



Department of Data Science

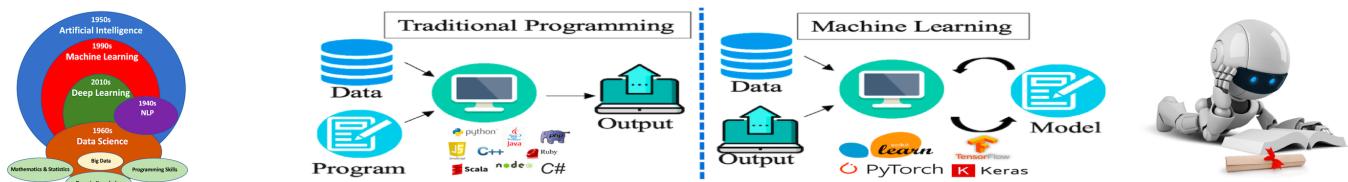
Course: Tools and Techniques for Data Science

Instructor: Muhammad Arif Butt, Ph.D.

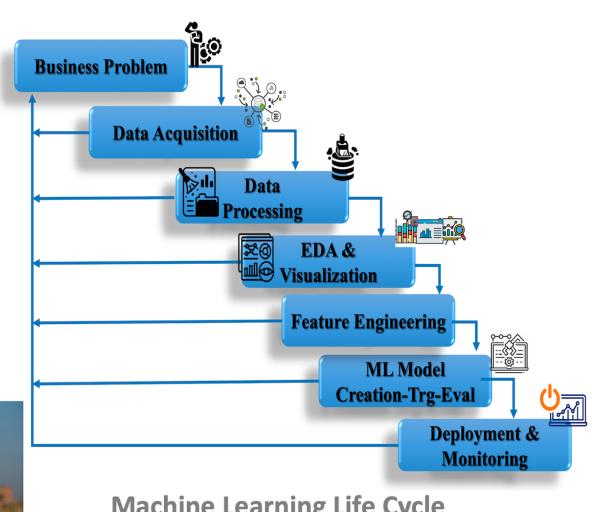
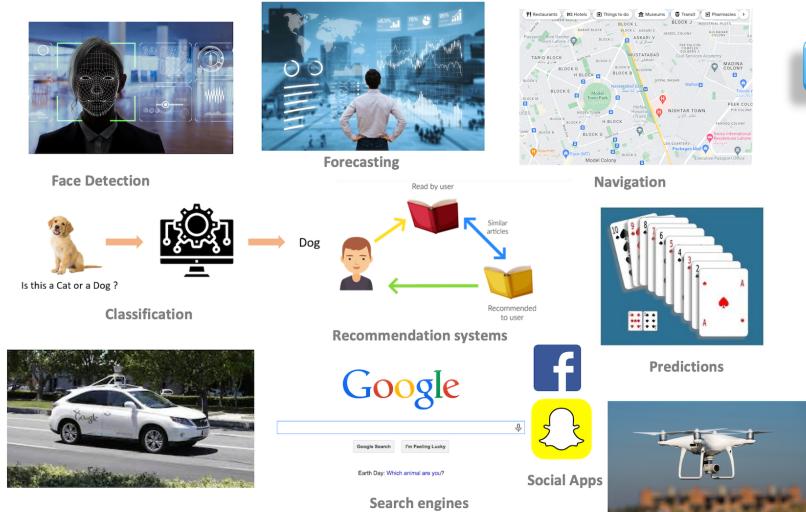
Lecture 6.4 (MLR using Gradient Descent: Math Behind The Curtain)

[Open in Colab](#)

([https://colab.research.google.com/github/arifpcit/data-science/blob/master/Section-4-Mathematics-for-Data-Science/Lec-4.1\(Descriptive-Statistics\).ipynb](https://colab.research.google.com/github/arifpcit/data-science/blob/master/Section-4-Mathematics-for-Data-Science/Lec-4.1(Descriptive-Statistics).ipynb))



ML is the application of AI that gives machines the ability to learn without being explicitly programmed



Learning agenda of this notebook

1. Overview of Calculus
 - What is Calculus?
 - Intuition Behind Derivatives
 - Using Differentiation Rules to Differentiate a Function

- Using First Principle to Differentiate a Function
- Python Code to Differentiate a function using First Principle
- Plotting the function and its derivative
- How to Find Minima & Maxima using Derivatives?
- Partial Derivatives and Gradient

2. Gradient Descent Algorithm

- Mountain Analogy for Gradient Descent Algorithm
- What is Gradient Descent and How to do Gradient Descent?
- Gradient Descent in Python
- Challenges of Gradient Descent Algorithm
 - Local Minima and Unfortunate Starting Point
 - Vanishing Gradient Problem
 - Exploding Gradient Problem

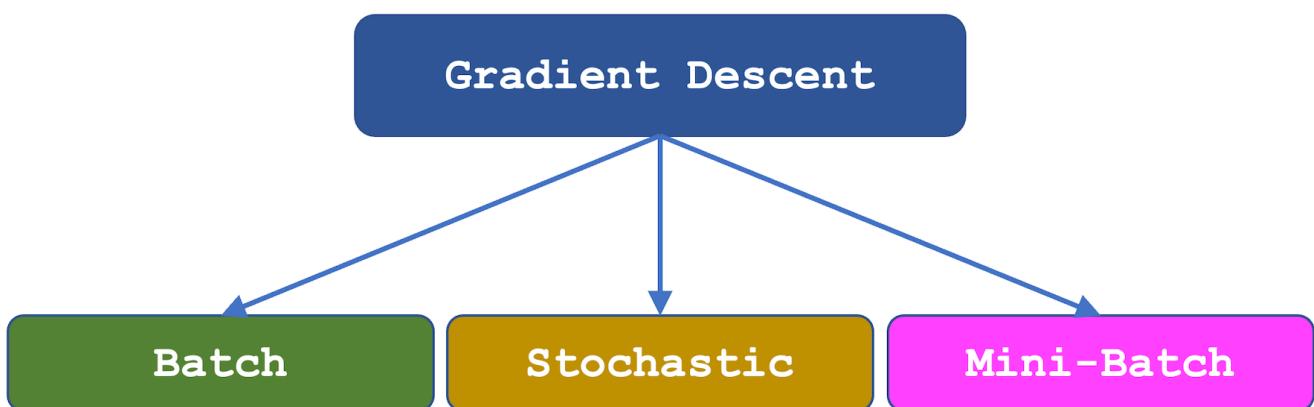
3. Linear Regression Using Gradient Descent Algorithm

- Recap of Linear Regression using OLS
- Example (Simple Linear Regression using Gradient Descent)
- Minimizing Error Function for Gradient Descent Algorithm
- Partial Derivatives of Error Function
- Calculate Regression Coefficients Using Gradient Descent Algorithm
- Fit the Line
- Carryout the Prediction

4. Task To Do

- Gradient Descent for n-D Dataset (Multiple Linear Regression using GD)
- Stochastic Gradient Descent (Multiple Linear Regression using SGD)
- Mini-Batch Gradient Descent (Multiple Linear Regression using Mini-Batch GD)

Gradient Descent is a *first order, iterative algorithm, used to find the local minimum of a differentiable function*



1. Overview of Calculus

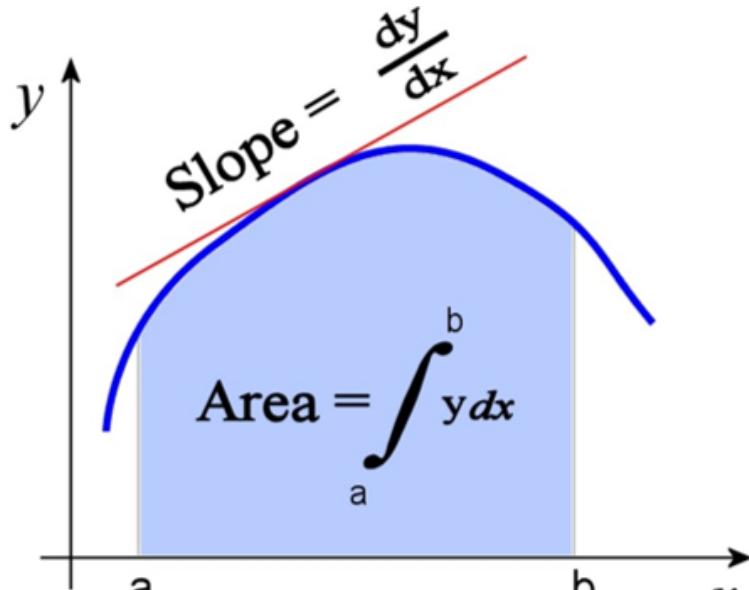
a. What is Calculus?

Calculus is the branch of mathematics that deals with calculating the derivatives and

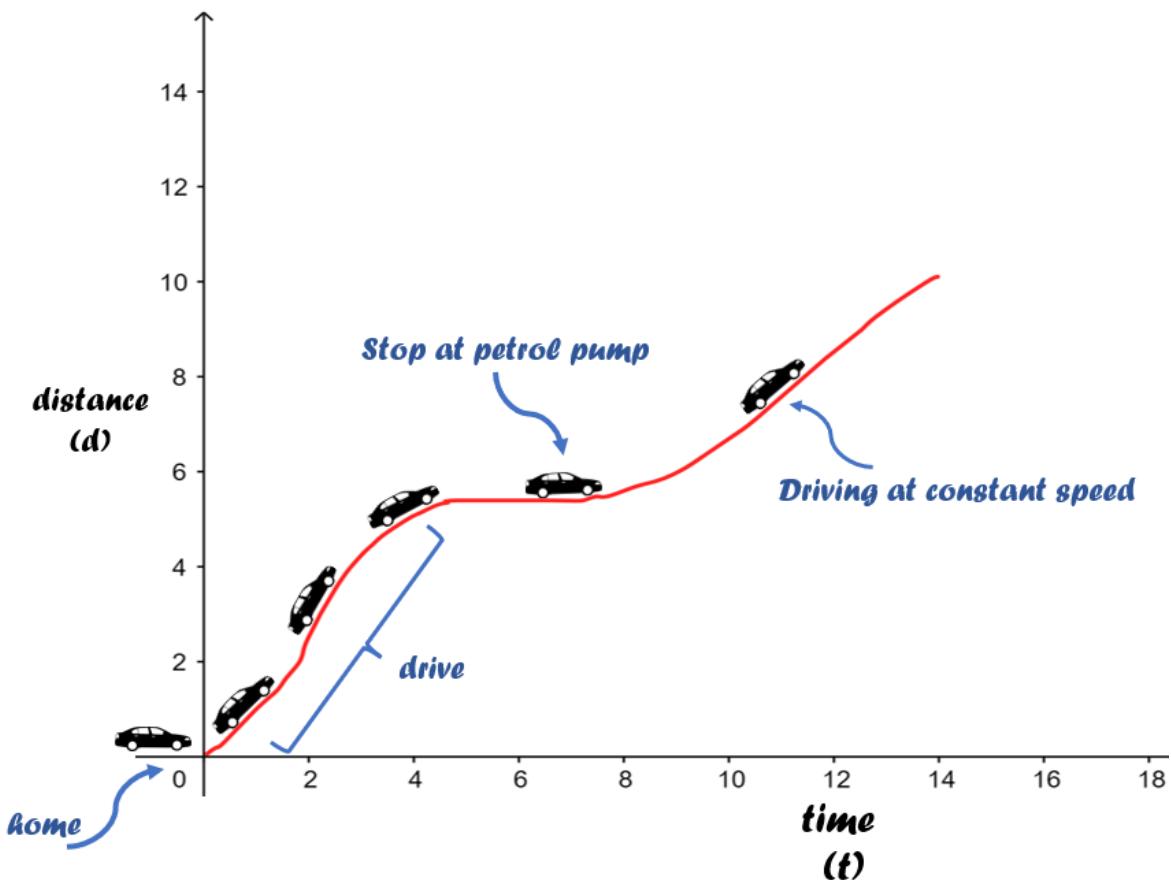
integrals of functions and understanding their properties.

Calculus has two major branches:

- Differential Calculus : Concerns instantaneous rates of change, and the slopes of curves.
- Integral Calculus : Concerns accumulation of quantities (summation of infinitely many small factors to determine some whole), and areas under or between curves.

