

1.1

Jensen's inequality: Let $f(x)$ be convex in interval I,

Then for any $x_1, \dots, x_3 \in I$

$$f(\text{avg. of } \{x_i\}) \leq \text{avg. of } \{f(x_i)\}$$

and if f were concave, the inequality is reversed i.e.

$$f(\text{avg. of } \{x_i\}) \geq \text{avg. of } \{f(x_i)\}$$

Given inequality to prove,

$$\frac{1}{x-1} + \frac{1}{x} + \frac{1}{x+1} > \frac{3}{x}$$

Let $f(x) = \frac{1}{x}$ which is convex in interval $(0, \infty)$

Let the set $X = \{x-1, x, x+1\} \setminus \{(x-1), x, (x+1)\}$

$$\therefore \text{Avg. of } \{f(x_i)\} = \frac{1}{3} (f(x-1) + f(x) + f(x+1))$$

$$= \frac{1}{3} \left[\frac{1}{x-1} + \frac{1}{x} + \frac{1}{x+1} \right] - \textcircled{1}$$

$$\text{and } f(\text{Avg. of } \{x_i\}) = f\left(\frac{1}{3}(x-1 + x + x+1)\right)$$

$$= f\left(\frac{3x}{3}\right)$$

$$= f(x)$$

$$= \frac{1}{x} - \textcircled{2}$$

From Jensen's inequality, ? $\textcircled{1}, \textcircled{2}$

$$\frac{1}{3} \left(\frac{1}{x-1} + \frac{1}{x} + \frac{1}{x+1} \right) > \frac{1}{x}$$

$$\therefore \frac{1}{x-1} + \frac{1}{x} + \frac{1}{x+1} > \frac{3}{x}$$

\rightarrow Given series, $S = 1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n} + \dots$

S can be written as,

$$S = \left(\frac{1}{2-1} + \frac{1}{2} + \frac{1}{2+1} \right) + \frac{1}{4} + \left(\frac{1}{6-1} + \frac{1}{6} + \frac{1}{6+1} \right) + \frac{1}{8} + \left(\frac{1}{10-1} + \frac{1}{10} + \frac{1}{10+1} \right) + \frac{1}{12} + \dots$$

Using Jensen's inequality for numbers in brackets we get.

$$S > \left(\frac{3}{2} + \frac{3}{6} + \frac{3}{10} + \dots \right) + \left(\frac{1}{4} + \frac{1}{8} + \frac{1}{12} + \dots \right)$$

$$\text{or } S > \frac{3}{2} \left(\frac{1}{2} + \frac{1}{3} + \frac{1}{5} + \frac{1}{7} + \dots \right) + \frac{1}{4} \underbrace{\left(\frac{1}{2} + \frac{1}{3} + \frac{1}{5} + \frac{1}{7} + \dots \right)}_S$$

$$\Rightarrow S > \frac{3}{2} \left(1 + \frac{1}{3} + \frac{1}{5} + \frac{1}{7} + \dots \right) + \frac{S}{4}$$

$$\Rightarrow \frac{3}{4} S > \frac{3}{2} \left(1 + \frac{1}{3} + \frac{1}{5} + \frac{1}{7} + \dots \right) - \quad (1)$$

Now, let $S_1 = 1 + \left(\frac{1}{3}\right) + \left(\frac{1}{5} + \frac{1}{7}\right) + \left(\frac{1}{9} + \frac{1}{11} + \frac{1}{13} + \frac{1}{15}\right) + \dots$

$$S_1 > 1 + \frac{1}{4} + \left(\frac{1}{8} + \frac{1}{8}\right) + \left(\frac{1}{16} + \frac{1}{16} + \frac{1}{16} + \frac{1}{16}\right) + \dots$$

$$= 1 + \frac{1}{4} + \frac{1}{4} + \dots$$

$$\text{Hence } \sum_{n=1}^{\infty} \frac{n}{4}$$

$$\therefore \text{as } n \rightarrow \infty \quad \sum \frac{n}{4} \rightarrow \infty$$

$\therefore S_1 \rightarrow \infty$ and hence S_1 is divergent.

Now from (1),

$$S > \frac{4}{2} S_1$$

$$\text{or } S > 2S_1$$

Since S_1 is divergent which implies S is also divergent.

Therefore the given series does not converge to real number.

