

$$\text{Primal} = \max \{ c^T x \mid Ax \leq b, x \geq 0 \}$$

$$\text{Dual} = \min \{ b^T y \mid A^T y \geq c, y \geq 0 \}$$

$\rightarrow$  Strong duality:  $c^T x^* = b^T y^*$

$\rightarrow$  Weak duality:  $c^T x \leq b^T y$   
if feasible  $x, y$ .

$$P = \max \{ c^T x \mid Ax \leq b, x \in \mathbb{R}^n \}$$

$$D = \min \{ b^T y \mid A^T y = c, y \geq 0 \}$$

$\rightarrow$  Let  $x^*$  be the primal optima

Let  $I$  be the set of constraints that are tight.

$$a_i^T x^* = b, i \in I$$

$$a_i^T x^* < b, i \notin I.$$

Claim:  $c$  is in  $K = \{x \mid x = \sum_{i \in I} \lambda_i a_i, \lambda_i \geq 0\}$

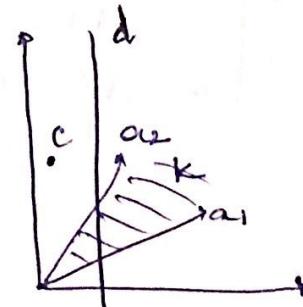
Proof by contradiction  $\Rightarrow$

[if  $c$  is not in  $K$ , then]  
[there exists "d"

$d$  separates  $c$  from  $K$  by

$$a_i^T d \leq 0, i \in I$$

$$c^T d > 0.$$



$$z = x^* + \epsilon d$$

$$a_i^T z \leq b_i, i \notin I$$

$$a_i^T z = a_i^T (x^* + \varepsilon d)$$

$$= a_i^T x^* + \varepsilon \underbrace{a_i^T d}_{-\text{ve}}$$

↗ feasible point.

$z \rightarrow$  Not going to violate for all constraints.

Objective =

$$c^T z = c^T [x^* + \varepsilon d] = c^T x^* + \varepsilon \underbrace{c^T d}_{+\text{ve}}$$

$z$  is superior to  $x^*$

(2) doesn't violate constraints  $\downarrow$

$\downarrow$  has a better obj ; it contradicts our statement,  $x^*$  is optimal.

$\rightarrow$  Hence, there exists no such  $d$  and  
our claim is true.

claim:  $c$  is in  $K = \{x \mid x = \sum_{i \in I} \lambda_i a_i, \lambda_i \geq 0\}$

$\lambda \geq 0$ ;  $\lambda_i = 0$ , when  $i \notin I$

$$A^T \lambda = \sum_i \lambda_i a_i = c$$

$$b^T \lambda = \sum_i \lambda_i b_i = \sum_i a_i^T x^* \lambda_i = c^T x^*$$

tight.

## Sparse coding & Dictionary Learning →

Compact representation → complete basis (Fourier.).

↓  
Instead of using a standard basis, can we learn it?

$$\text{available } \underset{n}{\underbrace{\gamma}} = n \underset{k}{\underbrace{D}} \underset{x}{\underbrace{|}}$$

$N \gg k$   
 $k \gg n$

$$\gamma = \{\gamma_1, \dots, \gamma_N\}$$

↓  
observations

$$\gamma_i = D x_i$$

Problem :- Find  $x_i$  (for each  $\gamma_i$ )?

Sparse coding :- Find the sparsest  $x$ , s.t  $\gamma_i = D x_i$

Dictionary learning :- Learn  $D$  (and  $x_i$ ) from  $\gamma$

$\gamma_i$  and  $D$  are given,

Sparse Coding ⇒

given set of samples.

and corresponding coeff.

$$\begin{array}{l} \text{Min } \|x_i\|_0 \\ \text{s.t. } \gamma_i = D x_i \end{array}$$

sparse.

greedy OMP.

relax the norm  $(0 \rightarrow 1)$

LP problem

$$\begin{array}{l} \text{Min } \|x_i\|_1 \\ \text{s.t. } \gamma_i = D x_i \end{array}$$

BP

Start with a completely sparse  $x$  vector and find the big most correlated column in  $D$  and take the residual and repeat the same process  $\Rightarrow$  approximate soln

under some conditions,  
both the optimas are same

"greedy"  
completely sparse vector

OMP  $\Rightarrow \text{Min } \|Ax - b\|$   $\curvearrowright$  minimise the error [residual]  
 s.t  $\|x\|_0 \leq T$  (atmost  $T$  elements)

- ① Select most correlated column, add to basis.  
start with  $y$
- ② Find  $x(i)$  corresponding to the new basis.  
and calculate the residual.