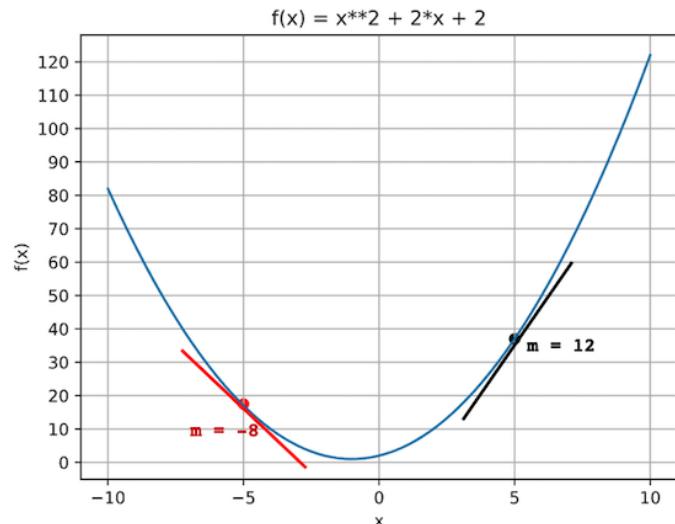
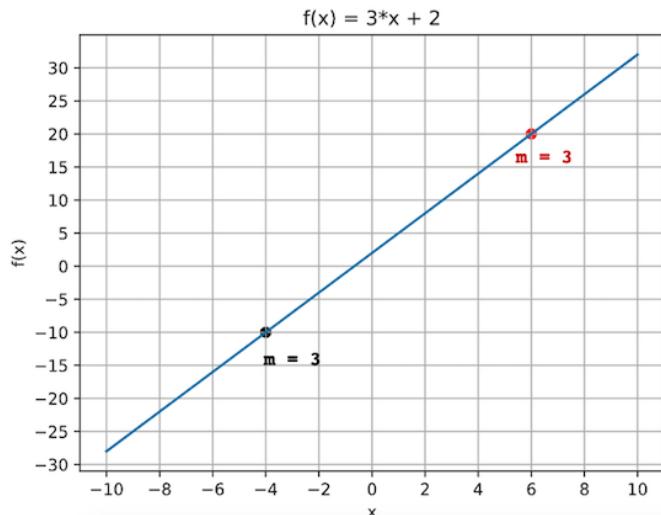


b. Intuition Behind Derivatives



The first derivative of a function represents the slope of a function at a point, and the second derivative describes how the slope changes over the independent variable in the graph.

- Usually, the second derivative of a given function corresponds to the curvature or concavity of the graph. If the second-order derivative value is positive, then the graph of a function is upwardly concave. If the second-order derivative value is negative, then the graph of a function is downwardly open.



Slope	Derivative
<ul style="list-style-type: none"> Used for Linear Equation of single variable It is a constant 	<ul style="list-style-type: none"> Used for Non-Linear Equation of single variable It is a function

Non-Differentiable Functions

- A function $f(x)$ is differentiable, if the derivative of the function exist at every point in its given domain
- Why the following functions are non-differentiable?

- $f(x) = |x|$

- $f(x) = \frac{x}{|x|}$

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

- First function is non-differentiable because there is a corner point in its graph ($x = 0$).
- Second function is non-differentiable because there is a gap in its graph ($x = 0$).
- Third function is non-differentiable because it goes to infinity at $x = 0$.

c. Using Differentiation Rules to Differentiate a Function

Derivative is the rate of change of a function $f(x)$ with respect to an independent variable x .

The process of finding derivative/slope is called differentiation.

(i) Constant Term Rule:

- Constant Term Rule states that derivative of any constant value c is always zero, as we already know that **A contant has no variation so its slope is nothing but zero.**

$$\frac{d}{dx}(c) = 0, \text{ For any value of } c \text{ where } c \in \mathbb{R}$$

(ii) Power Rule (Polynomial/Elementary power rule)

- Power rule is the most commonly used rules of differentiation, which can be used to compute the derivative of some variable x to the power of n e.g., x^n , for any real number $n \neq 0$ such as:

$$\frac{d}{dx} x^n = nx^{n-1}$$

(iii) Constant Multiple Rule

- Constant Multiple Rule is used to compute the derivative of some variable is multiplying with some constant value c . Their derivative can be computed by placing the constant in front of the derivative operator and applying derivation to the rest of the variables such as:

$$\frac{d}{dx}(cy) = c \frac{d}{dx}(y) = c \frac{dy}{dx}$$

(iv) The Sum Rule

- Sum rules states that the derivative of the sum of two terms y and w is the same as the sum of the derivative of y and the derivative of w . Same rule applies when there is a minus (-) sign present between both terms:

$$\frac{d(y + w)}{dx} = \frac{dy}{dx} + \frac{dw}{dx}$$

In [1]:

```
1 import sys
2 !{sys.executable} -m pip -q install sympy
```

Example 1: Differentiate the following function:

$$f(x) = 5x^3 + 3x^2 - 7x + 20$$

In [2]:

```
1 import numpy as np
2 import sympy as sym
3 from IPython.display import display # make the equations look nicer
4
5 x = sym.symbols('x') # create symbolic variables in sympy
6 fx = 5*x**3 + 3*x**2 - 7*x + 20 # create a function
7
8 # The sym.diff() is passed function and reference variable, and it returns the
9 df = sym.diff(fx, x)
10
11 display(df)
```

$$15x^2 + 6x - 7$$

Example 2: Find the derivative/slope of the following function at the point (2, 5)

$$f(x) = x^3 + 1$$

In [3]:

```
1 x = sym.symbols('x')
2 fx = x**3 + 1
3 df = sym.diff(fx, x)
4 display(df)
5 print("Slope of function at point (2,5): ", 3*2**2)
```

$$3x^2$$

Slope of function at point (2,5): 12

Example 3: Find the derivative/slope of the following function at the point (9, 12)

$$f(x) = 4\sqrt{x}$$

In [4]:

```
1 x = sym.symbols('x')
2 fx = 4 * (x**(1/2))
3 df = sym.diff(fx, x)
4 display(df)
5 print("Slope of function at point (9, 12): ", 2/(9**.5))
```

$$\frac{2.0}{x^{0.5}}$$

Slope of function at point (9, 12): 0.6666666666666666

(v) The Chain Rule

- Chain Rule has many applications in machine learning and one of them is Gradient descent algorithm. One of its important application is Backpropagation algorithm, which is used to train neural nets. It is based on the idea of **nested/composite functions**. It is an easy way to find derivative of nested functions. We can differentiate complex functions by splitting them in simple one and then applying differentiation using chain rule. For example:

$$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx}$$

Example 4: Find the derivative of the following function:

$$f(x) = y = (2x^2 + 8)^2$$

Solution:

- Suppose $u = (2x^2 + 8)$, so $y = u^2$
- First compute the derivative of (dy/du): $\frac{dy}{du} = 2u = 2(2x^2 + 8) = 4x^2 + 16$
- Now compute the derivative of (du/dx): $\frac{du}{dx} = 2x^2 + 8 = 4x$

$$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx} = (4x^2 + 16)(4x) = 16x^3 + 64x = 8x(2x^2 + 8)$$

In [5]:

```
1 x = sym.symbols('x')
2 fx = (2*x**2 + 8)**2
3 df = sym.diff(fx)
4 display(df)
```

$$8x(2x^2 + 8)$$

Example 5: Find the derivative of the following function:

$$f(x) = (3x^2 + x)^6$$

In [6]:

```
1 x = sym.symbols('x')
2 fx = (3*x**2 + x)**6
3 df = sym.diff(fx)
4 display(df)
```

$$(36x + 6)(3x^2 + x)^5$$

Example 6: Find the derivative of the following function:

$$f(x) = \sqrt{x^2 - x}$$

In [7]:

```
1 x = sym.symbols('x')
2 fx = (x**2 - x)**.5
3 df = sym.diff(fx)
4 display(df)
```

$$\frac{1.0x - 0.5}{(x^2 - x)^{0.5}}$$

To Do: Find the derivative of the following functions:

- $f(x) = (\sin(x))^2$

- $f(x) = \ln(\sin(x))$
- $f(x) = \ln(x)\sin(x)$

(vi) The Product Rule

- Product rule enable us to take the two variable (let say w and z) for which we want to compute derivative, and compute the derivative of those two variables separately by splitting them in two different terms such as:

$$\frac{d(wz)}{dx} = w \frac{dz}{dx} + z \frac{dw}{dx}$$

Example 7: Find the derivative of the following function:

$$f(x) = x\sqrt{2x+1}$$

Solution:

- Suppose $w = x$ and $z = (2x+1)^{1/2}$
- $\frac{dw}{dx} = 1$
- $\frac{dz}{dx} = \frac{1}{(2x+1)^{1/2}}$

$$\frac{dy}{dx} = w \frac{dz}{dx} + z \frac{dw}{dx} = x \frac{1}{(2x+1)^{1/2}} + (2x+1)^{1/2}$$

In [8] :

```

1 x = sym.symbols('x')
2 fx = x * (2*x+1)**.5
3 df = sym.diff(fx)
4 display(df)

```

$$\frac{1.0x}{(2x+1)^{0.5}} + (2x+1)^{0.5}$$

(vii) The Quotient (Fraction) Rule

- Quotient rule is much similar to product rule, but here we are dividing two variables instead of multiplying them.

$$\frac{d}{dx}\left(\frac{w}{z}\right) = \frac{z \frac{dw}{dx} - w \frac{dz}{dx}}{z^2}$$

Example 8: Find the derivative of the following function at the point (4, 6)

$$f(x) = \frac{3x}{x-2}$$

Solution:

- Suppose $w = 3x$ and $z = x - 2$
- $\frac{dw}{dx} = 3$
- $\frac{dz}{dx} = 1$

$$\frac{dy}{dx} = \frac{(x-2)*3 - 3x*1}{(x-2)^2} = \frac{3x-6-3x}{(x-2)^2} = \frac{-6}{(x-2)^2}$$

In [9]:

```

1 x = sym.symbols('x')
2 fx = 3 * x / (x-2)
3 df = sym.diff(fx, x)
4 display(df)
5 print("Slope of function at point (4,6): ", (-6)/((4-2)**2))

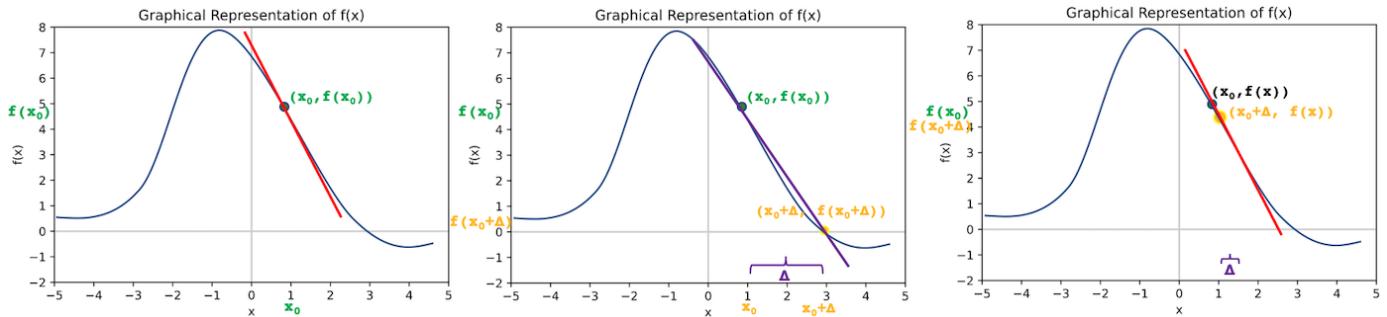
```

$$-\frac{3x}{(x-2)^2} + \frac{3}{x-2}$$

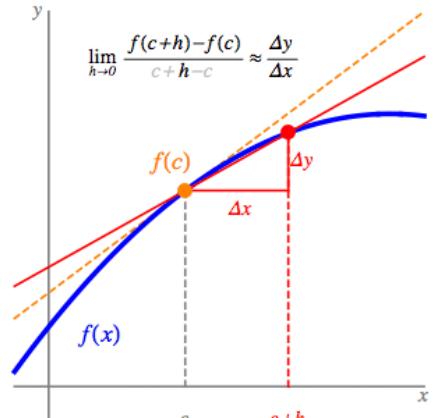
Slope of function at point (4,6): -1.5

d. Differentiating Functions using First Principle

Derivative by first principle (Delta method) is a technique that allows us to use limits to compute the slope of the line at any given point along the curve.



