

NORMS | DISTANCE METRICS

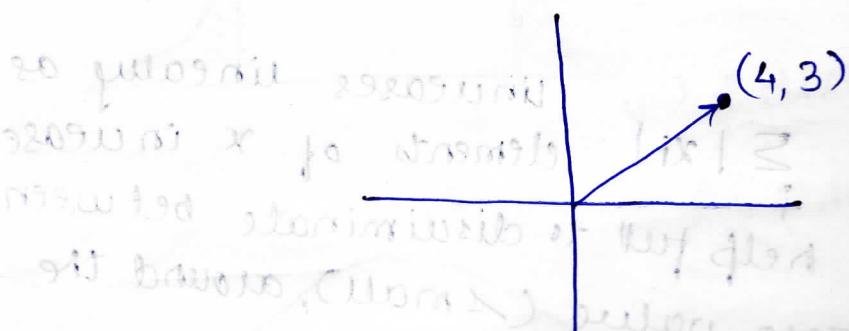
①

Let us say we have 2 variables (x_1, x_2).

i.e $x(x_1, x_2)$

obsr.	x_1	x_2
1	4	3
2		
3		

thus 1st observation can be viewed as a vector with elements {4, 3}.



Thus this vector becomes a vector with x & y coordinates.

length of this vector is the magnitude of the vector. This is called as NORM.

NORM:-

This Norm is used in Data Science to summarize or compact the information about a data vector to a single number. Thus each observation can be represented as a single number.

It is represented as 1×1 .

For the above example, we have used following formula to find the magnitude.

$$\begin{aligned}
 \|x\| &= \sqrt{4^2 + 3^2} \\
 &= \sqrt{x_1^2 + x_2^2} \\
 &= (x_1^2 + x_2^2)^{1/2}.
 \end{aligned}$$

in general

$$\|x\| = \left(\sum_{i=1}^N |x_i|^2 \right)^{1/2}$$

This is called as Euclidean norm (or distance from origin).

This is also called as L_2 norm.

In general; L.P. norm.

$$\|x\|_p = \left(\sum_{i=1}^N |x_i|^p \right)^{1/p}$$

$$= p \sqrt{|x_1|^p + |x_2|^p + \dots + |x_N|^p}$$

thus

when $p = 1$ L_1 norm $\therefore p = 1$

$$\|x\|_1 = \sum_{i=1}^N |x_i|$$

$$\|x\|_1 = |x_1| + |x_2| + \dots + |x_N|$$

Ex: Find L_1 norm for $x = (1, -1, 2)$

$$\begin{aligned}
 \Rightarrow \|x\|_1 &= \sum_{i=1}^3 |x_i| \\
 &= |1| + |-1| + |2| = 4
 \end{aligned}$$

$$\|x\|_1 = 4$$

AAM

when $p=2$ $L_2 \text{ Norm}$

$$\|x\|_2 = \sqrt{\sum_{i=1}^N |x_i|^2}$$
$$= \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

Ex:- find L_2 norm for $x = (1, -1, 2)$

\Rightarrow

$$\|x\|_2 = \sqrt{1^2 + (-1)^2 + (2)^2}$$
$$= \sqrt{6} \approx 2.449$$

$L_2 \text{ Norm} < L_1 \text{ Norm}$

use of Norm in Data science:-

① L_1 Norm :- { Manhattan distance }

We have taken Norm as the magnitude of a vector. It is the distance between the required point & the origin.

Now let us say instead of origin we have another vector i.e. another point $y(y_1, y_2)$. Then applying L_1 norm.

$$\therefore d(x, y) = \|x - y\|_1 = \sum_{i=1}^N |x_i - y_i|$$



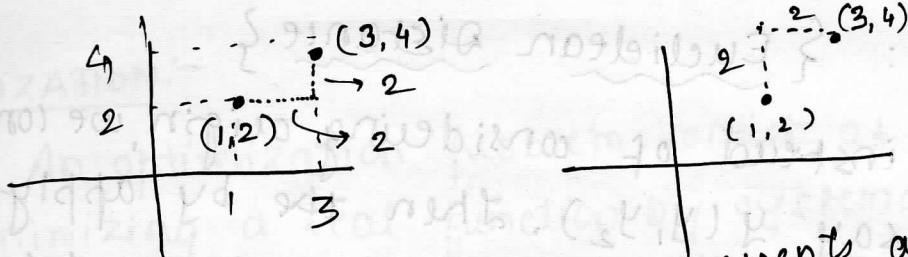
Ex:- if $x(1, 2)$ & $y(3, 4)$ are two data points

