outlook
about
Change
today



Calculus





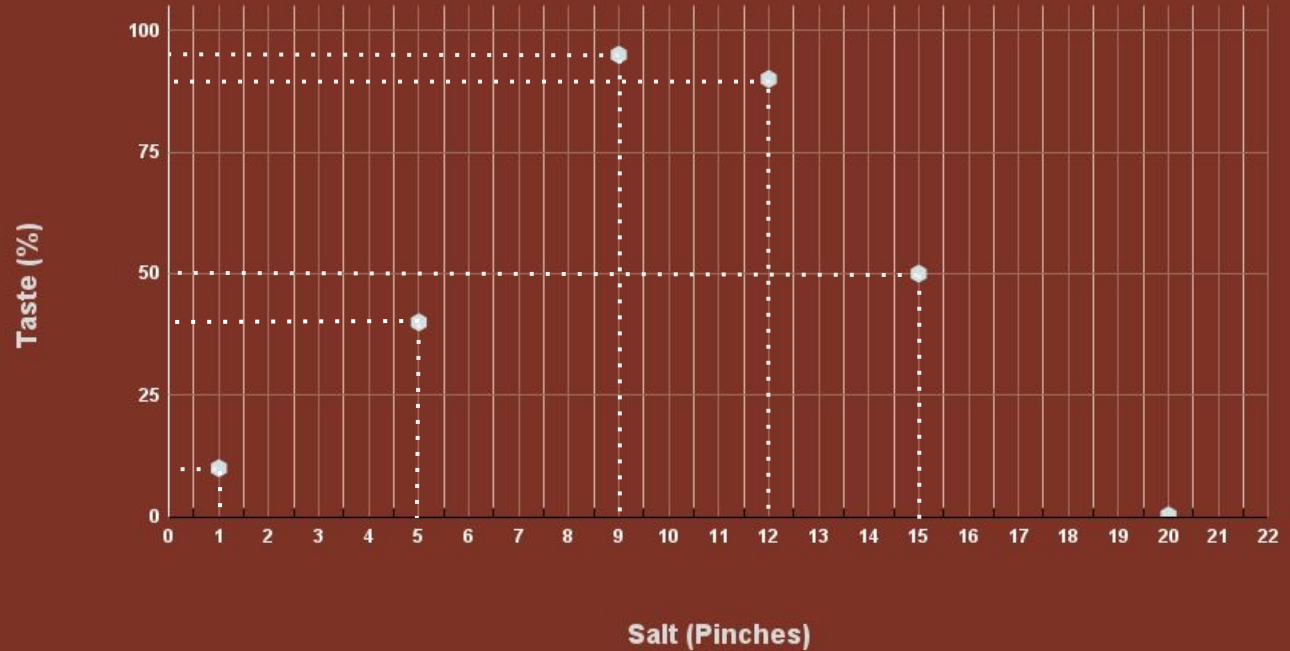
Understanding Functions

Salt (in Pinches)	Taste (in %)
1	10
5	40
9	95
12	90
15	50
20	0.3

Taste as a function of Salt

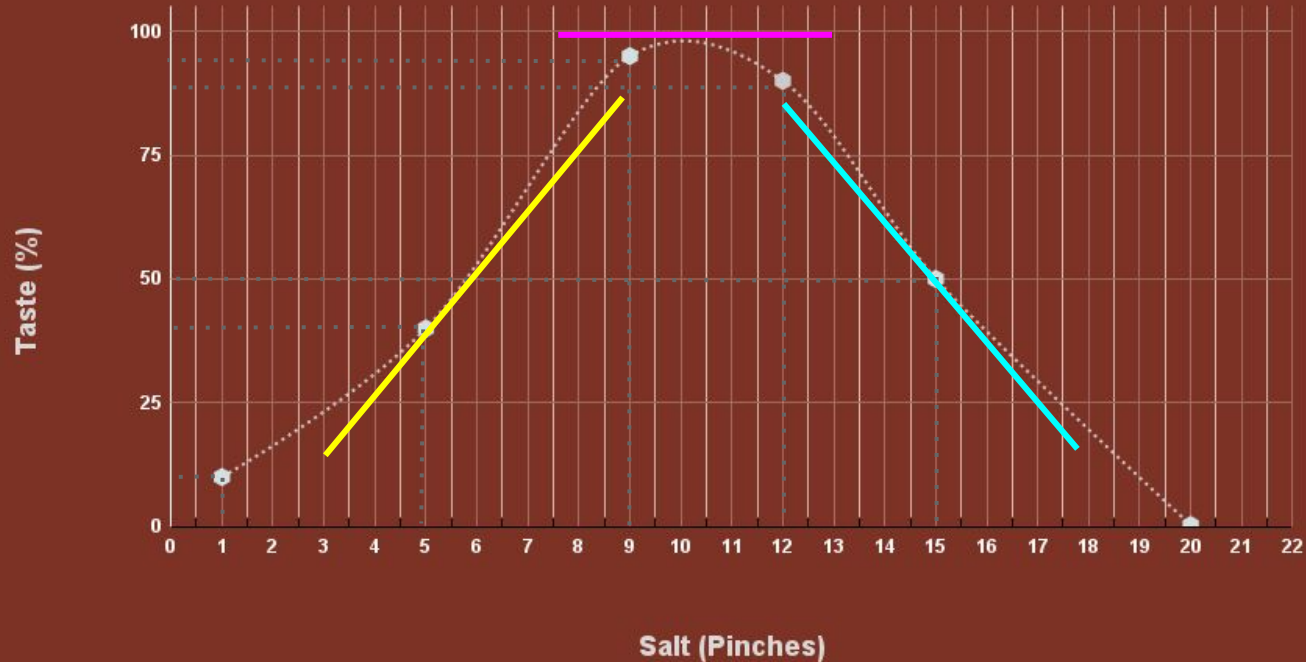