\* if  $T$  is larger, error goes lower and lower.

### K-SVD & K-Means:-

K-Means  $\rightarrow$  dictionary element  $\Rightarrow$  such that mean square error is small.

$\downarrow$   
sparsest code  
[extreme]  
 $\downarrow$   
representation  
of the elements.

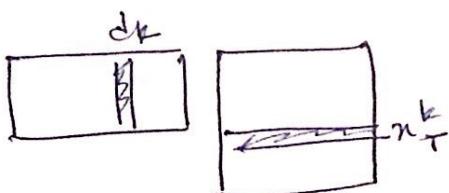
find the mean & encoding, given samples  
 ally, given  $y$ , find  $D$  and  $X$   
 $\uparrow$   
the right dictionary.

$$[y_1 \dots y_N] = [d_1 \dots d_K] \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_K \end{bmatrix}$$

Dict  $\quad \quad \quad$  code  
 $\downarrow$   $\quad \quad \quad$   $\downarrow$   
 $y \quad \quad \quad X$

$$y = Dx$$

$$\text{Min } \|y - Dx\|, \quad y - Dx = y - \sum_{i=1}^k d_i e_i^T$$



$$Y - DX = Y - \sum_{i=1}^k d_i x_i^T$$

↑ T representation: row vector.

$$= \underbrace{\left[ Y - \sum_{\substack{i=1 \\ i \neq j}}^k d_i x_i^T \right]}_E - d_j x_j^T$$

↑ assume, except  $j$ , all are fixed

$$= E - d_j x_j^T$$

Compute SVD;

on computing SVD,  $x_j^T \rightarrow$  not going to be sparse.

modify by retaining sparsity;  $Y_R$

$$E_R - d_j x_j^T$$

only non-zero elements are considered

- ① Initialize  $D$  (like <sup>means for</sup> K-Means)
- ② Do sparse coding.
- ③ Find  $D$  with the same sparsity (find new set of means)
- ④ Repeat ② & ③.

$$Y = DX$$

Sparse coding

given  $Y$  &  $D$ , find sparsest  $x$

→ greedy → OMP

→  $L1 \rightarrow$  BP (LP)

Dictionary Learning.

given  $Y$ ;

Find  $D$  and  $X$

# Non linear Optimisation →

## Convexity →

convex set,  $\equiv n, y \in C \text{ & } \theta \in \mathbb{R}$

$$\boxed{C} \quad \theta x + (1-\theta)y \in C \quad [\text{Convex combination}]$$

↓  
should also lie in  
the convex set

Convex function  $\equiv D(f) \rightarrow \text{domain of } f$  [convex]

if  $x, y \in D(f)$

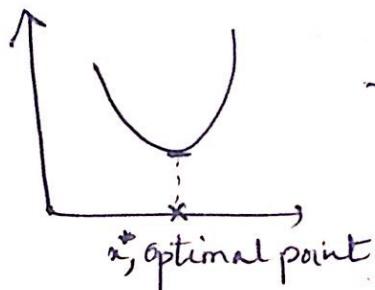
$$f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$$

line should lie  
above the  
function.



## Unconstrained -

$$\min_x f(x)$$



→ Necessary condition

$$\nabla_x f(x) = 0$$

↓  
generally, not alone  
enough.