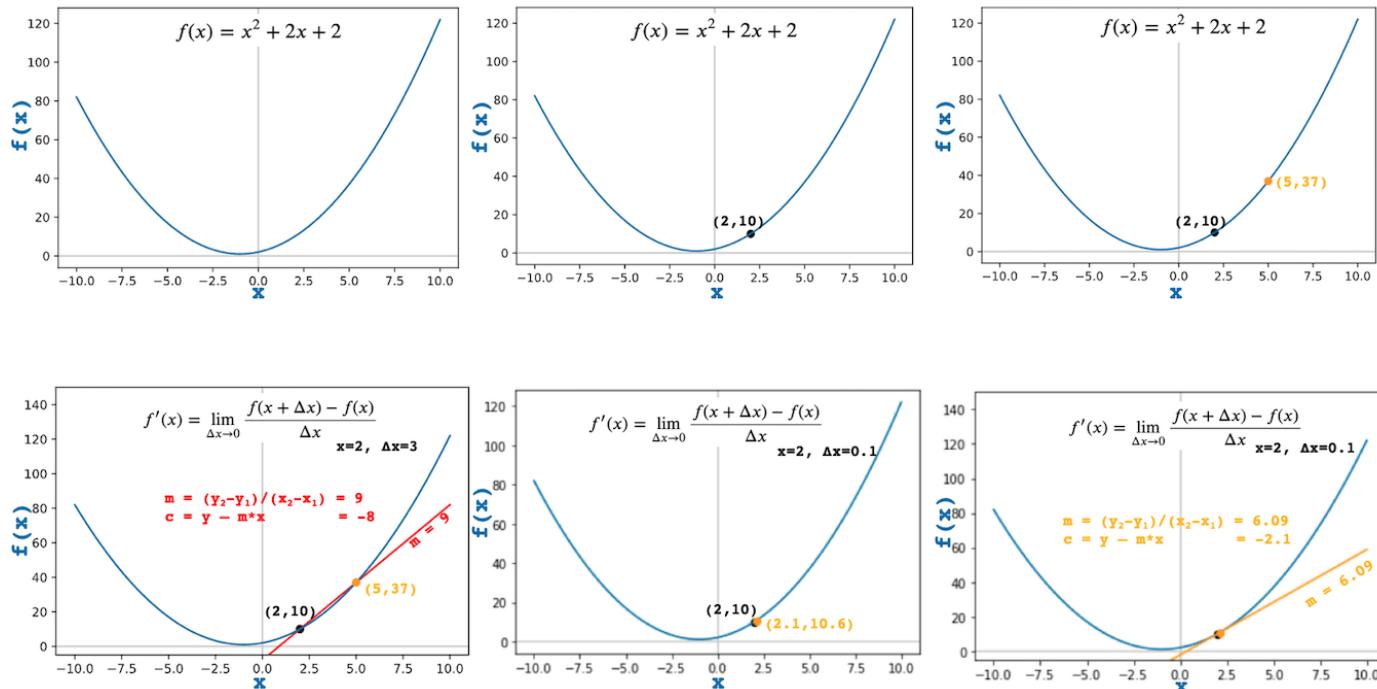
Limit Definition of the Derivative $f'(c)$



$$m = \frac{\text{change in } y}{\text{change in } x} = \frac{y_2 - y_1}{x_2 - x_1} \approx \frac{f(x_0 + \Delta) - f(x_0)}{(x_0 + \Delta) - x_0} \approx \frac{f(x_0 + \Delta) - f(x_0)}{\Delta}$$

Example: Find the derivative of the following function at the point (2, 10)

$$f(x) = x^2 + 2x + 2$$



Example 1A: Find the derivative of the following function at the point (2, 10) **using paper and pencil** using Differentiation Rules, and using First Principle

$$f(x) = x^2 + 2x + 2$$

$$f'(x_0) = \lim_{\Delta \rightarrow 0} \frac{f(x_0 + \Delta) - f(x_0)}{\Delta} \quad \dots \dots (i)$$

$$f'(2) = \lim_{\Delta \rightarrow 0} \frac{f(2 + \Delta) - f(2)}{\Delta} \quad \dots \dots (ii)$$

- Substitute, $f(2) = (2)^2 + 2 * 2 + 2$
and $f(2 + \Delta) = (2 + \Delta)^2 + 2(2 + \Delta) + 2$ in equation (ii):

$$f'(2) = \lim_{\Delta \rightarrow 0} \frac{(2+\Delta)^2 + 2(2+\Delta) + 2 - ((2)^2 + 2*2+2)}{\Delta}$$

$$f'(2) = \lim_{\Delta \rightarrow 0} \frac{\Delta^2 + 6\Delta}{\Delta}$$

$$f'(2) = \lim_{\Delta \rightarrow 0} \Delta + 6 \quad \dots \dots (iii)$$

- Apply the limit:

$$f'(2) = 6$$

Example 1B: Find the derivative of the following function at the point (2, 10) **using Python code** (Use different values of Δ)

$$f(x) = x^2 + 2x + 2$$

$$\frac{dy}{dx} = f'(x_0) = \lim_{\Delta \rightarrow 0} \frac{f(x_0 + \Delta) - f(x_0)}{\Delta}$$

In [10]:

```
1 # define your function (it is passed two values (x and y) and it returns f(x,y))
2 def f(x):
3     return x**2 + 2*x + 2
```

In [11]:

```
1 # function that calculates the derivative at a specific value of x, with a specific delta
2 def deriv(f, x, delta):
3     return (f(x+delta) - f(x)) / delta
```

In [12]:

```
1 deriv(f, 2, 0.01)
```

Out[12]:

6.00999999999849

In [13]:

```
1 # creating a list having values of Δx from 1 to smaller values of Δx (approaches zero)
2 deltas = [3, 2, 1, 0.1, 0.01, 0.001, 0.0001, 0.00001]
3 # calculate and display slopes for different values of delta
4 for delta in deltas:
5     print("Slope at x={0} and Δ={1} = {2:.5f}".format(2, delta, deriv(f, 2, delta)))
```

Slope at x=2 and Δ=3 = 9.0

Slope at x=2 and Δ=2 = 8.0

Slope at x=2 and Δ=1 = 7.0

Slope at x=2 and Δ=0.1 = 6.1

Slope at x=2 and Δ=0.01 = 6.01

Slope at x=2 and Δ=0.001 = 6.001

Slope at x=2 and Δ=0.0001 = 6.0001

Slope at x=2 and Δ=1e-05 = 6.0

In [14]:

```
1 # y = mx + c    ==> c = y - mx
2 x0 = 2
3 y0 = 10
4 m = 6
5 c = y0 - m * x0
6 c
```

Out[14]:

-2

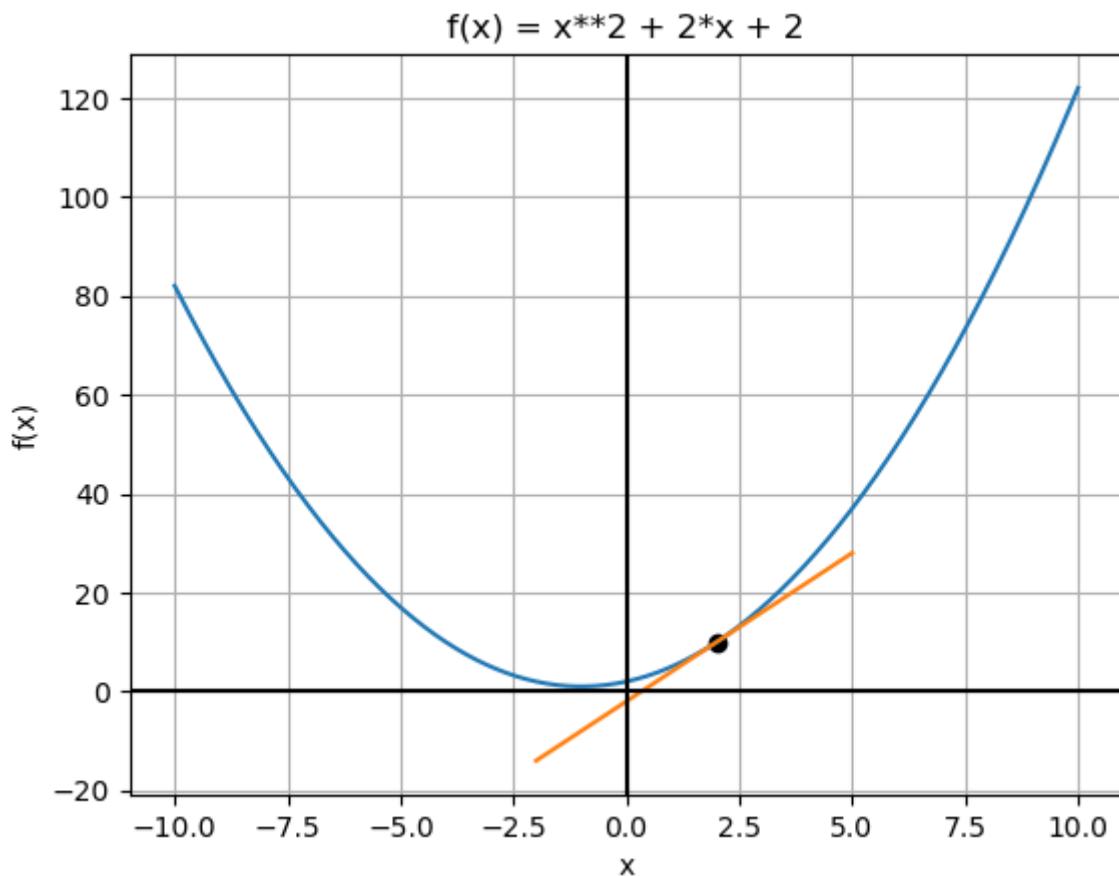
Example 1C: Given the following function

$$f(x) = x^2 + 2x + 2$$

The slope/derivative of this function at point (2, 10) is 6 and y-intercept is -2. Use matplotlib to plot the function, the point (2, 10) and the tangent line at that point.

In [15]:

```
1 import numpy as np
2 from matplotlib import pyplot as plt
3 %matplotlib inline
4 fig, ax = plt.subplots(nrows=1, ncols=1)
5
6 x = np.linspace(-10, 10, 1000) # x values for the function
7 y = x**2 + 2*x + 2           # corresponding y values of the function
8 ax.plot(x,y)                 # plot the function
9
10 ax.scatter(2, 10, c = 'black') # draw point P
11
12 m=6
13 c=-2
14 x1 = np.linspace(-2, 5, 100) # x values for the tangent line
15 y1 = m*x1 + c               # corresponding y values for the tangent line
16 ax.plot(x1,y1)
17
18
19 plt.xlabel('x')
20 plt.ylabel('f(x)')
21 plt.title("f(x) = x**2 + 2*x + 2")
22 plt.grid()
23 plt.axvline(x=0, color='black')
24 plt.axhline(y=0, color='black')
25 plt.show();
```



Task 1: Verify that the derivative of the following function at the point $(9, 12)$ is $\frac{2}{3}$. Do it by differentiation rules, and then by first principle on paper pencil, verify using Python code, and finally draw plot of the function and its derivative.

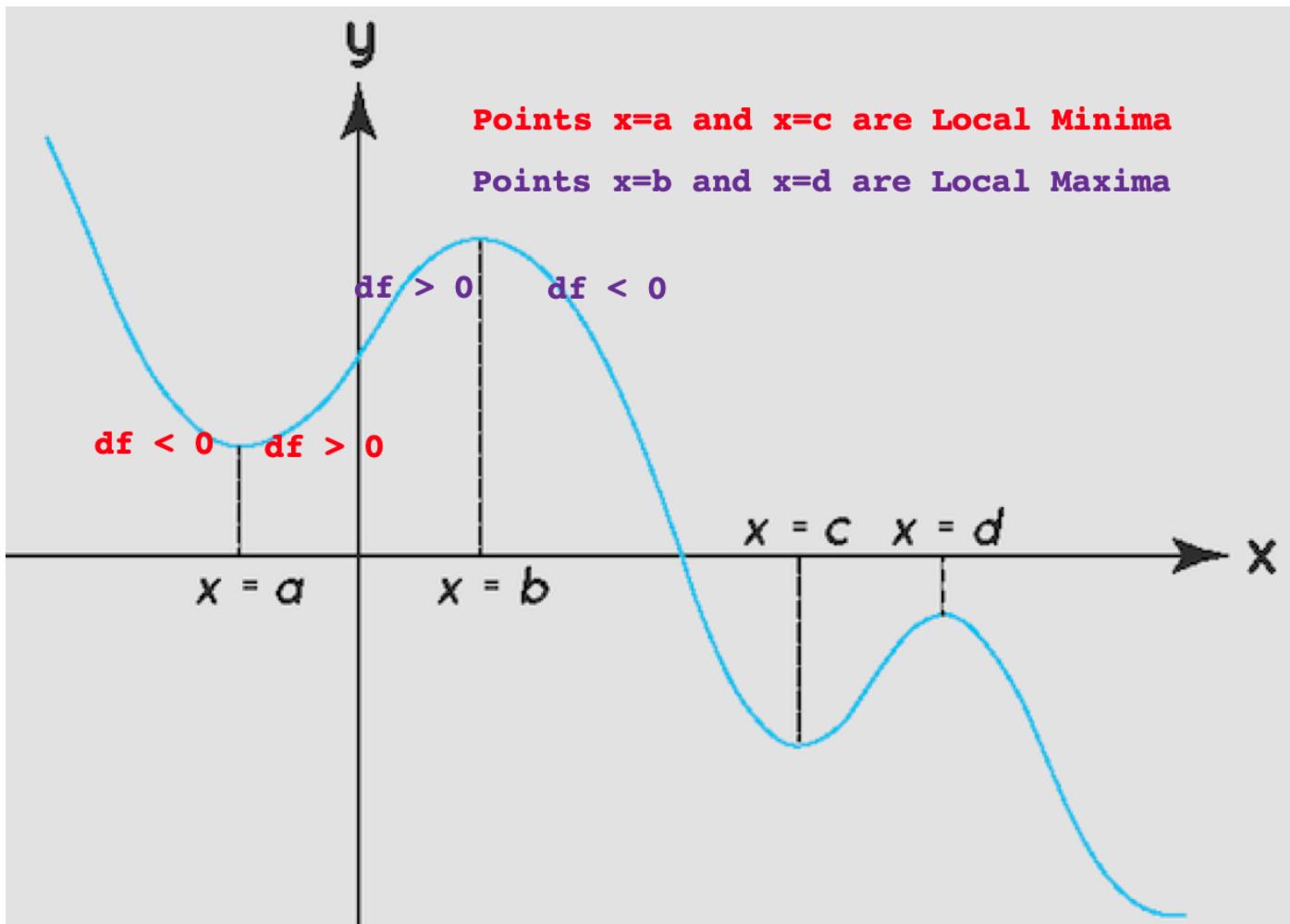
$$f(x) = 4\sqrt{x}$$

Task 2: Verify that the derivative of the following function at the point $(4, 6)$ is $\frac{-3}{2}$. Do it by differentiation rules, and then by first principle on paper pencil, verify using Python code, and finally draw plot of the function and its derivative.

$$f(x) = \frac{3x}{\underline{\hspace{1cm}}}$$

e. How to Find Minima & Maxima using Derivatives?