then

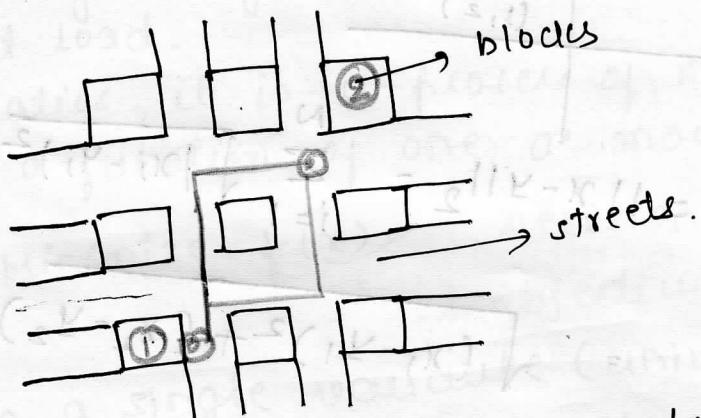
$$d(x, y) = |1-3| + |2-4|$$
$$= 4.$$

This distance which is calculated by L_1 norm is called as Manhattan or city block distance or Taxicab distance.

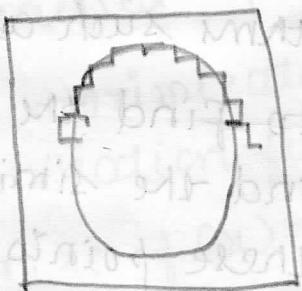
AAM



Thus only horizontal & vertical movements allowed.
 There is a Manhattan city in a suburb in New York which is having construction in the form of blocks (grid) (show ppt for geometric view of Manhattan distance).



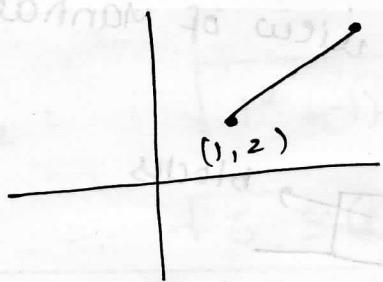
In Machine learning especially in computer vision where images are involved, this Manhattan distance is used. Image is made of pixels which are arranged or stored in the form of matrix. In general when we see any oval boundary, it is actually made up as shown below. Now when we want to go for shape matching etc we have to find the perimeter. In such cases we use Manhattan distance.



② L_2 Norm: { Euclidean distance }

⑤

If instead of considering origin, we consider any pt say $y(y_1, y_2)$. Then by applying L_2 Norm we can find the distance between these two vectors x & y , which is called as Euclidean distance.



$$d(x, y) = \|x - y\|_2 = \sqrt{\sum_{i=1}^N (x_i - y_i)^2}$$

$$= \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$$

Ex:- if $x(1,2)$ & $y(3,4)$ are two data points then

$$d(x, y) = \sqrt{(1-3)^2 + (2-4)^2}$$

$$= \sqrt{8}$$

In Machine learning algorithms such as classification or clustering we need to find the distance between two points to find the similarities & dissimilarities between these points. In such cases we use Euclidean distance to find the distance. e.g KNN. K-nearest neighbour.

Note:- Euclidean distance is the shortest distance i.e less than Manhattan distance. & we need to

③

AAM

OPTIMIZATION:-

An optimization problem consists of maximizing or minimizing a real function by systematically choosing input values from within an allowed set & computing the value of the function. It involves in giving the best solution.

e.g.: choosing stocks which gives maximum returns.

designing bridge that carry maximum possible desired load.

In mathematics, it is a process of maximizing or minimizing a function of one or more variables.

Maximize | minimize $f(x)$

→ decision variable
→ objective function

Here, x is a single variable (univariate optimization).
if x is a vector (multivariate optimization).

If we have full control on x without any restriction, it called as unconstrained optimization.