1.2 Three Chords Lemma: Let $f: I \rightarrow \mathbb{R}$. Then f is convex on I if and only if, for any points $a, b, c \in I$ with $a < b < c$ we have

$$\frac{f(b) - f(a)}{b-a} \leq \frac{f(c) - f(a)}{c-a} \leq \frac{f(c) - f(b)}{c-b}$$

Proof: Given that f is convex on I and $a < b < c$.

We can write $b = \lambda a + (1-\lambda)c$

$$\text{where } \lambda = \frac{c-b}{c-a} \in [0, 1]$$

Since f is convex,

$$f(b) \leq \left(\frac{c-b}{c-a} \right) f(a) + \left(\frac{b-a}{c-a} \right) f(c) \quad \text{--- (1)}$$

$$\Rightarrow f(b) - f(a) \leq \left(\frac{c-b}{c-a} \right) f(a) + \left(\frac{b-a}{c-a} \right) f(c)$$

$$\Rightarrow f(b) - f(a) \leq \frac{b-a}{c-a} (f(c) - f(a))$$

$$\Rightarrow \frac{f(b) - f(a)}{b-a} \leq \frac{f(c) - f(a)}{c-a} \quad \text{--- (2)}$$

This proves the 1st inequality.

From eq (1) we have

$$f(b) \leq \lambda f(a) + (1-\lambda)f(c)$$

$$\Rightarrow \lambda [f(a) - f(c)] + [f(c) - f(b)] \geq 0$$

$$\Rightarrow \frac{c-b}{c-a} [f(c) - f(a)] \leq f(c) - f(b)$$

$$\Rightarrow \frac{f(c) - f(a)}{c-a} \leq \frac{f(c) - f(b)}{c-b} \quad \text{--- (3)}$$

Hence from eqn (2) and (3) we prove,

$$\frac{f(b) - f(a)}{b-a} \leq \frac{f(c) - f(a)}{c-a} \leq \frac{f(c) - f(b)}{c-b}$$

1.3 Given: x, y, z are positive real numbers with $x+y+z=1$
 To show: $\left(1+\frac{1}{x}\right)\left(1+\frac{1}{y}\right)\left(1+\frac{1}{z}\right) \geq 64$

Soluⁿ: Let $f(x) = \log\left(1+\frac{1}{x}\right)$ which is convex in $(0, \infty)$

Consider average of functions, $f(x), f(y), f(z)$

$$\begin{aligned} &\Rightarrow \frac{1}{3} (f(x) + f(y) + f(z)) \\ &= \frac{1}{3} \left(\log\left(1+\frac{1}{x}\right) + \log\left(1+\frac{1}{y}\right) + \log\left(1+\frac{1}{z}\right) \right) \\ &= \frac{1}{3} \log \left[\left(1+\frac{1}{x}\right) \left(1+\frac{1}{y}\right) \left(1+\frac{1}{z}\right) \right] - \textcircled{1} \end{aligned}$$

Now, consider function f over average of x, y, z
 i.e. $f\left(\frac{x+y+z}{3}\right) = \log\left(1+\frac{3}{x+y+z}\right)$

$$\begin{aligned} &= \log\left(1+\frac{3}{1}\right) \quad [\because x+y+z=1] \\ &= \log 4 - \textcircled{2} \end{aligned}$$

Now, from Jensen's inequality,

$$\frac{1}{3} (f(x) + f(y) + f(z)) \geq f\left(\frac{x+y+z}{3}\right)$$

From $\textcircled{1}$ & $\textcircled{2}$

$$\frac{1}{3} \log \left[\left(1+\frac{1}{x}\right) \left(1+\frac{1}{y}\right) \left(1+\frac{1}{z}\right) \right] \geq \log 4$$

Taking exponentiation on both sides,

$$\frac{1}{3} \left(1+\frac{1}{x}\right) \left(1+\frac{1}{y}\right) \left(1+\frac{1}{z}\right) \geq 64.$$

1.4 Given: $x^3 - 2x - 5 = 0$, $x_0 = 2$, $x_1 = 3$

Solⁿ: The secant method is a root-finding algorithm that uses a succession of roots of secant lines to better approximate root of a function f .