so, we need →

$$\nabla_x^2 f(x) \geq 0$$

\* contour curve - along  
the curve value doesn't  
change

## Constrained -

$$\min_x f(x)$$

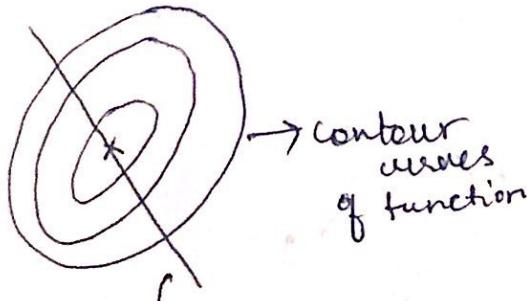
$$\text{s.t. } h(x) = 0$$

→ Necessary condition

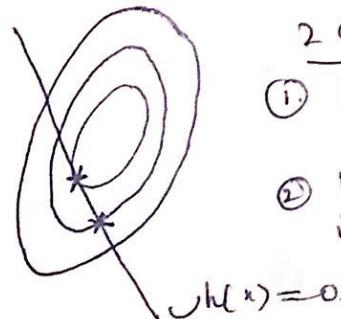
$$h(x) = 0$$

$$\nabla_x f(x) = 0$$

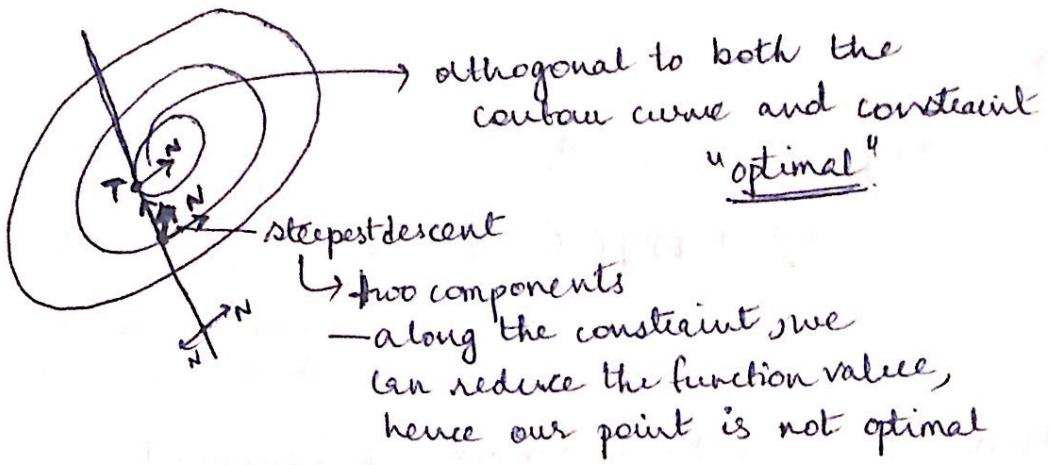
in this case →



In this case,  
constraint doesn't matter.



- 2 conditions
- ①  $h(x)$  tangent to contours
  - ②  $h(x)$  intersecting contours



Tangent point  $\Rightarrow h(x) = 0$

$$\nabla_x f(x) + \lambda \nabla_x h(x) = 0.$$

$$\begin{array}{l} \min_x f(x) \\ \text{s.t. } h(x) = 0 \end{array} \xrightarrow{\text{unconstrained problem}} \min_x f(x) + \lambda h(x)$$

multiple constraints

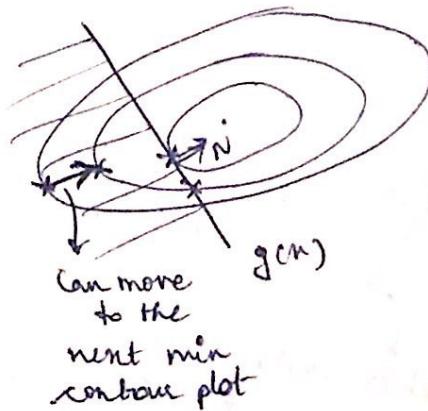
$$\begin{array}{l} \text{s.t. } h_i(x) = 0 \\ i = 1, 2, \dots, m \end{array} \xrightarrow{} \min_x f(x) + \sum_i \lambda_i h_i(x)$$

$$\nabla_x f(x) + \sum_i \lambda_i h_i(x) = 0.$$

Inequality constraint :-

$$\begin{array}{l} \min_x f(x) \\ \text{s.t. } g(x) \leq 0 \end{array}$$

$$\begin{aligned} \nabla_x f(x) &= 0 \\ \nabla_x f(x) + \mu \nabla_x g(x) &= 0 \\ \mu = 0 \& g(x) = 0 \Rightarrow \boxed{\mu g(x) = 0} \end{aligned}$$

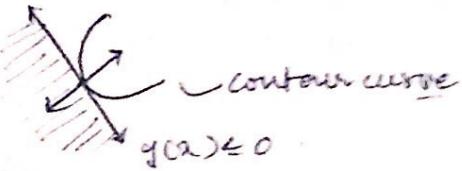


$$g(x) \leq 0$$

$$\nabla_x f(x) + \mu \nabla_x g(x) = 0$$

$$\boxed{\mu > 0}$$

gradients can only  
be antiparallel  
[not parallel]