If we have some restriction on x , it is called

as constrained optimization.

All the maximization problems can be talked as minimization function. i.e.

$\text{Minimize } -f(x)$.

Maximize $f(x) \Rightarrow$ Minimize $-f(x)$.

e.g.: A carpenter makes book cases in two sizes, large & small. It takes 6 hours to make a large book case & 2 hours to make a small one. The profit on a large bookcase is \$50 & the profit on a small bookcase is \$20. The carpenter can spend only 24 hours per week.

7

making bookcases & must make at least 2 of each size per week. Your job as a data scientist is to help your carpenter maximise his revenue.

Now, your initial thought can be the large bookcase will be most profitable.

But the "constraint" is at least 2 of each size per week.

Therefore you can suggest 2 small size (\min^m required). You have total 24 hrs.
 $2 \times \text{small size} = 4 \text{ hrs}$.
 \therefore You are left with 20 hrs. Whereas 3 large size
 $3 \times \text{large size} = 18 \text{ hrs}$

or $3 \times \text{small size} = 6$
 $3 \times \text{large size} = 18$

e.g: IPL fantasy league.

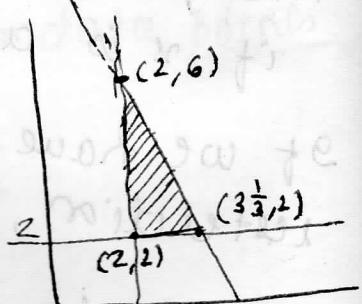
i.e Revenue :-

$$3 \times 20 + 3 \times 50 = \$210$$

Our objective function is

$$P = 50x + 20y$$

if we get a/c carpenter for 2 large bookcase & 6 small bookcases.



constraints are
 $x \geq 2$; $y \geq 2$; $6x + 2y \leq 24$
 large small
 hrs.

then Revenue is $50(2) + 20(6)$

corner-point principle - The maximum & minimum values of P = $\$220$.

obj.fⁿ each occur at one of the vertices of the feasible region

This is the optimum value with constraint.

It is constraint optimization. If there are no constraint like \min^m of at least 2 then it becomes unconstrained optimisation.

AAM

Why optimization for Machine Learning? ⑧
Almost all ML algorithms can be viewed as solutions to optimization problems.

For any optimization problem there are following components:

- (i) objective function (minimization problem)
- (ii) decision variables
- (iii) constraints.

If $f(x)$ is linear and all constraints are also linear, variables are continuous then it is linear programming problem.

If $f(x)$ or constraint are nonlinear then called nonlinear programming problem.

Nonlinear optimization unconstrained univariate.