$$\text{taste} = f(\text{salt})$$



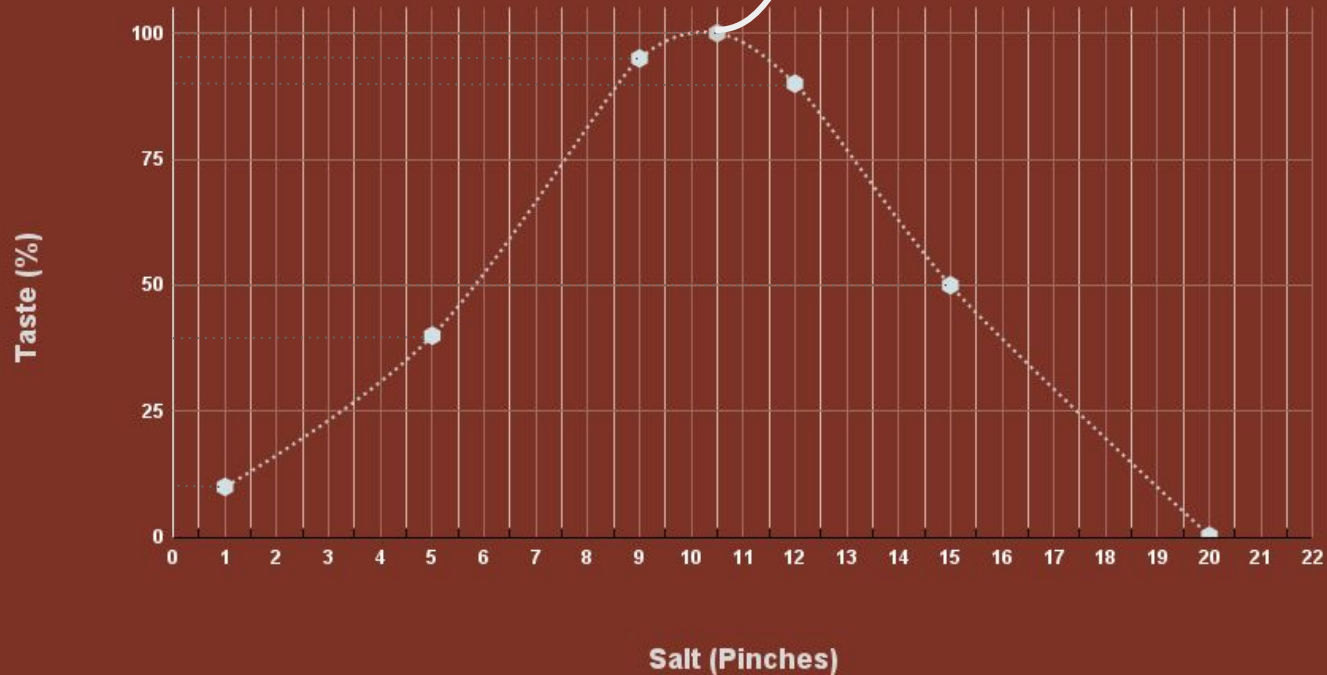
Understanding Derivatives

Taste as a function of Salt



Understanding Derivatives

Taste as a function of Salt



Limits