Multiple constraints:

$$g_j(x) \leq 0 \quad \forall j; \quad h_i(x) = 0 \quad \forall i$$

$$\nabla_x f(x) + \sum_j \mu_j \nabla g_j(x) + \sum_i \lambda_i \nabla_x h_i(x)$$

$$\mu_j \geq 0 \quad \forall j$$

$$\mu_j g_j(x) = 0 \quad \forall j$$

$g(x) \rightarrow$  inequality  
cond's

$h_i(x) \rightarrow$  equality  
cond's

→ KKT Conditions

↓  
not always necessary  
for the optimal sol's

Opt problems have to satisfy

↓  
Sufficiency conditions ← for KKT to be sufficient

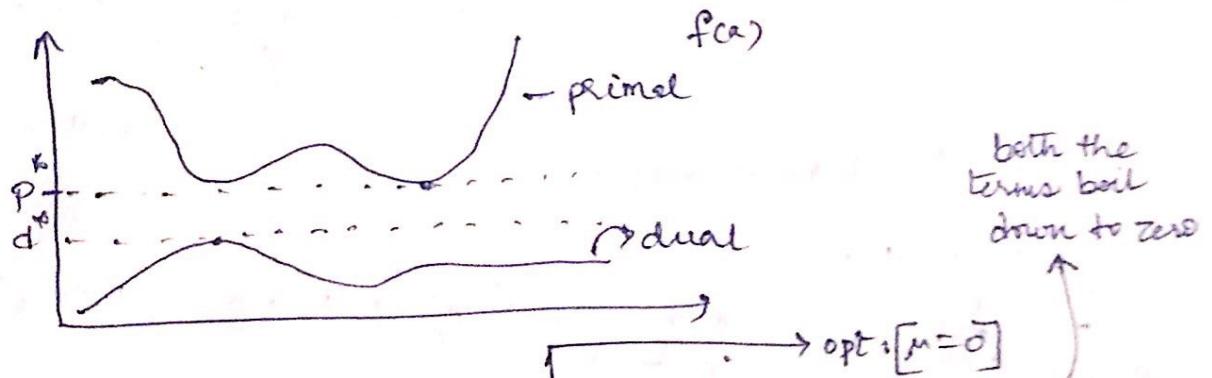
→ Convex optimisation ⇒ Slater conditions

Duality :-  $P \equiv \min_x f(x)$

$$\text{st } g_j(x) \leq 0 \quad j=1, \dots, n$$

$$h_i(x) = 0 \quad j=1, \dots, m.$$

} Primal



$$P \equiv \min_x \max_{\mu \geq 0, \lambda} f(x) + \sum_j \mu_j g_j(x) + \sum_i \lambda_i h_i(x)$$

$$\underbrace{\quad}_{L(x, \mu \geq 0, \lambda)} \Rightarrow \min_x \begin{cases} f(x), & x \text{ is feasible} \\ \infty, & \text{otherwise} \end{cases}$$

$$P = \min_{\mathbf{x}} \max_{\boldsymbol{\mu} \geq 0, \boldsymbol{\lambda}} f(\mathbf{x}) + \sum_j \mu_j g_j(\mathbf{x}) + \sum_i \lambda_i h_i(\mathbf{x})$$

↓

$$D = \max_{\boldsymbol{\mu} \geq 0, \boldsymbol{\lambda}} \min_{\mathbf{x}} f(\mathbf{x}) + \sum_j \mu_j g_j(\mathbf{x}) + \sum_i \lambda_i h_i(\mathbf{x})$$

Lagrangian dual problem.

Dual lower bound of Primal :-

$$\max_{\boldsymbol{\mu} \geq 0} \min_{\mathbf{x}} L(\cdot) \leq \min_{\mathbf{x}} \max_{\boldsymbol{\mu} \geq 0} L(\cdot)$$

$$\beta = [\mu, \lambda]$$

(dual variable)

whereas,

$\mathbf{x}$  = primal variable.