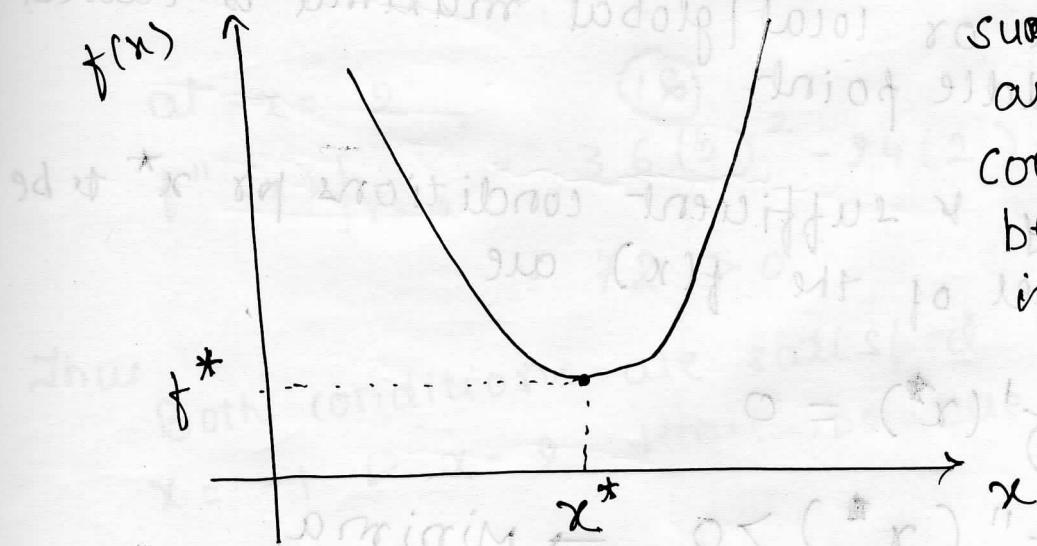
Problem is

$$\min_x f(x)$$

$x \in \mathbb{R}$

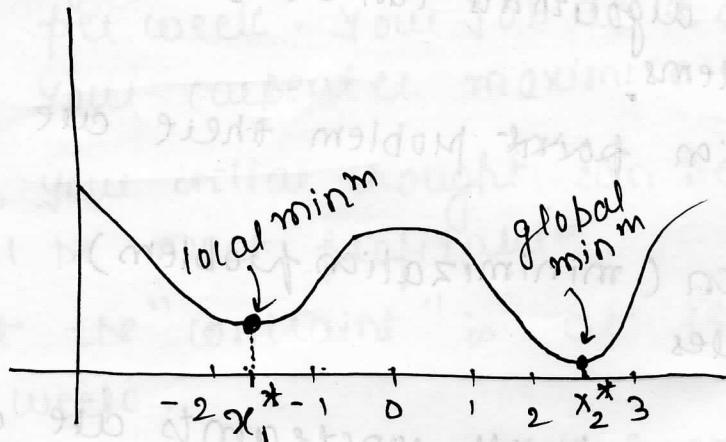
x is continuous any value from real nos.

univariate :- only x is the decision variable let the nonlinear function $f(x)$ is given as follows.



such functions are called as convex functions because there is only one minimum value i.e.

The nonlinear function can be following form as
⑨



such functions are called as non convex functions because there are multiple minima. Job of the optimizer is to find the best solution from these local optimal solutions.

Local minimum \rightarrow In vicinity to this point we cannot find a point better than x^* .

But if you go far away, there would be another local minimum. But this minimum is also the minimum point in the entire region hence called as Global minima.

And a point which is not local/global minima or local/global maxima is called as saddle point. ⑨

Necessary & sufficient conditions for x^* to be minimized of the $f(x)$ are.

$$f'(x^*) = 0$$

& $f''(x^*) > 0 \rightarrow \text{Minima}$

$< 0 \rightarrow \text{Maxima}$
 $= 0 \rightarrow \text{saddle point.}$

⑨

AAM

Ex- $f(x) = 3x^4 - 4x^3 - 12x^2 + 3$. find optimal solun (10)

$$\min_x f(x)$$

Solu:- for minimal value

(i) $f'(x) = 0$

$$f'(x) = 12x^3 - 12x^2 - 24x = 0$$
$$x^3 - x^2 - 2x = 0$$

$$x(x^2 - x - 2) = 0$$

$$x(x-2)(x+1) = 0$$

=) $x = 0, -1, 2$ → possibilities

(ii) $f''(x) = 36x^2 - 24x - 24$

∴ at $x = 0$

$$f''(x)|_{x=0} = -24 < 0$$

at $x = -1 = 36(-1)^2 - 24(-1) - 24$

$$f''(x)|_{x=-1} = 36(-1)^2 - 24(-1) - 24$$
$$= 36 > 0$$

at $x = 2$

$$f''(x)|_{x=2} = 36(2)^2 - 24(2) - 24$$
$$= 72 > 0$$

Thus both conditions are satisfied by values
 $x = -1 \& x = 2$. Hence both are local minima

But at $x = -1$

$$f(-1) = 3(-1)^4 - 4(-1)^3 - 12(-1)^2 + 3$$

$$f(-1) = -2$$

at $x = -1$ (10) AAM

(11)

at $x = 2$

$$f(2) = 3(2)^4 - 4(2)^3 - 12(2)^2 + 3$$

$f(2) = -29$

at $x = 2$ Thus out of $x = -1$ & $x = 2$, we get minimum $f(x)$ at $x = 2$ Hence $x^* = 2$ is the global minima.& $x^* = -1$ is the local minima.HWMinimize $f(x) = x^2 + 3x + 2$.

Non linear unconstrained multivariate optimization:

Let us take minimization problem with two variables.