Maxima and minima (Extrema) of a function are the largest and smallest value of the function respectively, either within a given range or on the entire domain.



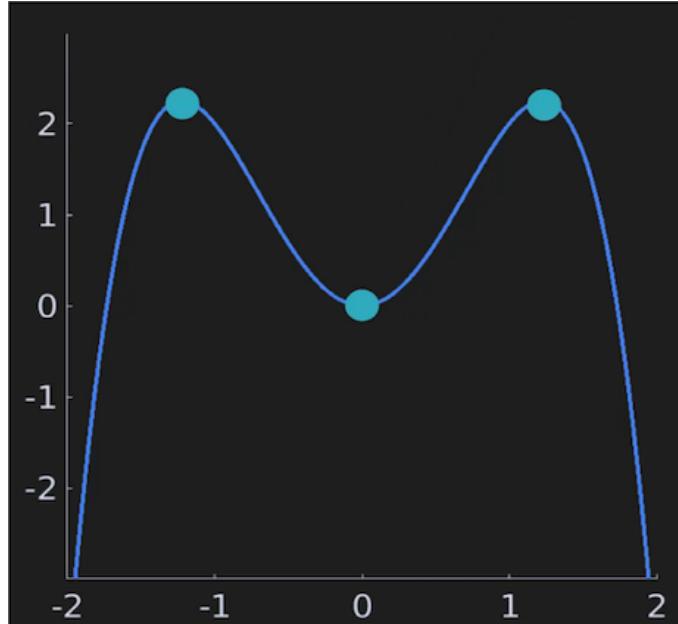
- Second derivative of a function tells us the curvature of the function's graph. If the second-order derivative value is positive, then the graph of a function is upwardly concave. If the second-order derivative value is negative, then the graph of a function is downwardly concave. If the second-order derivative value is also zero, then we cannot conclude anything about the curvature.

- If $f'(x) = 0$ and $f''(x) > 0$ at $x = a$, then the function $f(x)$ has a local minimum at a . For example $f(x) = x^2$

- If $f'(x) = 0$ and $f''(x) < 0$ at $x = a$, then the function $f(x)$ has a local maximum at a . For example $f(x) = -x^2$
- If $f'(x) = 0$ and $f''(x) = 0$ at $x = a$, then we cannot conclude. For example $f(x) = x^4$ or $f(x) = -x^4$. For this we have to analyze the first derivative.

- **Local minima:** If $f'(x)$ changes sign from negative to positive, as x increases via point a , then $f(a)$ gives the minimum value of the function in that range.
- **Local maxima:** If $f'(x)$ changes sign from positive to negative, as x increases via point a , then $f(a)$ gives the maximum value of the function in that range.

Example Function:

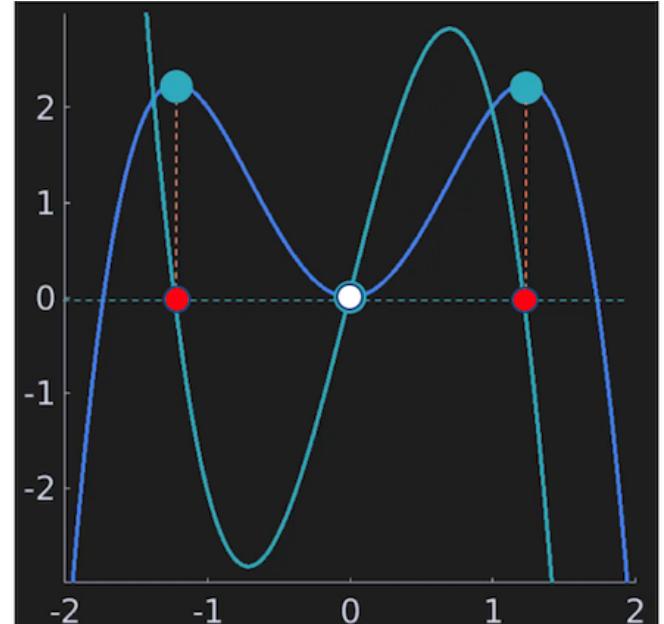


$$f(x) = -x^4 + 3x^2$$

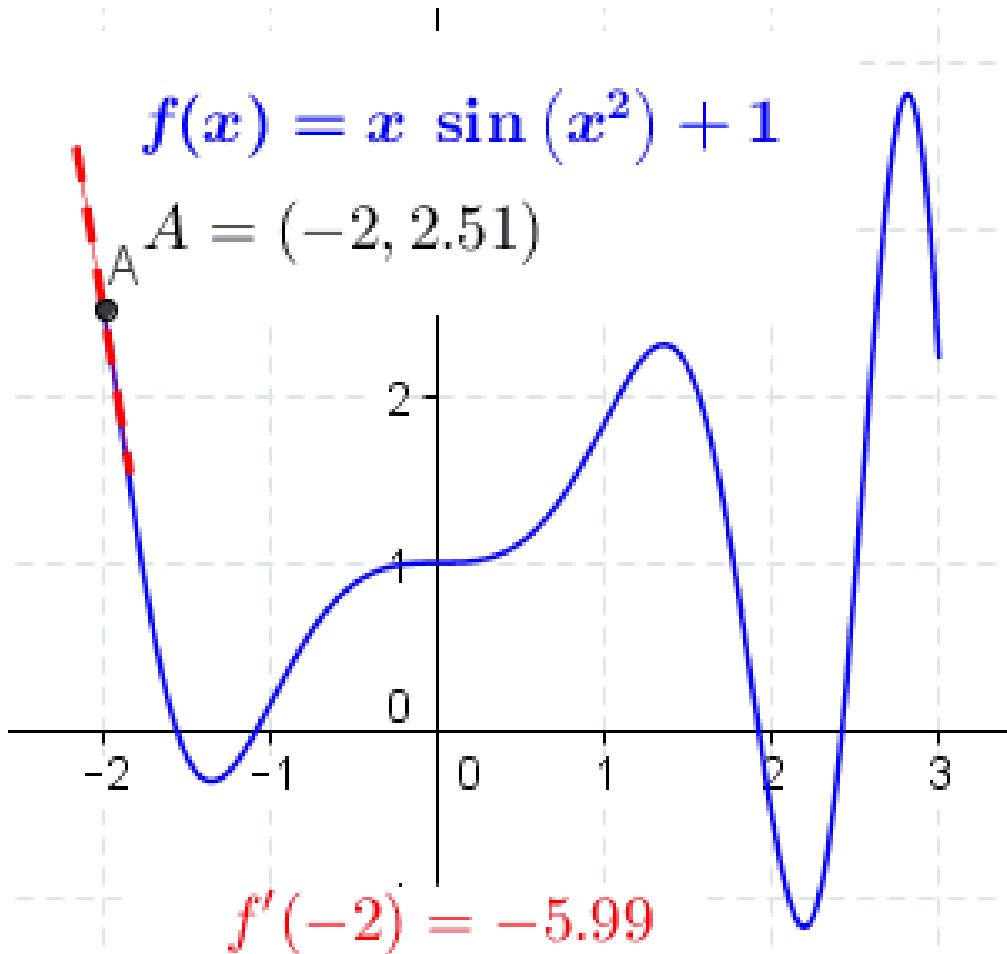
$$f'(x) = -4x^3 + 6x$$

$$-4x^3 + 6x = 0$$

$$x = -\sqrt{\frac{3}{2}}, 0, \sqrt{\frac{3}{2}}$$



Animation:



f. Partial Derivatives (Gradient) of a Function

Gradient of a function at any point is just a collection of all its partial derivatives with respect to all the dimensions of that function.

- In case of function of two or more variables we take partial derivatives.
- To find the partial derivative of a function with respect to one variable, all the other variables are treated as constant.

- **Example 1:**
$$f(x, y) = x^2 - 3xy$$
$$\frac{\partial f}{\partial x} = 2x - 3y \quad \frac{\partial f}{\partial y} = -3x$$

- **Example 2:**
$$f(x, y) = 2x^2 + 4xy$$
$$\frac{\partial f}{\partial x} = 4x + 4y \quad \frac{\partial f}{\partial y} = 4x$$

- **Example 3:** $f(x, y, z) = z^3 - x^2y$
 $\frac{\partial f}{\partial x} = -2xy$ $\frac{\partial f}{\partial y} = -x^2$ $\frac{\partial f}{\partial z} = 3z^2$
- **Example 4:** $f(x, y, z) = x^4 - 3xyz$
 $\frac{\partial f}{\partial x} = 4x^3 - 3yz$ $\frac{\partial f}{\partial y} = -3xz$ $\frac{\partial f}{\partial z} = -3x$

(i) 3-D Scatter Plot

Example: Draw a scatter plot of some random values of the following function:

$$f(x, y) = x^2 - 3xy$$

In [43]:

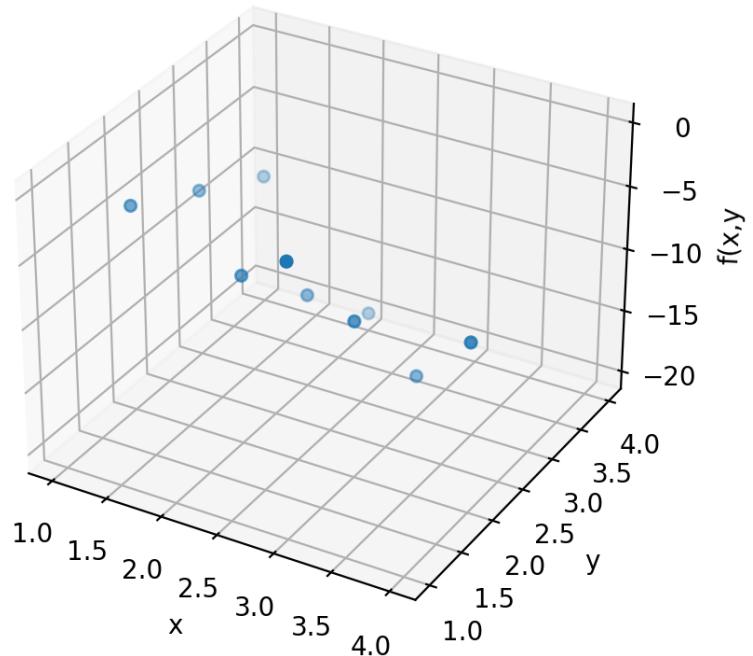
```

1 import numpy as np
2 from matplotlib import pyplot as plt
3
4 # define your function (it is passed two values (x and y) and it returns f(x,y))
5 def func(x,y):
6     return x**2 - 3*x*y
7
8 x = np.array([2, 1, 1, 3, 1, 3, 3, 2, 2, 4])
9 y = np.array([2, 3, 2, 3, 4, 1, 2, 3, 4, 2])
10 fxy = func(x,y)
11 print("x: ", x)
12 print("y: ", y)
13 print("f(x,y): ", fxy)

x: [2 1 1 3 1 3 3 2 2 4]
y: [2 3 2 3 4 1 2 3 4 2]
f(x,y): [-8 -8 -5 -18 -11  0 -9 -14 -20 -8]
```

In [44]:

```
1 %matplotlib notebook
2 ax = plt.axes(projection='3d')
3 ax.set_xlabel('x')
4 ax.set_ylabel('y')
5 ax.set_zlabel('f(x,y)')
6
7 ax.scatter3D(x,y,fxy);
```



(ii) 3-D Surface Plot

Example: Draw a surface plot of some random values of the following function:

$$f(x, y) = x^2 - 3xy$$

In [45]:

```
1 import numpy as np
2 from matplotlib import pyplot as plt
3
4 def func_xy(x,y):
5     return x**2 - 3*x*y
6
7 x = np.linspace(-5,5,20)
8 y = np.linspace(-5,5,20)
9 fxy = func_xy(x,y)
10 print("x.shape: ", x.shape)
11 print("y.shape: ", y.shape)
12 print("f(x,y).shape: ", fxy.shape)
```

```
x.shape: (20,)
y.shape: (20,)
f(x,y).shape: (20,)
```

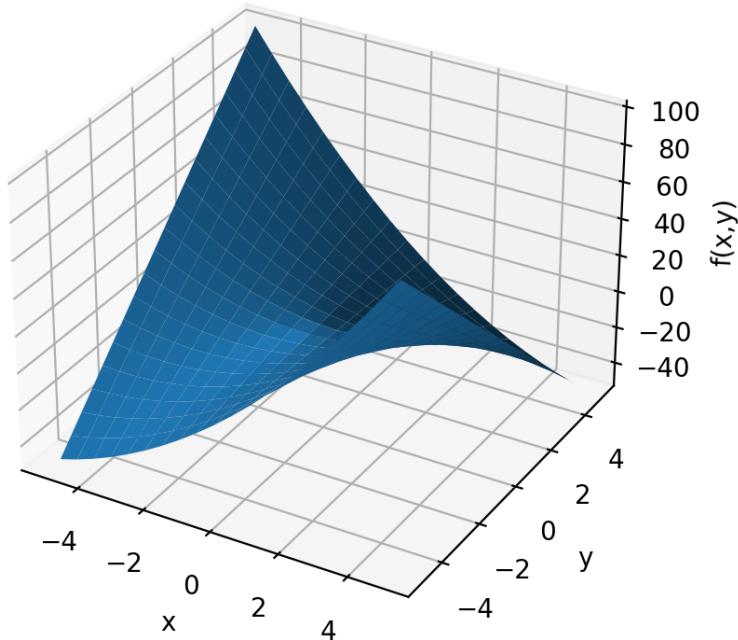
In [46]:

```
1 X,Y = np.meshgrid(x,y)
2 fxy = func_xy(X,Y)
3
4 print("X.shape: ", X.shape)
5 print("Y.shape: ", Y.shape)
6 print("f(x,y).shape: ", fxy.shape)
```

```
X.shape: (20, 20)
Y.shape: (20, 20)
f(x,y).shape: (20, 20)
```