Lagrangian dual problem = always concave  
in  $\underline{g(f(x))}$

↓      ↓  
lower bound to primal

$$\max_{\boldsymbol{\mu} \geq 0} \min_{\mathbf{x}} [f(\mathbf{x}) + \sum_j \mu_j g_j(\mathbf{x}) + \sum_i \lambda_i h_i(\mathbf{x})]$$

↓

$$\max_{\boldsymbol{\mu} \geq 0} \underbrace{[k_1 + \mu_1 k_2 + \lambda_1 k_3]}_{\text{linear function}} \rightarrow \text{maximising over all the linear functions.}$$

$$g(\mu, \lambda)$$

Lagrangian function :-  $L(\mathbf{x}, \beta)$

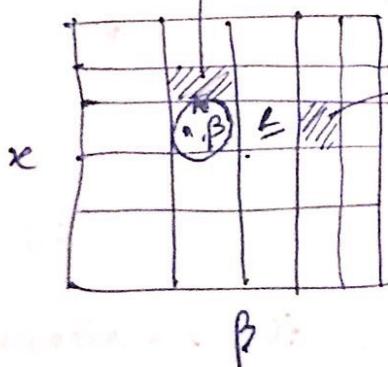
$$P: \min_{\mathbf{x}} \max_{\beta} L(\mathbf{x}, \beta)$$

first rows  
then cols

$$D: \max_{\beta} \min_{\mathbf{x}} L(\mathbf{x}, \beta)$$

first  $\mathbf{x}$  and then  $\beta$   
columns rows

$\mathbf{x}, \beta^*$  (dual optima)



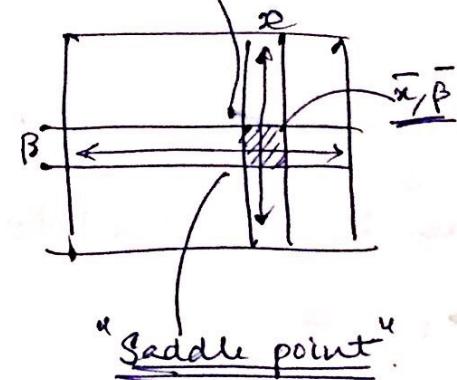
$$L(\mathbf{x}^*, \beta) \geq L(\mathbf{x}, \beta) \geq L(\mathbf{x}, \beta^*)$$

Primal obj  $\geq$  Dual obj

primal optima ( $p^*$ )  $\geq$  dual optima ( $d^*$ )  
 $p^* \geq d^*$

$p^* > d^*$  (weak duality)  
 $p^* = d^*$  (strong duality)

red and green  
in the same box → saddle point



Lagrangian [with saddle point]

$$\Rightarrow L(\bar{x}, \bar{\beta}) \leq L(\bar{x}, \bar{\beta}) \leq L(x, \bar{\beta})$$

$$L(\bar{x}, \bar{\beta}) \leq \max_{\beta} L(x, \beta)$$

At the saddle point,  $L(\bar{x}, \bar{\beta}) = \min_x \max_{\beta} (L(x, \beta)) = p^*$

↳ primal optima      ↓      primal

$$L(\bar{x}, \bar{\beta}) \geq \min_x L(x, \beta)$$

At the saddle point,  $L(\bar{x}, \bar{\beta}) = \max_{\beta} \min_x L(x, \beta) = d^*$

↳ dual optima      ↓      Dual

$L(\bar{x}, \bar{\beta}) \Rightarrow$   
 $p^* = d^*$   
 $\bar{x} = x^*, \bar{\beta} = \beta^*$   
[at the saddle point]

## Nonlinear Equations and Optimisation :-

$$\min \|f(x)\|$$



$$f(x) = 0$$

$$f(x_1, x_2, \dots, x_n) = 0.$$

multivariable

$$f_1(x_1, x_2, \dots, x_n) = 0.$$

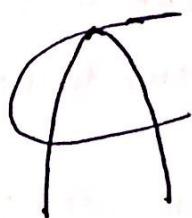
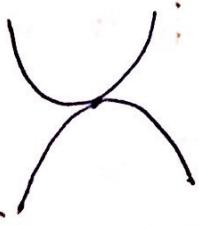
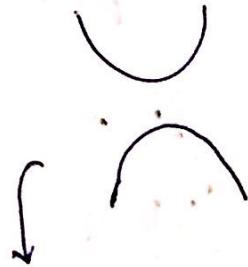
⋮

$$f(\bar{x}) = 0 \Rightarrow \mathbb{R}^n \rightarrow \mathbb{R}^m.$$

$$x^0 \rightarrow x^1 \rightarrow x^2 \rightarrow \dots \rightarrow x^k \rightarrow x^{k+1}$$

Iterative :-

- \* Initialize
- \* Update
- \* Iterate/Repeat until convergence.



possible cases of solutions.

$$\rightarrow f(x) = 0.$$

$\rightarrow$  How many times function is evaluated?

$\rightarrow$  How fast does it converge?

$\rightarrow$  Local or global?