$$\lim_{x \rightarrow a} f(x) = C$$

Function

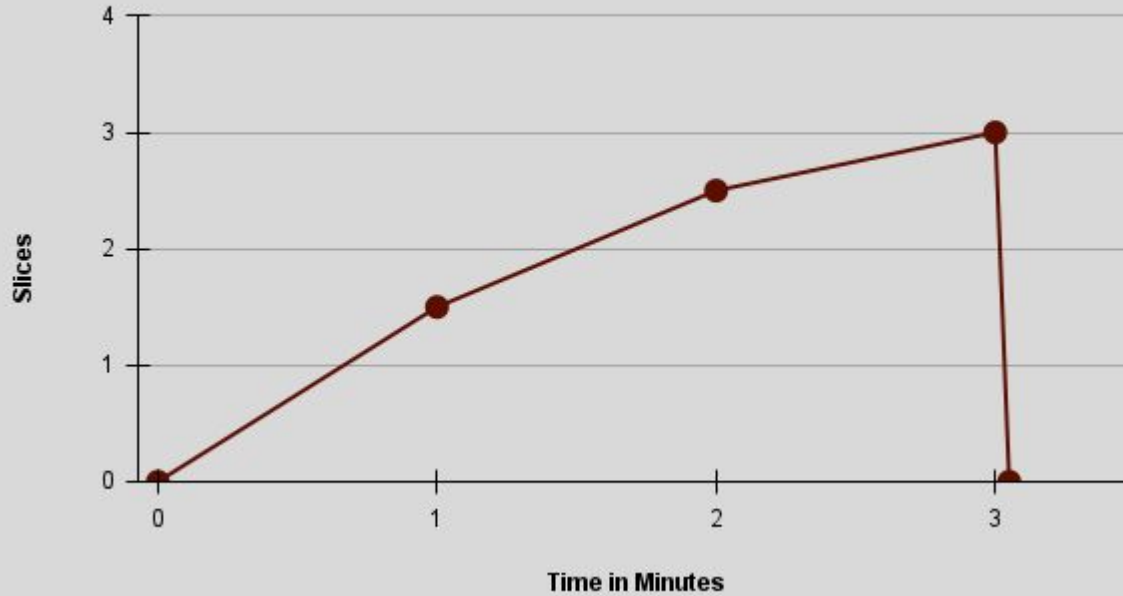
Approaching the input (a) value on the x-axis

Output value the function is approaching after substituting the value of (a) inside the function