In [47]:

```
1 %matplotlib notebook
2 ax = plt.axes(projection='3d')
3 ax.set_xlabel('x')
4 ax.set_ylabel('y')
5 ax.set_zlabel('f(x,y)')
6
7 ax.plot_surface(X,Y,fxy);
```



(iii) Calculate the gradient of a Function with Two Variables

Example 1: Find the gradient of the following function at the point $(4, 2, -8)$

$$f(x, y) = x^2 - 3xy$$

Using Differentiation Rules:

$$\frac{\partial f}{\partial x} = 2x - 3y$$

$$\frac{\partial f}{\partial y} = -3x$$

Gradient at $x=4$ and $y=2$: $\nabla = [2, -12]$

Using First Principle:

$$\frac{dy}{dx} = f'(x_0) = \lim_{\Delta \rightarrow 0} \frac{f(x_0 + \Delta) - f(x_0)}{\Delta}$$

In [21]:

```
1 # define your function (it is passed two values (x and y) and it returns f(x,y))
2 def f(x, y):
3     return x**2 - 3*x*y
```

In [22]:

```
1 # grad() function returns a vector containing partial derivatives of function f
2 # idx is the index of *args list, specifying the variable with respect to which
3 def grad(f, idx, *args):
4     delta = 0.00001
5     fxy = f(*args)          #value of function at the values x and y
6
7     args = list(args)
8     args[idx] += delta
9     fxy1 = f(*args)         #value of function at the values x+Δ and y OR x and y+Δ
10
11    return (fxy1-fxy)/delta      # First Principle
```

In [23]:

```
1 #calculate the partial derivative w.r.t x, at x=4,y=2
2 grad(f, 0, 4, 2)
```

Out[23]:

2.0000099997474763

In [24]:

```
1 #calculate the partial derivative w.r.t y, at x=4,y=2
2 grad(f, 1, 4, 2)
```

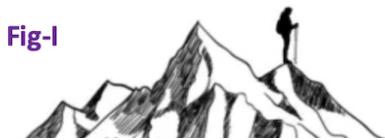
Out[24]:

-12.000000000256248

2. Gradient Descent Algorithm

Gradient Descent algorithm is the most important and most fundamental algorithm in all of DL and is also used in ML & AI

a. Mountain Analogy for Gradient Descent Algorithm

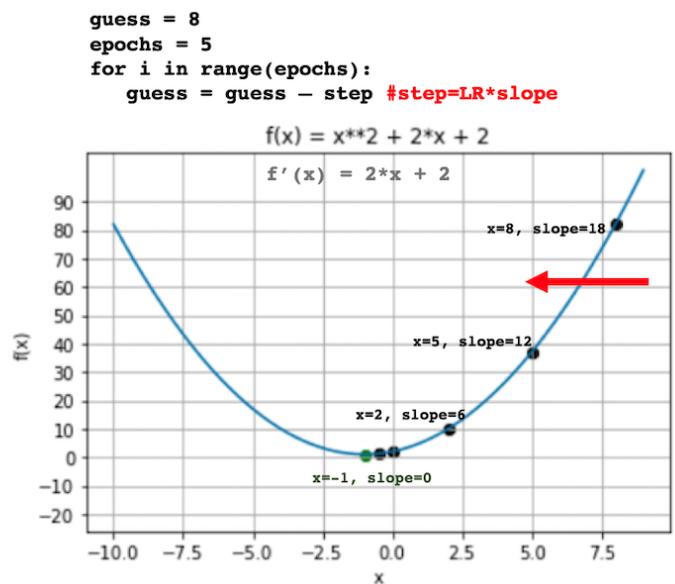
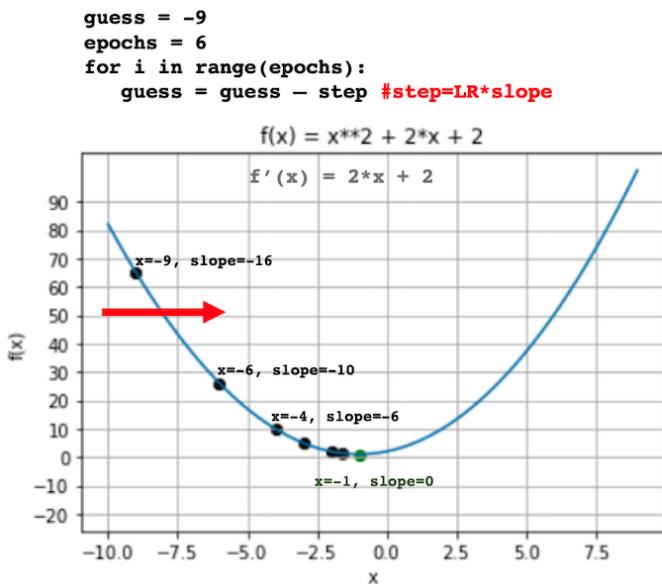


100\$ Questions

1. From where to start?
2. Which direction to move?
3. Size of step to take?
4. When to stop?

b. What is Gradient Descent and How to do Gradient Descent?

Gradient Descent is a first-order iterative optimization algorithm for finding a local minimum of a differentiable function. The idea is to take repeated steps in the opposite direction of the gradient of the function at the current point.



c. Gradient Descent (1-D) in Python

Example Function:

$$f(x) = 3x^2 - 3x + 4$$

$$f'(x) = 6x - 3$$

$$6x - 3 = 0 \implies x = \frac{1}{2}$$

In [25]:

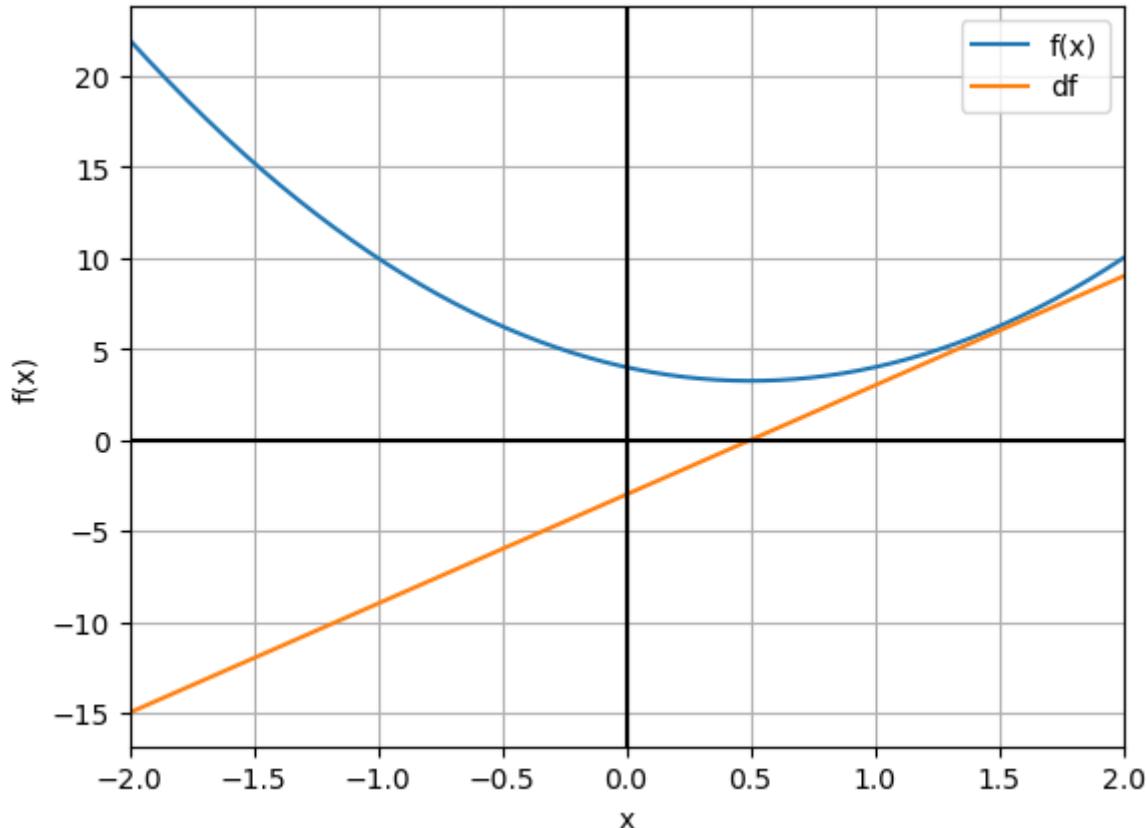
```
1 import numpy as np
2 import matplotlib.pyplot as plt
```

In [26]:

```
1 # function (as a function)
2 def fx(x):
3     return 3*x**2 - 3*x + 4
4
5 # derivative function
6 def deriv(x):
7     return 6*x - 3
```

In [27]:

```
1 %matplotlib inline
2 fig, ax = plt.subplots(nrows=1, ncols=1)
3 x = np.linspace(-2,2,2000)
4 plt.plot(x, fx(x))
5 plt.plot(x, deriv(x))
6 plt.xlim(x[[0,-1]])
7 plt.grid()
8 plt.axvline(x=0, color='black')
9 plt.axhline(y=0, color='black')
10 plt.xlabel('x')
11 plt.ylabel('f(x)')
12 plt.legend(["f(x)", "df"])
13 plt.show()
```



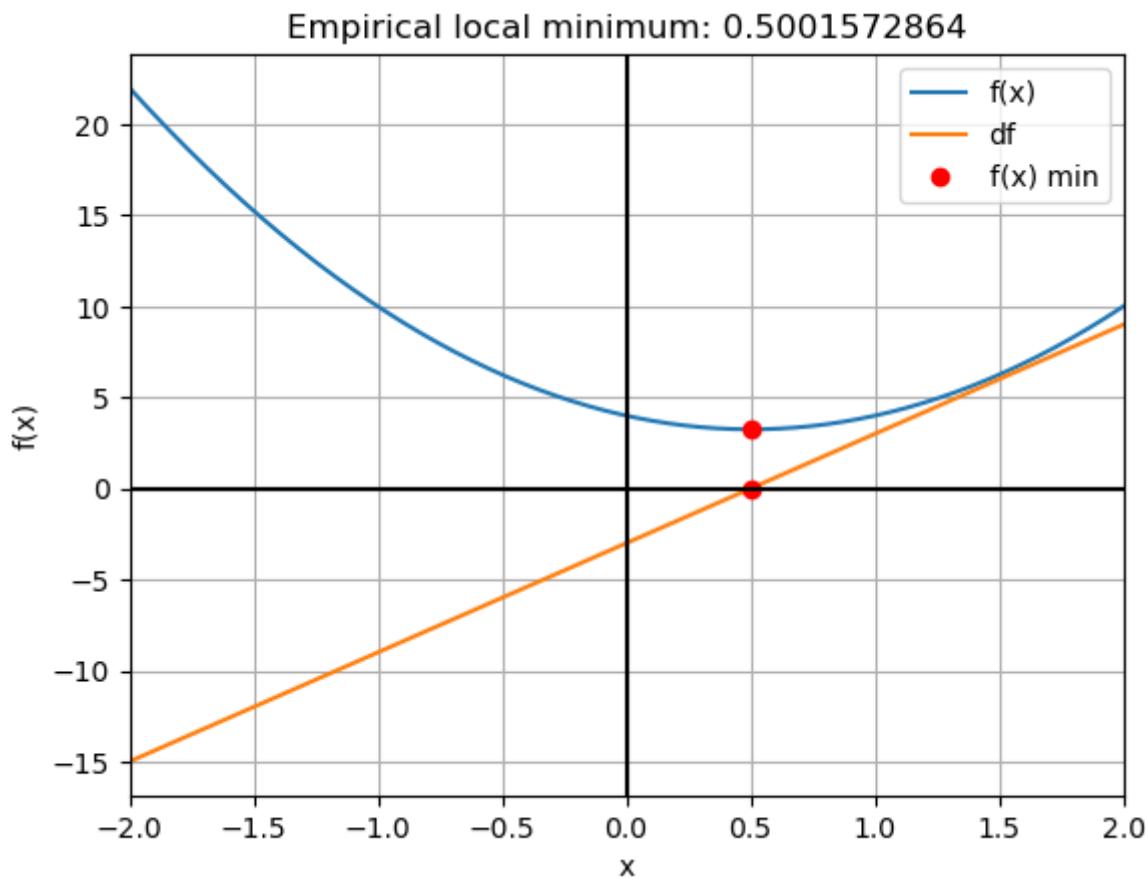
In [28]:

```
1 guess_min = 2          # See impact of changing this parameter of GD
2 print("Initial Guess: ", guess_min)
3 # learning parameters
4 alpha = 0.1           # See impact of changing this parameter of GD
5 epochs = 10            # See impact of changing this parameter of GD
6
7 # run through training
8 for i in range(epochs):
9     guess_min = guess_min - alpha * deriv(guess_min)
10
11 localmin = guess_min
12 print("Local Minimum: ", localmin)
```

```
Initial Guess:  2
Local Minimum:  0.5001572864
```