$$x_2 = x_0 - \frac{f(x_0)}{f(x_1) - f(x_0)} (x_1 - x_0)$$

1st iteration: $x_0 = 2$, $x_1 = 3$

$$f(x_0) = f(2) = (2)^3 - 2(2) - 5 = -1$$

$$f(x_1) = f(3) = 16$$

$$\therefore x_2 = x_0 - \frac{f(x_0)}{f(x_1) - f(x_0)} (x_1 - x_0)$$

$$= 2 - \frac{(-1)}{16 - (-1)} (3 - 2)$$

$$x_2 = 2.0588$$

$$\therefore f(x_2) = f(2.0588) = -0.3908$$

2nd iteration: $x_1 = 3$, $x_2 = 2.0588$

$$\therefore x_3 = x_1 - \frac{f(x_1)}{f(x_2) - f(x_1)} (x_2 - x_1)$$

$$= 3 - \frac{16}{-0.3908 - 16} (2.0588 - 3)$$

$$\Rightarrow x_3 = 2.0813$$

$$\therefore f(x_3) = f(2.0813) = -0.1472$$

1.5 Given: $xe^x = 1$, $x \in [0, 1]$

Solⁿ: Bisection method: It is a root-finding method that repeatedly bisects an interval and then selects a subinterval in which a root must lie for further processing.

Here $xe^x - 1 = 0$

Let $f(x) = xe^x - 1$

1st iteration: $f(0)$

$$f(0) = -1 < 0 \quad \text{and} \quad f(1) = 1.7183 > 0$$

∴ Root lies between 0 and 1

$$x_0 = \frac{0+1}{2} = 0.5$$

$$f(x_0) = f(0.5) = 0.5e^{0.5} - 1 = -0.1756 < 0$$

2nd iteration:

$$f(0.5) = -0.1756 < 0 \quad \text{but} \quad f(1) = 1.7183 > 0$$

∴ Root lies between 0.5 and 1

$$x_1 = \frac{0.5+1}{2} = 0.75$$

$$f(x_1) = f(0.75) = 0.75e^{0.75} - 1 = 0.5878 > 0$$

1.6 Given - $x^2 - y^2 = 3$, $x^2 + y^2 = 13$, $x_0 = y_0 = \sqrt{6.5}$

To find - Real roots of given equations by Newton's Method (2 ites)

$$\text{Sol}^n; \text{ Let } f_1 = x^2 - y^2 - 3 = 0 \quad - (1)$$

$$f_2 = x^2 + y^2 - 13 = 0 \quad - (2)$$

Determine functional form of the partial derivatives,

$$J_{1,1} = \frac{\partial f_1}{\partial x} = 2x, \quad J_{1,2} = \frac{\partial f_1}{\partial y} = -2y$$

$$J_{2,1} = \frac{\partial f_2}{\partial x} = 2x, \quad J_{2,2} = \frac{\partial f_2}{\partial y} = 2y$$

$$\therefore J = \begin{bmatrix} 2x & -2y \\ 2x & 2y \end{bmatrix}$$

with $x_0 = \sqrt{6.5} \approx 2.5495$, we get

$$\begin{pmatrix} x_{n+1} \\ y_{n+1} \end{pmatrix} = \begin{pmatrix} x_n \\ y_n \end{pmatrix} - J^{-1} f \begin{pmatrix} x_n \\ y_n \end{pmatrix}$$

$\therefore n = 1^{\text{st}}$ iteration,

$$\begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} 2.5495 \\ 2.5495 \end{bmatrix} - \begin{bmatrix} 5.099 & -5.099 \\ 5.099 & 5.099 \end{bmatrix}^{-1} \times f \begin{pmatrix} 2.5495 \\ 2.5495 \end{pmatrix} \rightarrow \begin{bmatrix} -3 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} 2.8437 \\ 2.2553 \end{bmatrix}$$

$n = 2^{\text{nd}}$ iteration, $x_1 = 2.8437$, $y_1 = 2.2553$

$$\begin{bmatrix} x_2 \\ y_2 \end{bmatrix} = \begin{bmatrix} 2.8437 \\ 2.2553 \end{bmatrix} - \begin{bmatrix} 5.6874 & -4.5107 \\ 5.6874 & 4.5107 \end{bmatrix}^{-1} \times \begin{bmatrix} 0 \\ 0.1731 \end{bmatrix}$$

$$\begin{bmatrix} x_2 \\ y_2 \end{bmatrix} = \begin{bmatrix} 2.82847 \\ 2.23615 \end{bmatrix}$$

\therefore After 2 iterations, approx roots are $x = 2.82847$, $y = 2.23615$

1.7

Fixed point iteration method: Given transcendental equation $f(x) = 0$ can be converted into the form $x = g(x)$ & then using the iterative scheme with the recursive relation,

$$x_{i+1} = g(x_i), \quad i=0, 1, 2, \dots$$

with some initial guess x_0 is called fixed point method.

