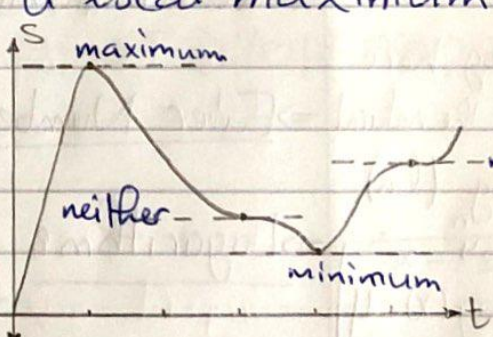


Calculus For Machine Learning and Data Science

Week 1: Derivatives and Optimization:

Importance	Calculus in machine Learning depends largely on derivatives for optimizations (minimizing or maximizing functions).
Definition	Derivative is the instantaneous change of function. This is expressed as the value of the tangent at the meant point. of the slope
	# The maximum or the minimum of a specific function are located where the function has tangent of slope equals zero (derivative equals zero).
	# The slope of the tangent can be zero, but the value doesn't have to be a local maximum or minimum.
	# horizontal line has slope zero.
	
	<div> <div> <p>Notation</p> <p>Lagrange's Notation</p> $f'(x)$ </div> <div> <p>Leibniz's Notation</p> $\frac{dy}{dx} = \frac{d}{dx} f(x)$ </div> </div>
	<div> <div> <p>⇒ <u>Line</u>: $f(x) = ax + b$</p> <p>$f(x) = \text{constant}$</p> </div> <div> <p>$f'(x) = a$</p> <p>$f'(x) = \text{zero}$</p> </div> </div>
	<div> <div> <p>⇒ <u>quadratic</u>: $f(x) = ax^n + b$</p> <p>$f(x) = \frac{1}{x} = x^{-1}$</p> </div> <div> <p>$f'(x) = an x^{n-1}$</p> <p>$f'(x) = (-1)(x^{-2}) = -\frac{1}{x^2}$</p> </div> </div>

definition # The inverse function undoes what the original function does:

→ if $g(x)$ is the inverse of $f(x)$, it is denoted as:

$$g(x) = f^{-1}(x) \neq \frac{1}{f(x)}$$

$$g(f(x)) = x$$

example: $f(x) = x^2$ for $x \geq 0$

$$g(x) = \sqrt{x} \quad \text{for } x \geq 0$$

⇒ Inverse: if $g(x) = f^{-1}(x)$, then $g'(x) = \frac{1}{f'(x)}$ knowing that the coordinates are either $\rightarrow (x, f(x))$ or $\rightarrow (y, g(y))$ or $\rightarrow (g(f(x)), f(x))$ or $\rightarrow (f(x), g(f(x)))$, so the points have to be set according to the plotting used.

⇒ Trigonometric: $f(x) = \sin(x)$ $f'(x) = \cos(x)$
 $f(x) = \cos(x)$ $f'(x) = -\sin(x)$

definition # The exponential (e) = $(1 + \frac{1}{\infty})^\infty = 2.71828182...$ (euler number), and $f(x) = e^x$, then $f^{-1}(y) = \log(y)$

log here is the natural

log (ln)

$$\boxed{e^y = x}$$

$$\log(x) = y$$

⇒ Euler Number Exponential: $f(x) = e^x$ $f'(x) = e^x$

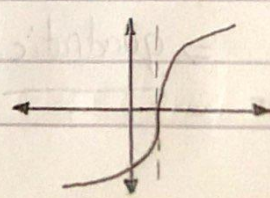
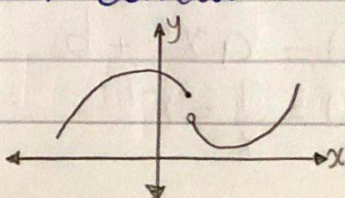
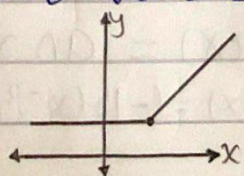
⇒ Logarithm: $f(x) = e^x$
 $f'(x) = e^x$

$$f'(y) = \log(y)$$

$$\frac{d}{dy} f'(y) = \frac{1}{f'(f^{-1}(y))}$$

$$\frac{d}{dy} \log(y) = \frac{1}{e^{\log(y)}} = \frac{1}{y}$$

The function can be non-differentiable at a corner or a cusp, jump discontinuity or has a vertical horizontal.