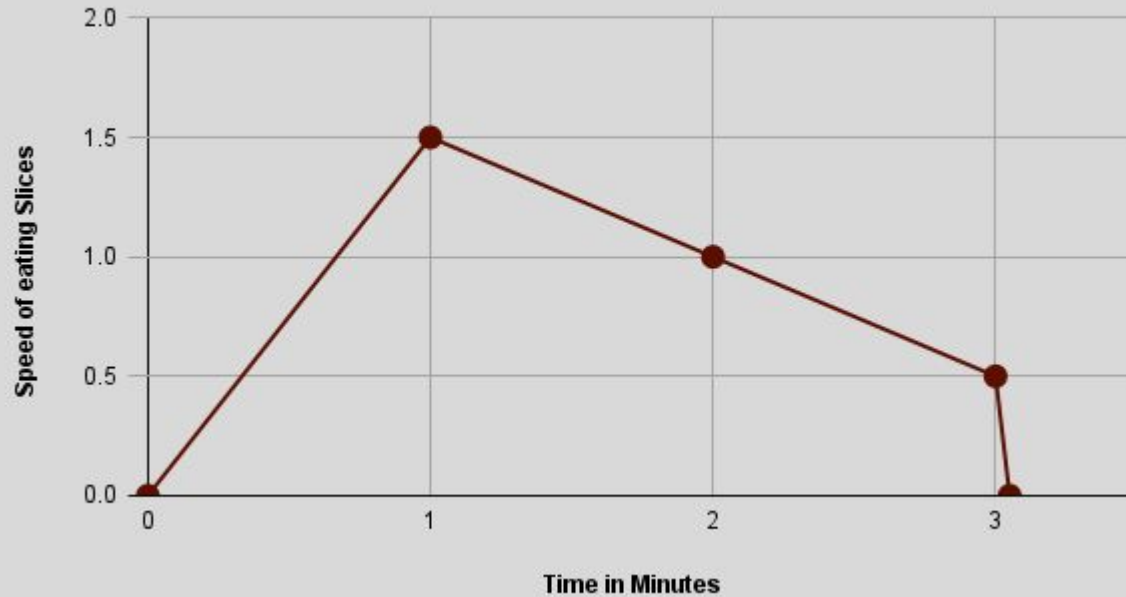
Let's eat Pizza

Slices vs Time in Minutes



Second Derivatives

Speed of eating Slices vs Time



Derivatives

Common Derivatives

$$\frac{d}{dx}(x) = 1$$

$$\frac{d}{dx}(\sin x) = \cos x$$

$$\frac{d}{dx}(\cos x) = -\sin x$$

$$\frac{d}{dx}(\tan x) = \sec^2 x$$

$$\frac{d}{dx}(\sec x) = \sec x \tan x$$

$$\frac{d}{dx}(\csc x) = -\csc x \cot x$$

$$\frac{d}{dx}(\cot x) = -\csc^2 x$$

$$\frac{d}{dx}(\sin^{-1} x) = \frac{1}{\sqrt{1-x^2}}$$

$$\frac{d}{dx}(\cos^{-1} x) = -\frac{1}{\sqrt{1-x^2}}$$

$$\frac{d}{dx}(\tan^{-1} x) = \frac{1}{1+x^2}$$

$$\frac{d}{dx}(a^x) = a^x \ln(a)$$

$$\frac{d}{dx}(e^x) = e^x$$

$$\frac{d}{dx}(\ln(x)) = \frac{1}{x}, \quad x > 0$$

$$\frac{d}{dx}(\ln|x|) = \frac{1}{x}, \quad x \neq 0$$

$$\frac{d}{dx}(\log_a(x)) = \frac{1}{x \ln a}, \quad x > 0$$

Differentiation Rules

Constant Rule	$\frac{d}{dx} [C] = 0$
Power Rule	$\frac{d}{dx} x^n = nx^{n-1}$
Product Rule	$\frac{d}{dx} [f(x)g(x)] = f'(x)g(x) + f(x)g'(x)$
Quotient Rule	$\frac{d}{dx} \left[\frac{f(x)}{g(x)} \right] = \frac{g(x)f'(x) - f(x)g'(x)}{[g(x)]^2}$
Chain Rule	$\frac{d}{dx} [f(g(x))] = f'(g(x)) g'(x)$