Minimize $f(x)$; $x = \{x_1, x_2, \dots, x_n\}$ Here we are taking $x = \{x_1, x_2, \dots, x_n\}$ two variables.In univariate optimization we had a condition $f'(x) = 0$ let $Z = f(x)$ $\frac{\partial Z}{\partial x} = f'(x) = 0$ In multivariate it is denoted as 'gradient' ∇f

& is given as

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} \text{ i.e. } \nabla f = 0$$

(11)

AAM

(12)

similarly for univariate $f''(x) > 0$.
on equivalent basis for multivariate optimization with n -variables we use $\nabla^2 f$, given by $n \times n$ Hessian matrix.

This Hessian matrix is given by

$$\nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

condition for multivariate x^* to be minima.

& $\nabla f(\bar{x}^*) = 0$
 $\nabla^2 f(\bar{x}^*) \rightarrow$ positive definite

i.e

Hessian matrix

- $H(x^*) \rightarrow$ positive definite i.e all eigen values are +ve (> 0) then x^* is minima.
- \rightarrow negative definite i.e all eigen values are -ve (< 0) then x^* is maxima
- \rightarrow indefinite i.e both +ve & -ve values then x^* is saddle point
- \rightarrow semidefinite (+ve $\lambda \geq 0$) or -ve $\lambda \leq 0$) then no conclusion

Ex:- find minimum $\Rightarrow f(x) = x_1 + 2x_2 + 4x_1^2 - x_1 x_2$ (B)

Soln:- finding $\nabla f = 0$

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 1 + 8x_1 - x_2 \\ 2 - x_1 + 4x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 8x_1 - x_2 + 1 \\ -x_1 + 4x_2 + 2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} x_1^* \\ x_2^* \end{bmatrix} = \begin{bmatrix} -0.1935 \\ -0.548 \end{bmatrix}$$

Second cond'n:-

$$\nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 8 & -1 \\ -1 & 4 \end{bmatrix}$$

This is Hessian matrix. Now finding the eigen value of Hessian matrix.

$$\begin{bmatrix} 8-\lambda & -1 \\ -1 & 4-\lambda \end{bmatrix} = 32 - 12\lambda + \lambda^2 - 1$$

$$= \lambda^2 - 12\lambda + 31$$

$$\lambda^2 - 12\lambda + 31 = 0$$

$$\Rightarrow \lambda = 3.76, 8.236$$

Both eigen values are positive, hence it is a positive definite matrix hence \bar{x}^* is the required minima

$$\bar{x}^* = \begin{bmatrix} -0.1935 \\ -0.548 \end{bmatrix}$$

(14)

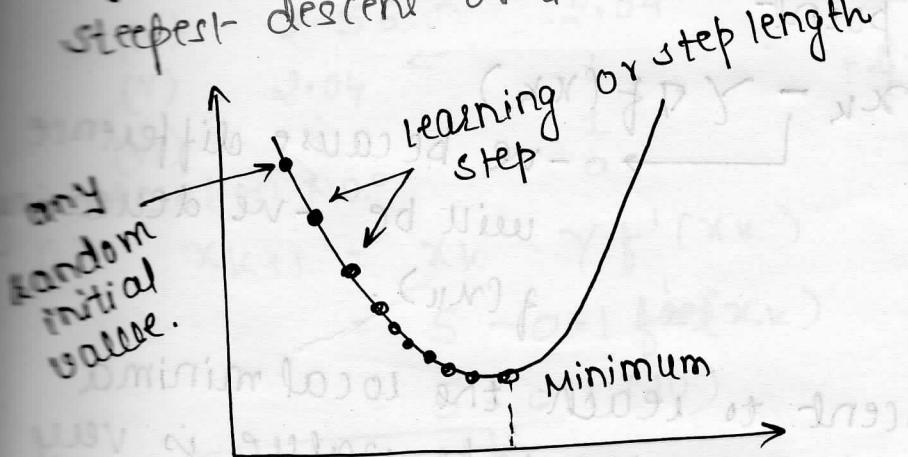
$$f(\bar{x}^*) = -0.1935 + 2(-0.548) + 4(-0.1935)^2 - (-0.1935 \times -0.548) + 2(-0.548)^2$$

$$f(\bar{x}^*) = 7$$

Ans: function $f(x,y) = x^3 + y^3 + 3xy$.