In [29]:

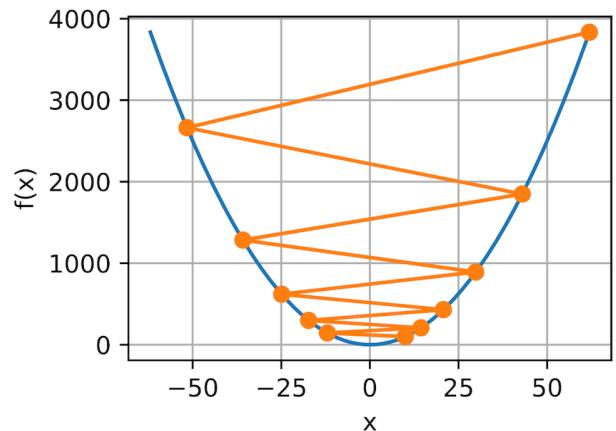
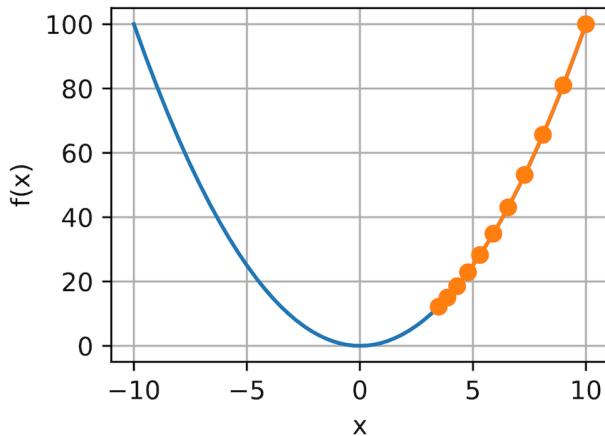
```
1 %matplotlib inline
2 fig, ax = plt.subplots(nrows=1, ncols=1)
3 plt.plot(x,fx(x), x,deriv(x))
4 plt.plot(localmin,deriv(localmin), 'ro')
5 plt.plot(localmin,fx(localmin), 'ro')
6
7 plt.xlim(x[[0,-1]])
8 plt.grid()
9 plt.axvline(x=0, color='black')
10 plt.axhline(y=0, color='black')
11 plt.xlabel('x')
12 plt.ylabel('f(x)')
13 plt.legend(['f(x)', 'df', 'f(x) min'])
14 plt.title('Empirical local minimum: %s' % localmin)
15 plt.show()
```



d. Challenges of Gradient Descent Algorithm

- Gradient Descent algorithm is not guaranteed to give us the correct solution. There are several circumstances that can make Gradient Descent algorithm either go completely wrong or just a bit wrong:
 - Poor choice of model parameters
 - Stuck in a local minimum due to unfortunate starting point
 - Vanishing Gradient
 - Exploding Gradient

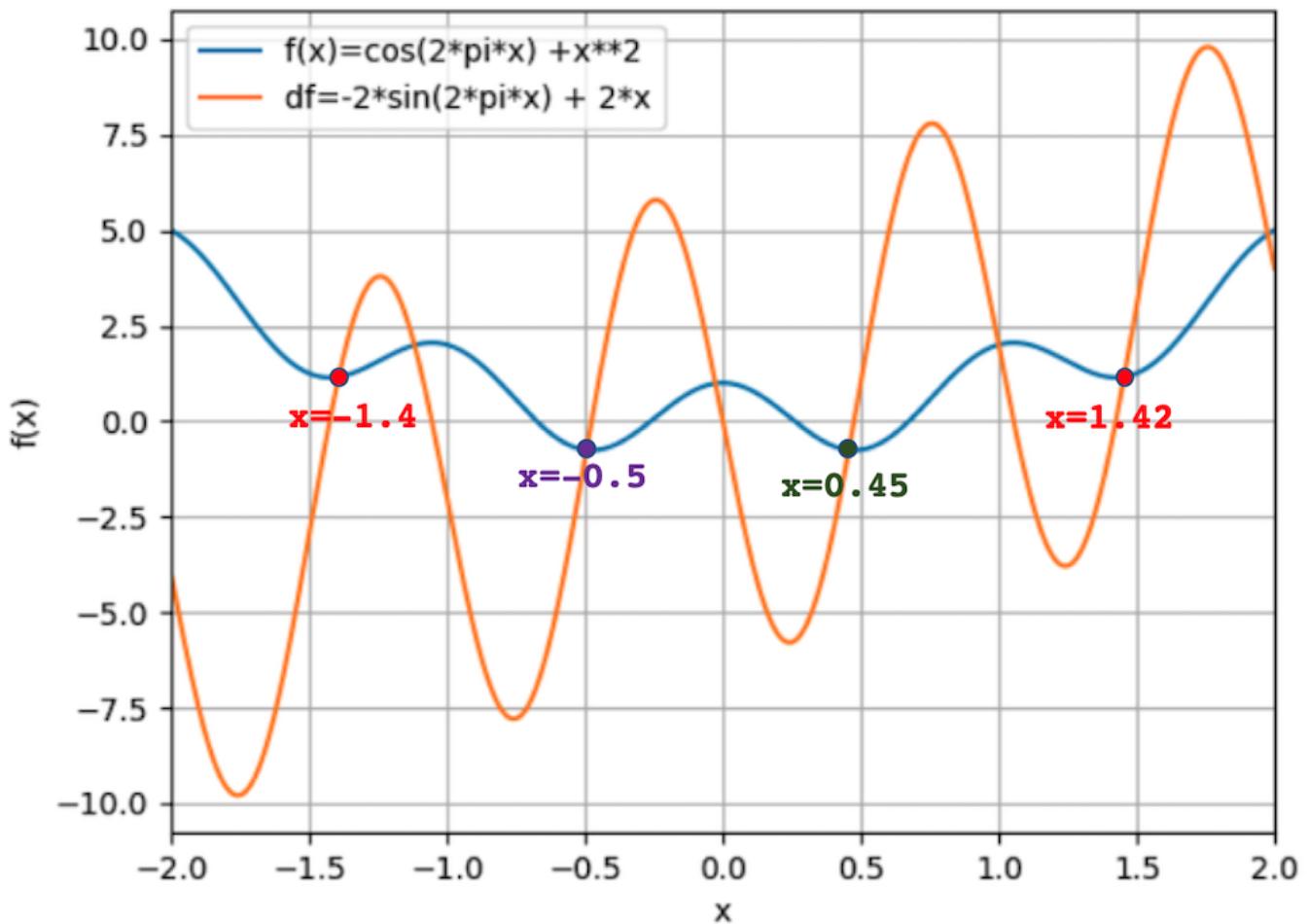
(i) Small vs Large Learning Rate



(ii) Local Minima and Unfortunate Starting Point

$$f(x) = \cos(2\pi x) + x^2$$

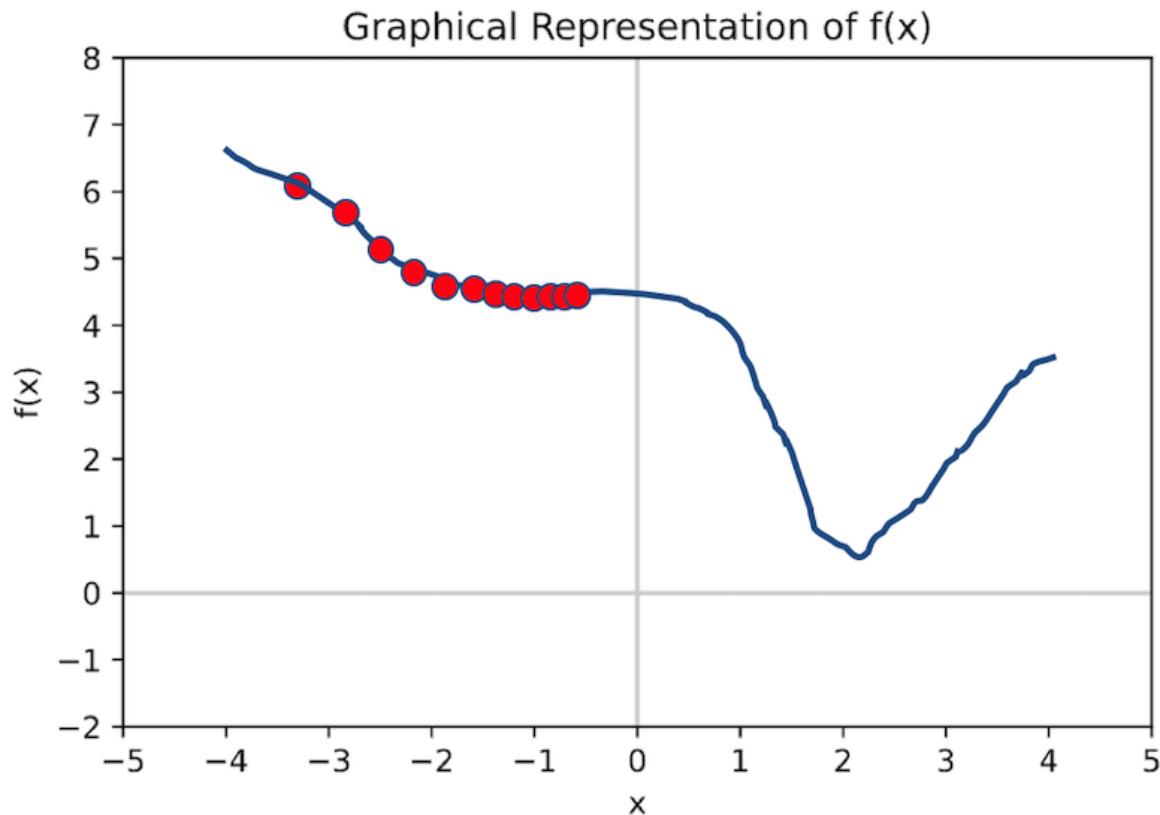
$$f'(x) = -2\pi \sin(2\pi x) + 2x$$



- There are various methods that help in overcoming the local minimum and unfortunate starting point:
 - Use stochastic Gradient Descent

(iii) Vanishing Gradient Problem

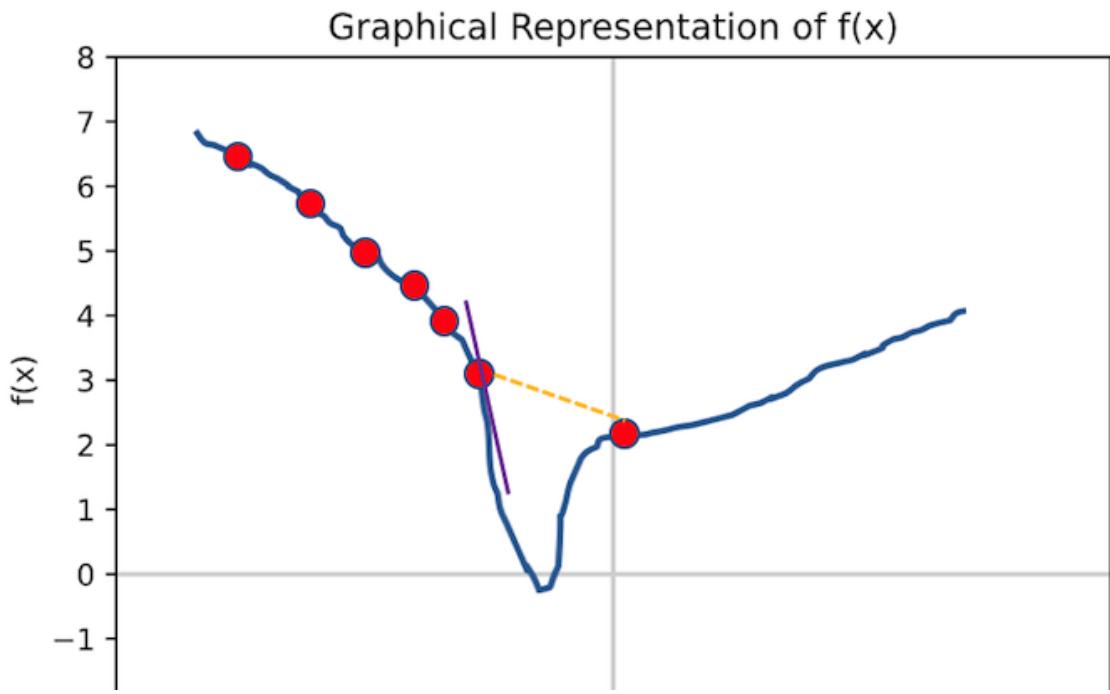
Vanishing gradient descent is a problem that occurs when the gradient become exponentially small and when multiplied with learning rate the step size becomes so much small that you do not reach the function minimum in specified epochs.