Given: $f(x) = 2x - \cos x - 3 = 0$

$$g(x) = \frac{\cos x + 3}{2}, \quad x_0 = \frac{\pi}{2}$$

From $f(x)$ we have, $x = \frac{\cos x + 3}{2} = g(x)$

Iteration-1 :- $x_0 = \frac{\pi}{2}, f(x_0) = \pi - 3 = 0.14159$

$$\begin{aligned} x_1 &= g(x_0) \\ &= \frac{\cos\left(\frac{\pi}{2}\right) + 3}{2} \end{aligned}$$

$$x_1 = 1.5, f(x_1) = 0.3472$$

Iteration-2 : $x_1 = 1.5$

$$\begin{aligned} x_2 &= g(1.5) \\ &= \frac{\cos(1.5) + 3}{2} = \frac{0.0707 + 3}{2} \end{aligned}$$

~~$x_2 = 1.53537$~~

$$x_2 = 1.53537, f(x_2) = -0.0173$$

1.8

Pseudocode for gradient ascent:

Let $X \rightarrow$ Data of m samples and n features

$y \rightarrow$ Output

$\alpha \rightarrow$ learning rate.

weights \rightarrow To learn.

Update step of gradient ascent is given by

$$w = w + \alpha$$

Gradient-descent ($X, y, t, \text{num_iters}$)

$$m, n = \text{shape}(X)$$

$$\alpha = t$$

$$\text{weights} = \text{zeros}(n)$$

for $i = 1 : \text{num_iters}$:

$$\hat{y} = X * \text{weights}$$

$$\text{error} = y - \hat{y}$$

$$\text{weights} = \text{weights} + \alpha * \text{error} * X$$

Maximum of $6 - (\gamma_1^2(\gamma_1^2 - 16) + \gamma_2^2(\gamma_2^2 - 9))$

\rightarrow Starting from $(0,0)$, the gradient does not update and the maximum value of function stays at 6.

- If the starting point is changed to $(0.1, 0.1)$, the maximum value the function reaches is 90.25

1.9

Batch gradient descent : In this, entire training set is used to perform one iteration of gradient descent. The average of the gradients of all the training example is taken and this mean gradient is used to update parameters.

Mini-batch GD - This uses a small subset of training data to compute the gradient.

Stochastic GD - A single training sample is used to update final gradient and update parameters.

Batch GD is used because it is good for convex or relatively smooth error manifolds as it moves somewhat directly towards an optimum solution.

But using SGD or mini-batch GD, the optimization path taken is erratic and gives inaccurate gradients but with advantage that computing the gradient is lot faster. Hence there is trade-off between accuracy in optimization steps versus the speed.

Batch GD takes faster

Batch GD gives more accurate weights but is slow if the data size is large in which case requires to high, slow computation.

1.10 Convergence proof of gradient descent (with fixed step size)

Theorem: Suppose the function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and differentiable, and that its gradient is Lipschitz continuous with constant $L > 0$ i.e. $\|\nabla f(x) - \nabla f(y)\|_2 \leq L\|x-y\|_2$ for any x, y . Then if grad. descent is run for k iterations with a fixed step size $t \leq 1/L$, it will yield a solution which satisfies,

$$f(x^{(k)}) - f(x^*) \leq \frac{\|x^{(k)} - x^*\|_2^2}{2tk} \quad - \textcircled{1}$$

where $f(x^*)$ is optimal value. This means GD is guaranteed to converge with rate $1/k$.

Q

Proof: Since ∇f is Lipschitz continuous with const. L

$$\Rightarrow \nabla^2 f(x) \leq L I,$$

or $\nabla^2 f = -L I$ is negative semidefinite matrix.