In practice, computing & storing full Hessian matrix takes large memory which is infeasible for high-dimensional functions. For such situations we use Gradient descent algorithm.

steepest descent or Gradient descent:



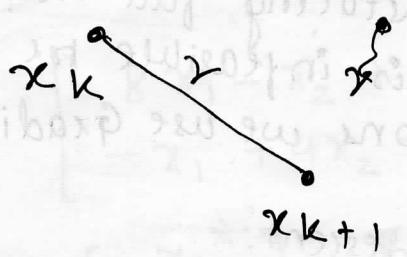
suppose you are at the top of a mountain & want to reach the base camp which is all the way down at the lowest point of the mountain. Due to bad visibility because

of weather you cannot see the path at all. How would you reach the base camp? You can use your feet to know where the land tends to descend. This will give an idea about what direction, slope & the length of your step. If you follow the descending path where you are getting -ve slope (gradient) until you encounter a plain area or an ascending path where the gradient becomes +ve, it is very likely you would reach the base camp.

(14)

AAM

This algorithm is based on the above mentioned concept. It is basically based on the fact that at any given point x_k , the function $f(x)$ decreases fast in the direction of negative gradient & increases in the opposite direction.



If one goes from x_k to x_{k+1} ,
the step length is γ .
 $\nabla f(x_k)$ is the rate
of change of function,

Hence the new point

$$x_{k+1} = x_k - \gamma \nabla f(x_k)$$

-ve because difference
will be -ve decreasing
of $f(x_k)$.

To gradient descent to reach the local minima the selection of this learning rate value is very important. If you select big value as your learning rate, then it may not reach the local minimum because it bounces back & forth between the convex function of gradient descent. If you select very low value of your learning rate then it will take more iteration & hence more time to find the minima.

(15)

AAM

MAA

Ex:- For given univariate function $f(x) = x^2$, apply gradient-descent algorithm.

Solutn:-

$$f(x) = x^2 \quad \text{let } \gamma = 0.1 \text{ & initial point } x_1 = 5.$$

$$f'(x) = 2x$$

By gradient descent algorithm

$$x_{k+1} = x_k - \gamma f'(x_k)$$

Iteration No:- $x_1 = 5 \quad \gamma = 0.1 \quad f'(x_1) = 10.$

1:

$$x_2 = x_1 - 0.1 f'(x_1)$$

$$= 5 - 0.1 (10)$$

$$\boxed{x_2 = 4} \quad \boxed{f(x_2) = 16} \quad f'(x_2) = 8$$

2:

$$x_3 = x_2 - 0.1 f'(x_2)$$

$$= 4 - 0.1 (8)$$

$$\boxed{x_3 = 3.2} \quad \boxed{f(x_3) = 10.24} \quad f'(x_3) = 6.4$$

3:

$$x_4 = x_3 - 0.1 f'(x_3)$$

$$= 3.2 - 0.1 (6.4)$$

$$\boxed{x_4 = 2.56} \quad \boxed{f(x_4) = 6.55} \quad f'(x_4) = 5.12$$

4:

$$x_5 = x_4 - 0.1 f'(x_4)$$

$$= 2.56 - 0.1 (5.12)$$

$$\boxed{x_5 = 2.04} \quad \boxed{f(x_5) = 4.08} \quad f'(x_5) = 4.08$$

5:

$$x_6 = x_5 - 0.1 f'(x_5)$$

$$= 2.04 - 0.1 (4.08)$$

$$\boxed{x_6 = 1.632} \quad \boxed{f(x_6) = 2.66} \quad f'(x_6) = 3.264$$

thus we can see that $f(x)$ is reducing, the iteration continues & we get at $x^* = 0$, $f(x^*) = 0$. Thus the algorithm converges to optimum value 0.

