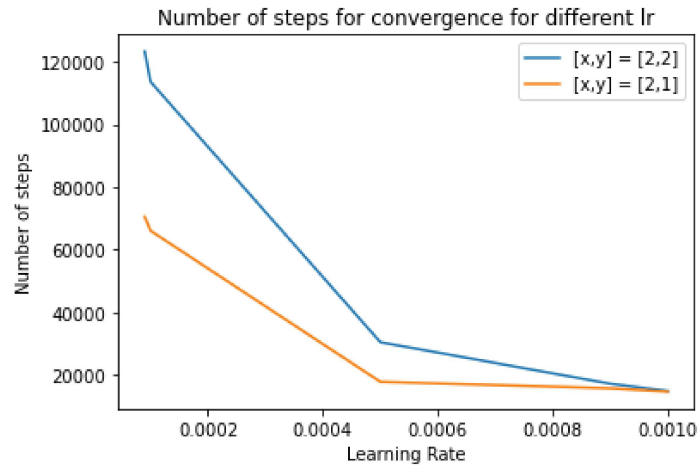


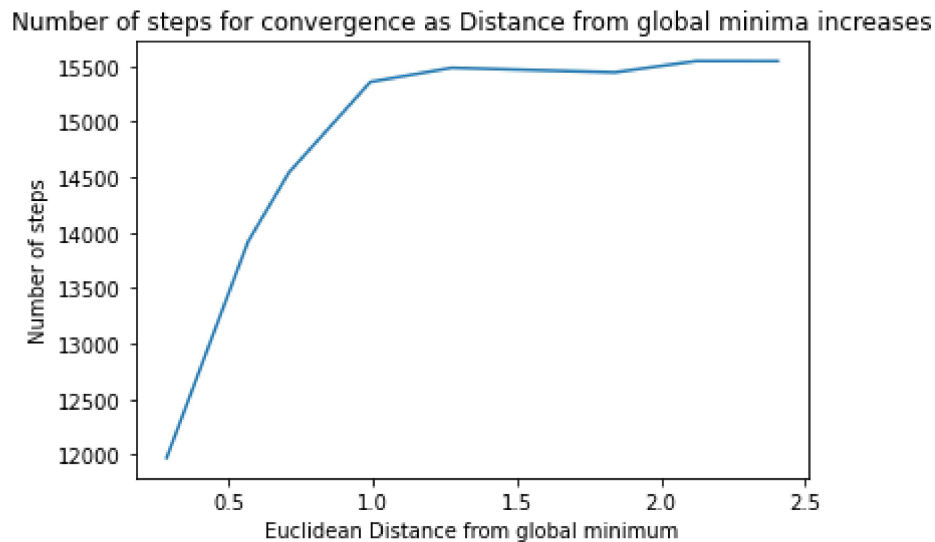
Problem Sheet 3:

Provide a graph which demonstrates the effect of learning rate on convergence of steepest descent algorithm (learning-rate(x) vs num iterations for convergence(y)). What do you think will be the optimal learning rate for this function?



- As we see in the above graph, as the learning rate increases from $2e-4$ to $1e-3$, the number of steps required to convergence is reduced.
- Increasing learning rate further leads to oscillation and it fails to converge.
- Also, if the initial guess is closer to global minimum, the number of steps is lesser.
- Optimal learning rate is $1e-3$

Provide a graph which demonstrates the convergence rate of steepest descent for various initial guesses. (The distance of initial guess from global minimum $\|x - x^*\|$ vs number of iterations y).



- As distance of initial guess from global minimum increases, no of steps to converge increases. Learning rate of $1e-3$ with tolerance $1e-6$ was used for initial guess = $[1.2, 1.2]$, $[1.4, 1.4]$, $[1.5, 1.5]$, $[1.7, 1.7]$, $[1.9, 1.9]$, $[2.3, 2.3]$, $[2.5, 2.5]$, $[2.7, 2.7]$, $[2.9, 2.9]$.

State true or false. The steepest descent method (i.e., gradient descent) requires multiple random starts to ensure the global minimum has reached if the optimization function is convex in nature. (2 Point)

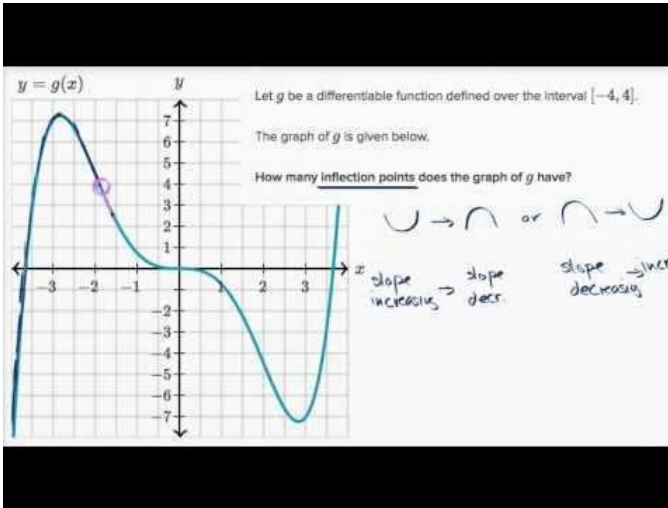
- False: If optimization function is convex, gradient descent will converge with any random initial guess. Convex functions are smooth and gradient direction will always point in direction towards global minima irrespective of starting point.