- In Deep Learning, there are various techniques that help in overcoming the vanishing gradient problems:
 - Instead of sigmoid activation function use activation functions that do not saturate like Rectified Linear Unit (ReLU, Leaky-ReLU, Randomized ReLU)
 - Long Short Term Memory (LSTM)
 - Gated Recurrent Unit (GRU)

(iv) Exploding Gradient Problem

Exploding gradient is a problem that occurs when the gradient is too steep/large and when multiplied with learning rate the step size becomes so much large that you may miss the minimum and cross the other side.



3. Linear Regression Using Gradient Descent Algorithm

a. Recap of Linear Regression using OLS

- Cost Function:

$$E(\beta_0, \beta_1, \beta_2, \dots, \beta_m) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i + \dots + \beta_m x_{i,m}))^2$$

- Regression Coefficients for Simple Linear Regression:

$$\beta_1 = \frac{(\bar{x}_i y_i) - \bar{x}_i \bar{y}_i}{\bar{x}_i^2 - \bar{x}_i^2} \quad \beta_0 = \bar{y}_i - \beta_1 \bar{x}_i$$

- Regression Coefficients for Multiple Linear Regression:

$$\boldsymbol{\beta} = (X^T X)^{-1} X^T Y$$

- Limitation of using Ordinary Least Squares:

- Taking the derivative of cost function and then solving it for zero to get the set of β coefficients becomes computationally expensive when the number of input features are in thousands or even millions.
- The Gauss Jordan Elimination technique of computing the inverse of a $n \times n$ matrix is

4. Example (Simple Linear Regression using Gradient Descent)

- Consider that a teacher has collected a sample data of seven students, their GPA and the number of hours they have studied on daily basis in the entire semester.
- Suppose, a teacher want to determine the relationship between the two variables `GPA` and `study_hours` is positive or negative.
- The teacher also wants to determine how much impact the dependent variable `study_hours` has on the independent variable `GPA`
- For this she has collected a sample data of seven students as tuple containing (`daily_study_hours` and `acquired_GPA`):

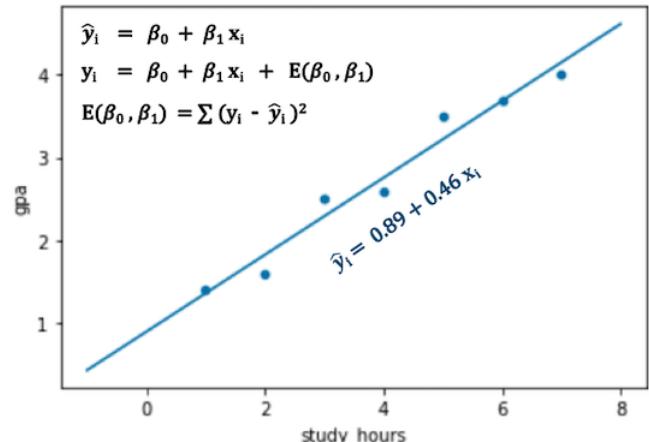
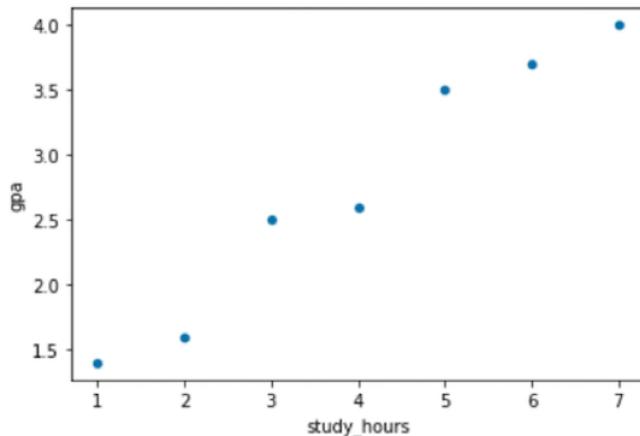
(1, 1.4), (2, 1.6), (3, 2.5), (4, 2.6), (5, 3.5), (6, 3.7), (7, 4.0)

In [30]:

```
1 import pandas as pd
2 df = pd.read_csv("datasets/study-hours.csv")
3 df
```

Out[30]:

	study_hours	gpa
0	1.0	1.4
1	2.0	1.6
2	3.0	2.5
3	4.0	2.6
4	5.0	3.5
5	6.0	3.7
6	7.0	4.0



(i) Minimizing Error Function for Gradient Descent Algorithm

Partial Derivative w.r.t β_0 (y-intercept):

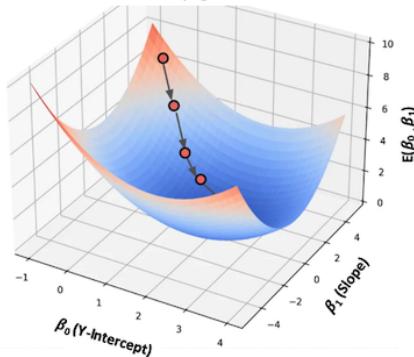
$$\frac{\partial E}{\partial \beta_0} = \frac{\partial}{\partial \beta_0} \frac{1}{2n} \sum (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial E}{\partial \beta_0} = \frac{2}{2n} \sum (y_i - \beta_0 - \beta_1 x_i)^2 (-1)$$

$$\frac{\partial E}{\partial \beta_0} = -\frac{1}{n} \sum (y_i - (\beta_0 + \beta_1 x_i))$$

$$\frac{\partial E}{\partial \beta_0} = -\frac{1}{n} \sum (y_i - \hat{y}_i)$$

$$E(\beta_0, \beta_1) = \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$



Partial Derivative w.r.t β_1 (slope):

$$\frac{\partial E}{\partial \beta_1} = \frac{\partial}{\partial \beta_1} \frac{1}{2n} \sum (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial E}{\partial \beta_1} = \frac{2}{2n} \sum (y_i - \beta_0 - \beta_1 x_i)^2 (-x_i)$$

$$\frac{\partial E}{\partial \beta_1} = -\frac{1}{n} \sum x_i (y_i - (\beta_0 + \beta_1 x_i))$$

$$\frac{\partial E}{\partial \beta_1} = -\frac{1}{n} \sum x_i (y_i - \hat{y}_i)$$

$$\frac{\partial E}{\partial \beta_m} = -\frac{1}{n} \sum x_{im} (y_i - \hat{y}_i)$$

$$\frac{\partial E}{\partial \beta_0} = -\frac{1}{n} \sum (y_i - \hat{y}_i)$$

$$\frac{\partial E}{\partial \beta_1} = -\frac{1}{n} \sum x_{i1} (y_i - \hat{y}_i)$$

$$\frac{\partial E}{\partial \beta_2} = -\frac{1}{n} \sum x_{i2} (y_i - \hat{y}_i)$$

$$\frac{\partial E}{\partial \beta_3} = -\frac{1}{n} \sum x_{i3} (y_i - \hat{y}_i)$$

(ii) Calculate Regression Coefficients Using GD

In [31]:

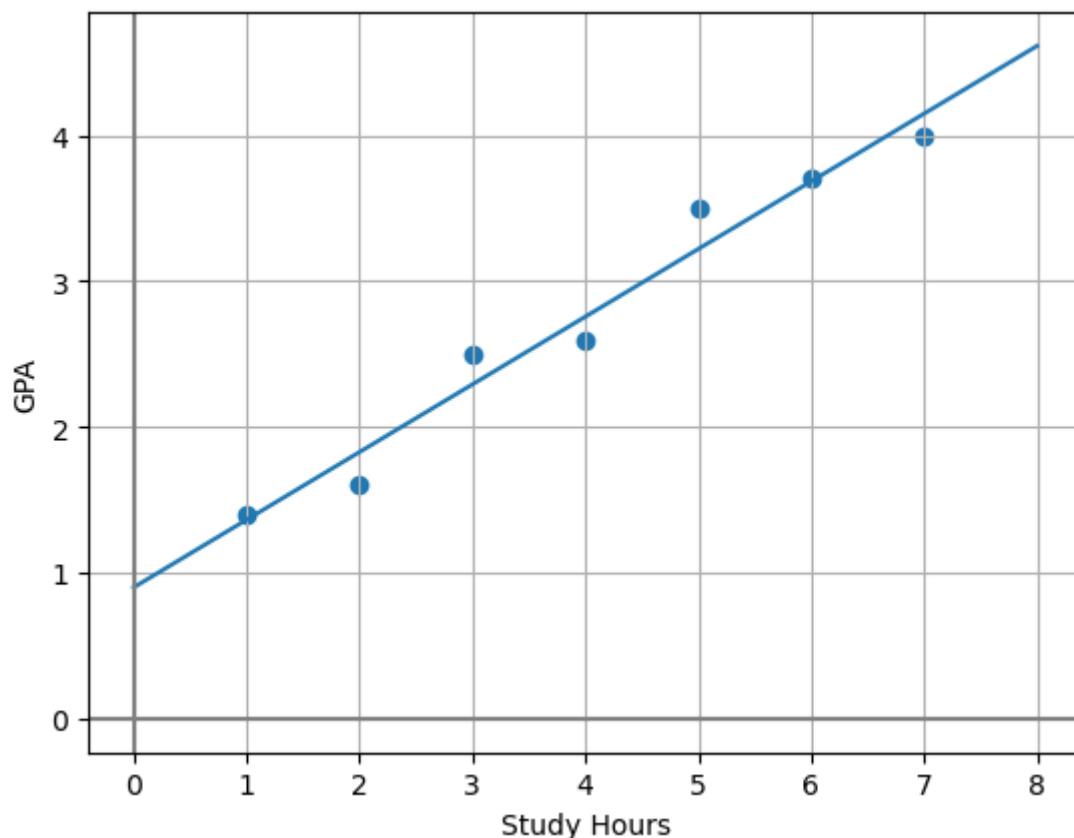
```
1 x = df['study_hours'].values
2 y = df['gpa'].values
3 n = df.shape[0]           # number of training examples
4
5 # Step 1: Start with some random values of y-intercept and coefficient(s) (Model)
6 beta0 = 0.0
7 beta1 = 1.0
8
9 # Step 2: Set appropriate values for Learning Rate and Training Epochs (Hyperparameters)
10 alpha = 0.01    # learning rate
11 epochs = 10000  # number of iteration/epochs
12
13 # Step 3: Update beta0 and betal in epoch iterations (Training)
14 for i in range(epochs):
15     yhat = beta0 + beta1*x      #calculate yhat for entire dataset(predictions)
16
17     # Calculate the gradients of the cost function
18     d_beta0 = (-1/n) * sum((y - yhat))          # wrt beta0
19     d_beta1 = (-1/n) * sum(x * (y - yhat))      # wrt beta1
20
21     # update the gradients
22     beta0 = beta0 - alpha * d_beta0    # Update beta0
23     beta1 = beta1 - alpha * d_beta1    # Update beta1
24
25
26 print("beta0 = {:.5f}\nbeta1 = {:.5f}\n".format(beta0, beta1))
```

```
beta0 = 0.90000
beta1 = 0.46429
```

(iii) Fit the Line

In [32]:

```
1 import numpy as np
2 from matplotlib import pyplot as plt
3 %matplotlib inline
4
5 x = np.array(df['study_hours'])
6 y = np.array(df['gpa'])
7
8 fig, ax = plt.subplots(nrows=1, ncols=1)
9 plt.scatter(x,y)
10
11 xpointsofline = np.linspace(0, 8, 1000)
12 ypointsofline = beta0 + beta1 * xpointsofline
13 plt.plot(xpointsofline, ypointsofline)
14
15 plt.grid()
16 plt.axvline(x=0, color='gray')
17 plt.axhline(y=0, color='gray')
18 plt.xlabel('Study Hours')
19 plt.ylabel('GPA')
20 plt.show()
```



(iv) Carry out the Prediction:

A student has studied for 4.6 hours per day in the entire semester. Can you predict his/her GPA?