Multiplication by a scalar \Rightarrow if $f(x) = c \cdot g(x)$, then $f'(x) = c \cdot g'(x)$

sum rule \Rightarrow if $f(x) = g(x) + h(x)$, then $f'(x) = g'(x) + h'(x)$

product rule \Rightarrow if $f(x) = g(x) \cdot h(x)$, then $f'(x) = g'(x) \cdot h(x) + g(x) \cdot h'(x)$

chain rule \Rightarrow if $f(x) = g(h(x))$, then $f'(x) = g'(h(x)) \cdot h'(x)$
 $\frac{d}{dx} f(x) = \frac{dg}{dh} \cdot \frac{dh}{dx} = \frac{dg}{dx}$

The rule is called a chain, because it can be chained to another function, and the process continues.

examples $f(x) = j(i(h(g(x))))$, then $\frac{df}{dx} = \frac{dj}{di} \cdot \frac{di}{dh} \cdot \frac{dh}{dg} \cdot \frac{dg}{dx}$
 $f'(x) = j'(i(h(g(x)))) \cdot i'(h(g(x))) \cdot h'(g(x)) \cdot g'(x)$

SymPy is a python library for symbolic computations. An expression can be defined as follows:

```
from sympy import *
x, y = symbols('x y')
expr = 2 * x**2 - x * y // expr = 2x^2 - xy
expr_manip = x * (expr + x * y + x**3) // expr_manip = x(x^3 + 2x^2 + xy)
expand(expr_manip)
→ x^4 + 2x^3
factor(expr_manip)
→ x^3(x + 2)
expr.evalf(subs = {x:-1, y:2}) // 2(3)^2 - (-1)(2)
→ 4.0
```

evaluate function

```
# import numpy as np
x_array = np.array([1, 2, 3])
f_symb = x**2
try:
    f_symb(x_array)
```

sympy functions can't operate

on numpy
arrays[objects]

```
except TypeError as err:  
    print(err)
```

→ 'Pow' object is not callable

make the
function
numpy-
friendly

```
from sympy.utilities.lambdify import lambdify  
f_symb_numpy = lambdify(x, f_symb, 'numpy')  
f_symb_numpy(x_array)  
→ [149]
```

differentiate
f_symb with
respect to x

```
# diff(f_symb, x)  
→ 2x
```

⇒ Autograd and JAX are the most commonly used frameworks to build neural networks. Their functionality depends on automatic differentiation which is breaking down the function into common basic functions (sin, cos, log, ...etc.) and construct computational graph from them, then chain rule is used to differentiate any node on the graph.

⇒ jax.numpy replaces numpy when jax is used, and can be imported as jnp or np, but jnp sometimes uses different approach:

```
from jax import grad, vmap  
import jax.numpy as jnp  
x_array = np.array([1, 2, 3])  
x_array_jnp = jnp.array([1, 2, 3])  
x_array_jnp[0]
```

→ 1

```
x_array[2] = 4 // [1, 2, 4]
```

```
y_array_jnp = x_array_jnp.at[2].set(4)
```

jnp arrays are like tuples; can be accessed not modified

Most of jax functions work on jnp and np arrays.

```
print(jnp.log(x_array))
```

```
→ [0. 0.6931472 1.0986123]
```

```
print(jnp.log(x_array-jnp))
```

```
→ [0. 0.6931472 1.0986123]
```

To find the derivative ^{of a function} at value of x :

$$P = x^{**2}$$

```
grad(P)(3.0) → always float only
```

```
→ 6.0
```

, and we can use vmap to solve the problem of using integers and float:

```
vmap(grad(P))(x_array-jnp)
```

```
→ [2. 4. 6.]
```

Optimization

Derivatives are used to get the maxima and minima:

→ local maxima have +ve slope before & -ve slope after.

→ local minima have -ve slope before & +ve slope after.

Then the local minima and maxima are then compared to get the absolute (global) maximum or minimum.

⇒ The Square Loss:

Minimize $(x-a_1)^2 + (x-a_2)^2 + \dots + (x-a_n)^2$

Solution $x = \frac{a_1 + a_2 + \dots + a_n}{n}$

⇒ The Log Loss:

If we have a large equation of probabilities, then we take the log of the equation, and use the logarithm properties:

$$\rightarrow \log(ab) = \log a + \log b$$

$$\rightarrow \log(x^a) = a \cdot \log(x)$$

$$\rightarrow \log\left(\frac{a}{b}\right) = \log(a) - \log(b)$$

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

\rightarrow The logarithm of very small numbers is very large negative number.

Logarithm is used to simplify equations; instead of:

- \rightarrow product rule, use sum rule.

- \rightarrow dealing with very tiny probabilities, we deal with large negative numbers.

IF $\log(g(p)) = G(p)$, then we use $-G(p)$, because $0 \leq \text{probability}(p) \leq 1$, so its log always ≤ 0 .