You are given some function $f(X): \mathbb{R}^n \rightarrow \mathbb{R}, X \in \Omega \subset \mathbb{R}^n$. • If the Hessian matrix $G = \nabla^2 f$ is positive definite, does f have a minimum? Justify your answer. • If G is negative definite, does f have a minimum value? Justify

- If Hessian matrix is positive definite, then function $f(x)$ is a convex function. If a function $f(x)$ is convex then f will have a minimal.
- Similarly, if hessian matrix is negative definite, then $f(x)$ is a concave function and f will have a maximum.

Consider a function which has a stationary point but also has zero Hessian at the stationary point. Can you use the gradient descent method to find the stationary point for such a function?

- Yes, by doing Gradient descent, our function value will decrease until we hit a point where gradient is zero. Since hessian is zero, the stationary point is the point of inflection. So, convergence or identification of stationary points depends on the initial guess. If initial guess is a point above the inflection point, stationary points will be identified.



Give a linear model which best approximates the following data (least square fit / linear regression). Use the analytical formula for linear regression

A simple linear regression model can be given as:

$$y = w_0 + w_1 x$$

$$Y = [1 \ x] [w_0 \ w_1]^T$$

$$Y = XW^T$$

$$\text{Loss} = (Y - XW^T)^T (Y - XW^T)$$

$$\text{Loss} = Y^T Y - 2WX^T Y + WX^T XW^T$$

$$D\text{Loss}/dW = 0 - 2X^T Y + 2WX^T X \Rightarrow 0$$

$$W = (X^T X)^{-1} X^T Y$$

$$\text{Here } X = [1 \ x] = [[1 \ 1], [1 \ 2], [1 \ 3]]$$

$$W = [[1.33 \ 0.33 \ -0.67], [-0.5 \ 0.0 \ 0.5]] * [1.5 \ 1.5 \ 3]$$

$$W = [0.48, 0.75]$$

$$w_0 = 0.48, w_1 = 0.75$$

Q.4 Usually quadratic function is given by
 $A(x) = x^T A x - b x + c$

a)

$$\nabla f(x) = \text{grad} = A x - b \quad \nabla^2 f(x) = A = H$$

$$= A x^{(1)} - b$$

$$= A x^* + A \lambda s - b$$

$$= b + A \lambda s - b$$

$$= \lambda A s$$

$$= \lambda H s$$

$$\boxed{g^{(1)} = \lambda H s}$$

$$\therefore x^{(1)} = x^* + \lambda s$$

$$\therefore A x^* = b \rightarrow \text{grad} = 0$$

$$\Rightarrow H = A$$

Q.4. (b) $f(x_0 + \alpha d_k) = f(x_0) + \alpha d_k^T \nabla f(x_0) + \frac{1}{2} \alpha d_k^T \nabla^2 f(x_0) \alpha d_k$

$d_k = -g$, $\nabla f(x_0) = g$, $\nabla^2 f(x_0) = H$
 $= f(x_0) + \alpha d_k$

$f(x_0 + \alpha d_k) = f(x_0) - \alpha g^T g + \frac{1}{2} \alpha^2 g^T H g$

For $f(x_0)$ is

Function converges if $f(x_0 + \alpha d_k) = f(x_0)$.
 diff. w.r.t. α .

i.e. $0 = -g^T g + \alpha g^T H g$

optimal α is $\boxed{\alpha = \frac{g^T g}{g^T H g}}$

As $x_{k+1} = x_k + \alpha_k d_k$

Similarly

$$\begin{aligned} e_{k+1} &= e_k + \alpha_k d_k \\ &= e_k - \frac{g^T g}{g^T H g} \cdot g \\ &= e_k - \frac{g^T g \cdot (He)}{g^T H e} \\ &= e_k - \frac{g^T g}{g^T g} \lambda e \\ &= e - e \end{aligned}$$

$\boxed{e_{k+1} = 0}$

$e_k \rightarrow$ error
 $(x_k - x^*)$

$$\begin{aligned} g &= Ax - b \\ g &= Ax - A x^* \\ g &= A e \\ g &= \lambda e \end{aligned}$$

Since error is 0, convergence is reached. For one dimensional x , convergence reached in 1 step.