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

In [33]:

```
1 sh = 4.6
2 gpa = beta0 + betal * sh
3 print("If a student studies 4.6 hours daily, his/her predicted GPA is: {:.3f}".f
```

If a student studies 4.6 hours daily, his/her predicted GPA is: 3.036

(v) A Deep Dive: Calculate Regression Coefficients Using GD

Calculate Partial Derivatives of Error Function

$$(1, 1.4), (2, 1.6), (3, 2.5), (4, 2.6), (5, 3.5), (6, 3.7), (7, 4.0)$$

Derivative of Error Function w.r.t β_0 :

$$\frac{\partial E}{\partial \beta_0} = -\frac{1}{n} \sum (y_i - \beta_0 - \beta_1 x_i)$$

$$\frac{\partial E}{\partial \beta_0} = -\frac{1}{7} [(1.4 - \beta_0 - \beta_1) + (1.6 - \beta_0 - 2\beta_1) + (2.5 - \beta_0 - 3\beta_1) + (2.6 - \beta_0 - 4\beta_1) + (3.5 - \beta_0 -$$

$$\frac{\partial E}{\partial \beta_0} = \beta_0 + 4\beta_1 - 2.758 \quad \dots (i)$$

Derivative of Error Function w.r.t β_1 :

$$\frac{\partial E}{\partial \beta_1} = -\frac{1}{n} \sum x_i (y_i - \beta_0 - \beta_1 x_i)$$

$$\frac{\partial E}{\partial \beta_1} = -\frac{1}{7} [(1.4 - \beta_0 - \beta_1) + 2 * (1.6 - \beta_0 - 2\beta_1) + 3 * (2.5 - \beta_0 - 3\beta_1) + 4 * (2.6 - \beta_0 - 4\beta_1) +$$

$$\frac{\partial E}{\partial \beta_1} = \beta_0 + 5\beta_1 - 3.223 \quad \dots (ii)$$

Python Functions to Calculate Partial Derivative at specific values:

In [34]:

```
1 def deriv_beta0(beta0,beta1):
2     return beta0 + 4*beta1 - 2.758
3
4 def deriv_beta1(beta0,beta1):
5     return beta0 + 5*beta1 - 3.223
```

- Here to keep things simple, I have hard coded every thing.
 - Please write your own function that is passed the error function name, the index of the argument/variable w.r.t which you want to differentiate and the values of β_0 and β_1 at which you want to calculate the slope. (This we have done above)

Another Implementation of GD Algorithm:

In [35]:

```

1 # Step 1: Start with some random values of y-intercept and coefficient(s) (Model)
2 beta0 = 0.0
3 beta1 = 1.0
4
5 # Step 2: Set appropriate values for Learning Rate and Training Epochs (Hyperparameters)
6 alpha = 0.1      # learning rate
7 epochs = 500     # number of iteration/epochs
8
9 # Step 3: Update beta0 and beta1 in epoch iterations (Training)
10 print('Guess:    beta0_grad {:.3}, beta1_grad {:.3}'.format(beta0, beta1))
11 for i in range(epochs):
12     beta0 = beta0 - alpha * deriv_beta0(beta0, beta1)
13     beta1 = beta1 - alpha * deriv_beta1(beta0, beta1)
14
15     print('Epoch {}, beta0_grad {:.3}, beta1_grad {:.3}'.format(i, beta0, beta1))

```

```
Guess: beta0_grad 0.0, betal_grad 1.0
Epoch 0, beta0_grad -0.124, betal_grad 0.835
Epoch 1, beta0_grad -0.17, betal_grad 0.757
Epoch 2, beta0_grad -0.18, betal_grad 0.719
Epoch 3, beta0_grad -0.173, betal_grad 0.699
Epoch 4, beta0_grad -0.16, betal_grad 0.688
Epoch 5, beta0_grad -0.143, betal_grad 0.68
Epoch 6, beta0_grad -0.125, betal_grad 0.675
Epoch 7, beta0_grad -0.107, betal_grad 0.671
Epoch 8, beta0_grad -0.0887, betal_grad 0.666
Epoch 9, beta0_grad -0.0706, betal_grad 0.663
Epoch 10, beta0_grad -0.0527, betal_grad 0.659
Epoch 11, beta0_grad -0.0352, betal_grad 0.655
Epoch 12, beta0_grad -0.018, betal_grad 0.652
Epoch 13, beta0_grad -0.00108, betal_grad 0.648
Epoch 14, beta0_grad 0.0155, betal_grad 0.645
Epoch 15, beta0_grad 0.0318, betal_grad 0.642
Epoch 16, beta0_grad 0.0478, betal_grad 0.638
Epoch 17, beta0_grad 0.0635, betal_grad 0.635
```

$$\beta_0 = 0.89 \quad \beta_1 = 0.46$$

5. Example (Multiple Linear Regression using Gradient Descent)

Dataset Description

- `number_courses` : Number of Courses Opted by the student
- `time_study` : Average Time Studied per day by the student
- `Marks` : Marks Obtained by the student

In [36]:

```

1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import warnings
6 warnings.filterwarnings('ignore')
7
8 df = pd.read_csv('datasets/Student_Marks.csv')
9 df

```

Out[36]:

	number_courses	time_study	Marks
0	3	4.508	19.202
1	4	0.096	7.734
2	4	3.133	13.811
3	6	7.909	53.018
4	8	7.811	55.299
...
95	6	3.561	19.128
96	3	0.301	5.609
97	4	7.163	41.444
98	7	0.309	12.027
99	3	6.335	32.357

100 rows × 3 columns

a. Batch Gradient Descent

- Batch gradient descent (SGD) calculates the gradient using a the entire dataset in each iteration, therefore, requires the entire training dataset to be in memory and available to the algorithm.

$$\beta_{t+1} = \beta_t - \alpha \nabla f(\beta_t)$$

where

- β_t is the weight vector at time t
- α is the learning rate
- $\nabla f(\beta_t)$ is the gradient of the loss function f with respect to β_t .

In [37]:

```
1 x1 = df.number_courses.values # first input feature
2 x2 = df.time_study.values    # second input feature
3 y = df.Marks.values         # output feature
4
5 # Step 1: Start with some random values of y-intercept and coefficients (Model Parameters)
6 beta0, beta1, beta2 = 0, 0, 0      # weights initialization
7
8 # Step 2: Set appropriate values for Learning Rate and Training Epochs (Hyperparameters)
9 alpha = 0.01                      # learning rate
10 epochs = 10000                     # number of iteration/epochs
11 n = df.shape[0]                   # number of training examples
12
13 # Step 3: Update betas using Batch/Vanilla Gradient Descent Algorithm (Training)
14
15 for i in range(epochs):
16     yhat = beta0 + beta1*x1 + beta2*x2           #calculate yhat (predictions)
17
18     # Calculate the gradients of the cost function
19     d_beta0 = (-1/n) * sum((y - yhat))           # wrt beta0
20     d_beta1 = (-1/n) * sum(x1 * (y - yhat))       # wrt beta1
21     d_beta2 = (-1/n) * sum(x2 * (y - yhat))       # wrt beta2
22
23     # update the gradients
24     beta0 = beta0 - alpha * d_beta0   # Update beta0
25     beta1 = beta1 - alpha * d_beta1   # Update beta1
26     beta2 = beta2 - alpha * d_beta2   # Update beta2
27
28
29 print("beta0 = {:.5f}\nbeta1 = {:.5f}\nbeta2 = {:.5f} \n".format(beta0, beta1, beta2))
```

beta0 = -7.45526
beta1 = 1.86390
beta2 = 5.39913

b. Stochastic Gradient Descent

- Stochastic gradient descent (SGD) calculates the gradient using a single sample at a time.
- In contrast to Batch GD, Stochastic GD updates the model parameters in each epoch by considering one record at a time. The advantage is it may converge faster and is not memory hungry like Batch GD

$$\beta_{t+1} = \beta_t - \alpha \nabla f_i(\beta_t)$$

where,

- β_t is the weight vector at time t
- α is the learning rate
- i is a randomly chosen index from the training set
- $\nabla f_i(\beta_t)$ is the gradient of the loss function f with respect to β_t for the training example indexed by i

In [38]:

```
1 x1 = df.number_courses.values # first input feature
2 x2 = df.time_study.values    # second input feature
3 y = df.Marks.values         # output feature
4
5 # Step 1: Start with some random values of y-intercept and coefficients (Model 1)
6 beta0, beta1, beta2 = 0, 0, 0      # weights initialization
7
8 # Step 2: Set appropriate values for Learning Rate and Training Epochs (Hyperparameters)
9 alpha = 0.01                      # learning rate
10 epochs = 10000                     # number of iteration/epochs
11
12 n = df.shape[0]                   # number of training examples
13
14 # Step 3: Update betas using Stochastic Gradient Descent Algorithm (Training)
15 for i in range(epochs):
16     # Select a random index from the training examples
17     idx = np.random.randint(0, n)
18     # Get the corresponding input values
19     x1_i, x2_i, y_i = x1[idx], x2[idx], y[idx]
20     # Get the prediction for the current input using the current weights
21     yhat_i = beta0 + beta1*x1_i + beta2*x2_i           #calculate yhat (predictions)
22
23     # Calculate the gradients of the cost function
24     d_beta0 = (-1) * (y_i - yhat_i)          # wrt beta0
25     d_beta1 = (-1) * x1_i * (y_i - yhat_i)    # wrt beta1
26     d_beta2 = (-1) * x2_i * (y_i - yhat_i)    # wrt beta2
27
28     # update the gradients
29     beta0 = beta0 - alpha * d_beta0        # Update beta0
30     beta1 = beta1 - alpha * d_beta1        # Update beta1
31     beta2 = beta2 - alpha * d_beta2        # Update beta2
32
33
34 print("beta0 = {:.5f}\nbeta1 = {:.5f}\nbeta2 = {:.5f} \n".format(beta0, beta1, beta2))
```

beta0 = -7.77313
beta1 = 1.56886
beta2 = 5.61978

c. Mini-batch Gradient Descent

- Suppose you have 100 observations you make four batches of 25 observations each and in turn you get four betas. Finally you select the best betas for your model.
- Mini-batch gradient descent is a combination of the concepts of SGD and batch gradient descent. It simply splits the training dataset into small batches and performs an update for each of those batches. This creates a balance between the robustness of stochastic gradient descent and the efficiency of batch gradient descent.

$$\beta_{t+1} = \beta_t - \alpha \nabla f_{B_t}(\beta_t)$$

where,

- β_t is the weight vector at time t
- α is the learning rate
- B_t is a randomly chosen mini-batch of examples from the training set

- $\nabla f_{B_t}(\beta_t)$ is the average gradient of the loss function f with respect to β_t for the examples in the mini-batch B_t

In [39]:

```

1 x1 = df.number_courses.values # first input feature
2 x2 = df.time_study.values    # second input feature
3 y = df.Marks.values         # output feature
4
5 # Step 1: Start with some random values of y-intercept and coefficients
6 beta0, beta1, beta2 = 0, 0, 0 # weights initialization
7
8 # Step 2: Set appropriate values for Learning Rate and Training Epochs
9 alpha = 0.01      # learning rate
10 epochs = 10000    # number of iteration/epochs
11 batch_size = 32   # set of batch size
12
13 n = df.shape[0]   # number of training examples
14
15 # Step 3: Update betas using Mini-Batch Gradient Descent Algorithm
16 for i in range(epochs):
17     # create a batch
18     random_indices = np.random.randint(n, size=batch_size) # select random indices
19     batch_x1 = x1[random_indices]
20     batch_x2 = x2[random_indices]
21     batch_y = y[random_indices]
22
23     # Get the prediction for the current input using the current betas
24     batch_yhat = beta0 + beta1*batch_x1 + beta2*batch_x2           # calculate yhat
25
26     # Calculate the gradients of the cost function
27     d_beta0 = (-1/batch_size) * sum((batch_y - batch_yhat))          # wrt beta0
28     d_beta1 = (-1/batch_size) * sum(batch_x1 * (batch_y - batch_yhat)) # wrt beta1
29     d_beta2 = (-1/batch_size) * sum(batch_x2 * (batch_y - batch_yhat)) # wrt beta2
30
31     # update the gradients
32     beta0 = beta0 - alpha * d_beta0      # Update beta0
33     beta1 = beta1 - alpha * d_beta1      # Update beta1
34     beta2 = beta2 - alpha * d_beta2      # Update beta2
35
36
37 print("beta0 = {:.5f}\nbeta1 = {:.5f}\nbeta2 = {:.5f} \n".format(beta0, beta1, beta2))

```

```

beta0 = -7.43077
beta1 = 1.91948
beta2 = 5.37702

```

d. Using Scikit-Learn's SGDRegressor

In [40]:

```
1 from sklearn.linear_model import SGDRegressor
2 import pandas as pd
3
4 df = pd.read_csv('datasets/Student_Marks.csv')
5 X = df.drop('Marks', axis=1)
6 y = df.Marks.values
7
8 model = SGDRegressor(loss = 'squared_loss', penalty='l2', learning_rate = 'const'
9 model.fit(X,y)
10
11 print("beta0 = {:.5f}".format(model.intercept_[0]))
12 print("beta1 = {:.5f}".format(model.coef_[0]))
13 print("beta2 = {:.5f}".format(model.coef_[1]))
```

```
beta0 = -5.64495
beta1 = 1.49874
beta2 = 5.47643
```

5. Task To Do (Gradient Descent on 4D Dataset)

Use Batch GD

Use Stochastic GD

Use Mini-Batch GD

a. Use Anscombe's Quartet Dataset (Simple Linear Regression)

In [41]:

```
1 import seaborn as sns
2 import numpy as np
3 import pandas as pd
4 df_anscombe = sns.load_dataset('anscombe')
5 df_anscombe.sample(10)
```

Out[41]:

	dataset	x	y
9	I	7.0	4.82
2	I	13.0	7.58
3	I	9.0	8.81
14	II	9.0	8.77
21	II	5.0	4.74
12	II	8.0	8.14
27	III	14.0	8.84
8	I	12.0	10.84
39	IV	8.0	5.25
28	III	6.0	6.08

b. Use Advertising Dataset (Multiple Linear Regression)

In [42]:

```
1 import pandas as pd
2 df = pd.read_csv("datasets/advertising4D.csv")
3 df
```

Out[42]:

	TV	radio	newspaper	sales
0	230.1	37.8	69.2	22.1
1	44.5	39.3	45.1	10.4
2	17.2	45.9	69.3	9.3
3	151.5	41.3	58.5	18.5
4	180.8	10.8	58.4	12.9
...
195	38.2	3.7	13.8	7.6
196	94.2	4.9	8.1	9.7
197	177.0	9.3	6.4	12.8
198	283.6	42.0	66.2	25.5
199	232.1	8.6	8.7	13.4

200 rows × 4 columns

In []:

```
1
```