Using this, quadratic expansion of f around $f(x)$ gives following inequality,

$$\begin{aligned} \nabla f(y) &\leq f(x) + \nabla f(x)^T(y-x) + \frac{1}{2}\nabla^2 f(x)\|y-x\|_2^2 \\ &\leq f(x) + \nabla f(x)^T(y-x) + \frac{1}{2}L\|y-x\|_2^2 \end{aligned}$$

By putting grad. desc. update in this, let $y = x^+ = x - t\nabla f(x)$,

$$\begin{aligned} f(x^+) &\leq f(x) + \nabla f(x)^T(x^+-x) + \frac{1}{2}L\|x^+-x\|_2^2 \\ &= f(x) + \nabla f(x)^T(x - t\nabla f(x) - x) + \frac{1}{2}L\|x - t\nabla f(x) - x\|_2^2 \\ &= f(x) + \nabla f(x)^Tt\nabla f(x) + \frac{1}{2}L\|t\nabla f(x)\|_2^2 \\ &= f(x) - t\|\nabla f(x)\|_2^2 + \frac{1}{2}L^2\|\nabla f(x)\|_2^2 \\ &= f(x) - (1 - \frac{1}{2}Lt) + \frac{1}{2}\|\nabla f(x)\|_2^2 \quad - \textcircled{2} \end{aligned}$$

Using $t \leq 1/L$, we know that,

$$-(1 - \frac{1}{2}Lt) = \frac{1}{2}Lt - 1 \leq \frac{1}{2}t\left(\frac{1}{L}\right) - 1 = \frac{1}{2} - 1 = -\frac{1}{2}$$

Putting this in $\textcircled{2}$

$$f(x^+) \leq f(x) - \frac{1}{2}t\|\nabla f(x)\|_2^2 \quad - \textcircled{3}$$

Since $\frac{1}{2}t\|\nabla f(x)\|^2$ will always be FOC unless $\nabla f(x) = 0$, this inequality implies that objective function value strictly decreases with each iteration of grad. desc. until it reaches the optimal value $f(x) = f(x^*)$. This holds only if t is chosen small enough i.e. $t \leq 1/L$

We can bound $f(x^+)$, the objective value at next iteration in terms of $f(x^*)$, the optimal objective value, since f is convex,

$$f(x^*) \geq f(x) + \nabla f(x)^T (x^* - x)$$

$$f(x) = f(x^*) + \nabla f(x)^T (x - x^*)$$

Putting this in ③

$$f(x^+) \leq f(x^*) + \nabla f(x)^T (x - x^*) - \frac{t}{2} \|\nabla f(x)\|_2^2$$

$$f(x^+) - f(x^*) \leq \frac{1}{2t} (2t \nabla f(x)^T (x - x^*) - t^2 \|\nabla f(x)\|_2^2)$$

$$f(x^+) - f(x^*) \leq \frac{1}{2t} (2t \nabla f(x)^T (x - x^*) - t^2 \|\nabla f(x)\|_2^2 - \|x - x^*\|_2^2)$$

$$f(x^+) - f(x^*) \leq \frac{1}{2t} (\|x - x^*\|_2^2 - \|x - t \nabla f(x) - x^*\|_2^2) \quad - ④$$

By definition, $x^+ = x - t \nabla f(x)$, putting this in ④

$$f(x^+) - f(x^*) \leq \frac{1}{2t} (\|x - x^*\|_2^2 - \|x^+ - x^*\|_2^2) = ⑤$$

This inequality holds for x^+ in every iteration of GD, summing over iterations

$$\sum_{i=1}^K f(x^{(i)} - f(x^*)) \leq \sum_{i=1}^K \frac{1}{2t} (\|x^{(i-1)} - x^*\|_2^2 - \|x^{(i)} - x^*\|_2^2)$$

$$= \frac{1}{2t} (\|x^{(0)} - x^*\|_2^2 - \|x^{(K)} - x^*\|_2^2)$$

$$\leq \frac{1}{2t} (\|x^{(0)} - x^*\|_2^2) \quad - ⑥$$

Now, f is decreasing on every iteration, we can conclude that.

$$f(x^{(k)} - f(x^*)) \leq \frac{1}{K} \sum_{i=1}^K f(x^{(i)}) - f(x^*)$$

$$\leq \frac{\|x^{(0)} - x^*\|_2^2}{2tk} \quad [\text{substituting from } ⑥]$$

Hence proves ①, the convergence of gradient descent.