Partial derivatives

Multivariable Function: $f(x, y) = x^2y$

$$\frac{\partial f}{\partial x} = \underbrace{\frac{\partial}{\partial x} x^2 y}_{\text{Treat } y \text{ as constant; take derivative.}} = 2xy$$

$$\frac{\partial f}{\partial y} = \underbrace{\frac{\partial}{\partial y} x^2 y}_{\text{Treat } x \text{ as constant; take derivative.}} = x^2 \cdot 1$$

Chain Rule

If $h(x) = f(g(x))$, then:

$$h'(x) = f'(g(x)) \cdot g'(x)$$

Given: $h(x) = (2x + 3)^4$

Let $f(u) = u^4$ and $g(x) = 2x + 3$.

Find $f'(u)$ and $g'(x)$:

$$f'(u) = 4u^3$$

$$g'(x) = 2$$

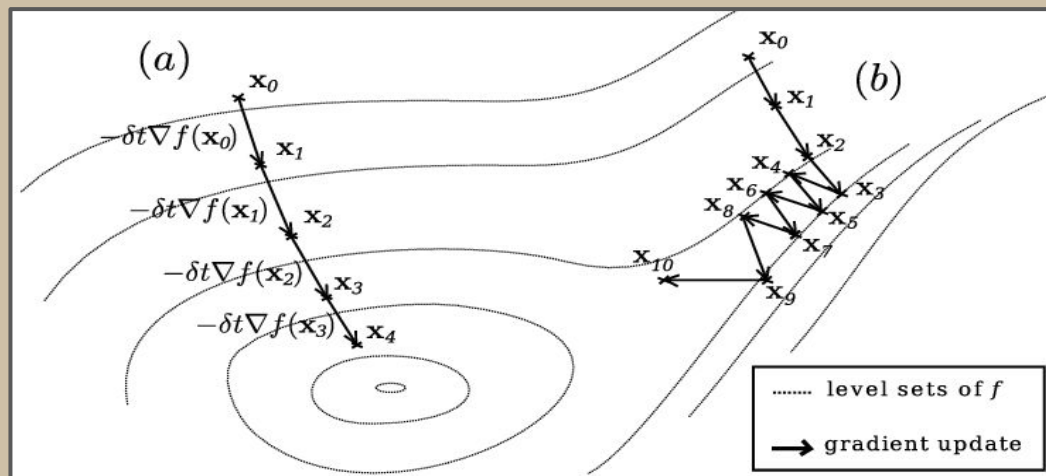
Now, using the chain rule $h'(x) = f'(g(x)) \cdot g'(x)$:

$$h'(x) = 4(2x + 3)^3 \cdot 2$$

$$h'(x) = 8(2x + 3)^3$$

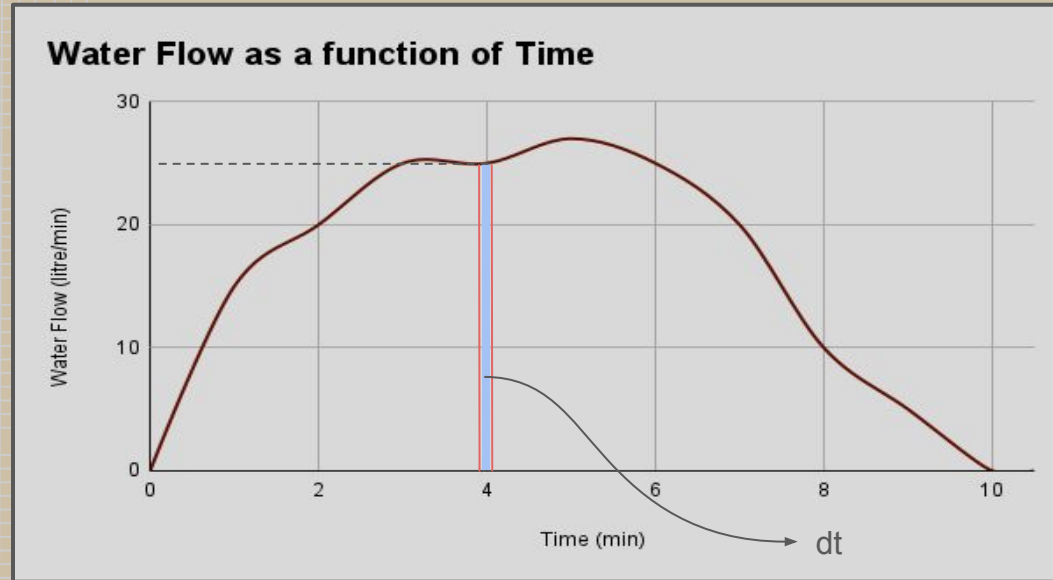
Gradient Descent

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \\ \vdots \end{bmatrix} \quad \nabla f(x_0, y_0) = \begin{bmatrix} \frac{\partial f}{\partial x}(x_0, y_0) \\ \frac{\partial f}{\partial y}(x_0, y_0) \end{bmatrix}$$





Integration

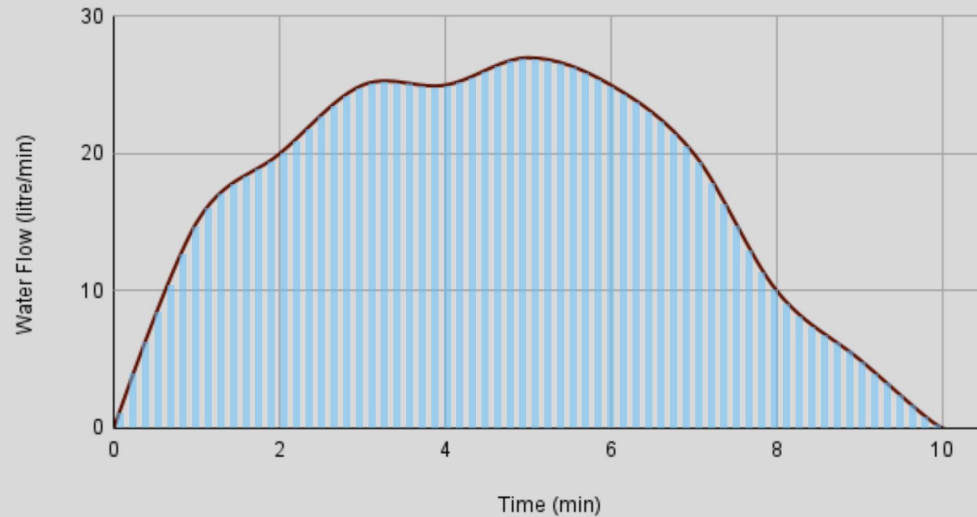


How much total water flowed in the time interval dt ?

$25dt$

Integration

Water Flow as a function of Time



Total water flowed from $t = 0$ to $t = 10$:

$$\text{Total Water} = \int_0^{10} f(t) dt$$

Sum over all intervals from $t = 0$ to $t = 10$:

$$\int_0^{10} f(t) dt \approx \sum (\text{flow rate in each interval}) \times dt$$

Integration

Given: Water flow rate as a function of time $f(t)$.

Total water flowed from $t = 0$ to $t = 10$:

$$\text{Total Water} = \int_0^{10} f(t) dt$$

In the diagram, if $f(t) \approx 25$ liters/min at a specific interval dt :

$$\text{Water in interval} = f(t) \cdot dt \approx 25 \cdot dt$$

Sum over all intervals from $t = 0$ to $t = 10$:

$$\int_0^{10} f(t) dt \approx \sum (\text{flow rate in each interval}) \times dt$$

Recap

- Fundamental concept of Derivatives
- Understanding Types of Slopes and the insights they give
- Second Derivatives
- Partial Derivatives
- Rules of Differentiation
- Chain Rule
- Gradient Descent
- Integration