Ex2: For given multivariate function with $x = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$ &

$$f(x) = 4x_1^2 + 3x_1x_2 + 2.5x_2^2 - 5.5x_1 - 4x_2.$$

Apply gradient descent algorithm. ($\gamma = 0.135$)

Solutn:- since it is multivariate, the variable x is a vector $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$.

$$f(x) = 4x_1^2 + 3x_1x_2 + 2.5x_2^2 - 5.5x_1 - 4x_2.$$

$$f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \end{bmatrix}$$

$$f'(x) = \begin{bmatrix} 8x_1 + 3x_2 - 5.5 \\ 3x_1 + 5x_2 - 4 \end{bmatrix}$$

$$\text{let } x_1 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \gamma = 0.135$$

By gradient descent

$$x_{k+1} = x_k - \gamma f'(x_k)$$

$$x_1 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$$f'(x_1) = \begin{bmatrix} 8x_1 + 3x_2 - 5.5 \\ 3x_1 + 5x_2 - 4 \end{bmatrix}$$

$$f'(x_1) = \begin{bmatrix} 16.5 \\ 12 \end{bmatrix}$$

Iteration No

$$1: x_2 = x_1 - 0.135 \begin{bmatrix} 16.5 \\ 12 \end{bmatrix}$$

$$= \begin{bmatrix} 2 \\ 2 \end{bmatrix} - 0.135 \begin{bmatrix} 16.5 \\ 12 \end{bmatrix}$$

$$x_2 = \boxed{\begin{bmatrix} -0.2275 \\ 0.38 \end{bmatrix}}$$

$$f'(x_2) = \begin{bmatrix} 8x_1 (-0.2275) + 3(0.38) - 5.5 \\ 3x_1 (-0.2275) + 5(0.38) - 4 \end{bmatrix}$$

$$f(x_2) = 4(-0.2275)^2 + 3(-0.2275)(0.38) + 2.5(0.38)^2 - 5.5(-0.2275) - 4(0.38)$$

$$= \boxed{\begin{bmatrix} 6.18 \\ -2.76825 \end{bmatrix}}$$

$$f(x_2) = 0.0399$$

AAM

$$2. \quad x_3 = x_2 - 0.135 \begin{bmatrix} 6.18 \\ -2.16825 \end{bmatrix}$$

$$= \begin{bmatrix} -0.2275 \\ 0.38 \end{bmatrix} - 0.135 \begin{bmatrix} -6.18 \\ -2.70825 \end{bmatrix}$$

$$\boxed{x_3 = \begin{bmatrix} 0.6068 \\ 0.7556 \end{bmatrix}}$$

$$f'(x_2) = \begin{bmatrix} 8 \times 0.6068 + 3 \times 0.7556 - 5.5 \\ 3 \times 0.6068 + 5 \times 0.7556 - 4 \end{bmatrix}$$

$$f'(x_2) = \begin{bmatrix} 1.6212 \\ 1.5984 \end{bmatrix}$$

$$f(x_2) = 4(0.6068)^2 + 3(0.6068)(0.7556) \\ + 2.5(1.5984 \times 0.7556)^2 - 5.5(0.6068) \\ - 4(0.7556)$$

$$\boxed{f(x_2) = -2.0841}$$

$$3. \quad x_4 = x_3 - \gamma f'(x_3)$$

$$= \begin{bmatrix} 0.6068 \\ 0.7556 \end{bmatrix} - 0.135 \begin{bmatrix} 0.6068 \\ 0.7556 \end{bmatrix}$$

$$\boxed{x_4 = \begin{bmatrix} 0.3879 \\ 0.5398 \end{bmatrix}}$$

$$\boxed{f(x_4) = -2.3342}$$

And the iteration continues, at every iteration the value of the objective function keeps coming down. continuing further we get

$$x_5 = \begin{bmatrix} 0.4928 \\ 0.5583 \end{bmatrix} \quad & f(x_5) = -2.3675$$

After few more iterations we get optimal solution

$$(\bar{x}^*) = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \quad f(\bar{x}^*) = -2.3750$$

This algorithm converges. Convergence can be decided by the user based on

- (i) vector difference between \bar{x}_k & \bar{x}_{k+1}
- (ii) difference between $f^*(\bar{x}_k)$ & $f^*(\bar{x}_{k+1})$ or
- (iii) $\nabla f(\bar{x}^{k+1})$ & $\nabla f(\bar{x}^k)$. Difference. Typically when optimal solution is obtained these differences become 0.

H.W: Minimize $f(x, y) = x^2 + y^2$; $\gamma = 0.1$ $x = 1$ & $y = 1$