Situations:-

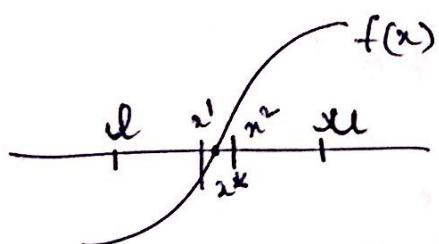
① We know the form of  $f(x)$ .

② We know only to evaluate  $f(x)$

$[f(x) - \text{black box}]$

$$\rightarrow f(x)$$

\* Bisection  $\rightarrow$



sign of  $f(x)$  differs at  $u$  and  $l$   
 this implies that there is a  
 solution in between  $u$  and  $l$ .  
 (exists a root)

Initialize  $\rightarrow u, l$

$$\rightarrow x^k = \frac{u+l}{2}$$

If  $f(x^k) f(u) < 0$ ,  $l = x^k$

If  $f(x^k) f(l) < 0$ ,  $u = x^k$

$k = k + 1 \rightarrow$  slowly converges to the true value.

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

$$x^* = -1, 2$$

Initialize;  $u=5, l=1$

$$x^1 = 3$$

$$x^2 = 2$$

$$x^3 = 2.5$$

what's good about this scheme?

→ we don't need to know the form of the function.

→ we don't need to differentiate

→ there is a bound on the number of steps

Initialising  $u$  and  $l \rightarrow$  might be a problem.

Local optimum Sol<sup>n</sup>

~~Hausman~~ Iterations → interval is becoming half.

$$|x^k - x^{k-1}| < \epsilon \rightarrow k \geq \log_{\frac{1}{2}} \left[ \text{in terms of the initial interval } [u, l] \text{ and } \epsilon \right]$$

\* Fixed point :-

$p$  is a fixed point.

→ if  $p = g(p)$

→ Initialize

Iterate

$$x^{k+1} \leftarrow g(x^k)$$

until convergence.

} fixed point iteration scheme.

Iterations

$$x^1 = g(x^0) = 3$$

$$x^2 = g(x^1) = 1.6$$

$$x^3 = g(x^2) = 2.9$$

$$\xrightarrow{\text{converged to } 2} x^4 = g(x^3) = 1.9.$$

$$\text{or } x^2 - x - 2 = 0$$

$$\downarrow x = x^2 - 2$$

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

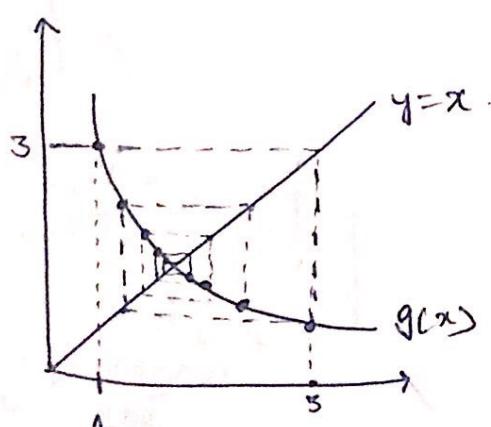
$$x^0 = 1$$

$$x^1 = -1$$

$$x^2 = -1$$

↓ converges.  
to  $\underline{-1}$

(good scheme)



$$\text{if } x^0 = 2.1 \\ x^1 = 2.4 \\ x^2 = 2.5$$

↓ diverging  
(bad scheme)

## Fried point iteration scheme $\rightarrow$

$$\Rightarrow x^2 = a ; x = a/x$$

$$g(x) = a/x$$

$$x^0 = 1$$

$$x^1 = g(x^0) = a$$

$$x^2 = g(x^1) = 1$$

$$x^3 = g(x^2) = a$$

continuously  
oscillating

$$g(x) = \frac{1}{2} \left( \frac{a}{x} + 1 \right)$$

$$x^0 = 2$$

$$x^1 = g(x^0) = \frac{a+2}{4}$$

$$x^2 = g(x^1) = \frac{5a+2}{2(a+2)}$$

$\downarrow$

doesn't  
converge.

→ Not all  $g(x)$ 's  
work / converge.

if  $a = 9$

$$x^0 = 2$$

$$x^1 = 3.25$$

$$x^2 = 3.009$$

$\downarrow$   
converging.

\* If  $|g'(x^*)| < 1$ , the scheme is locally convergent

→ Many  $g(x)$ 's can be created from the original problem.

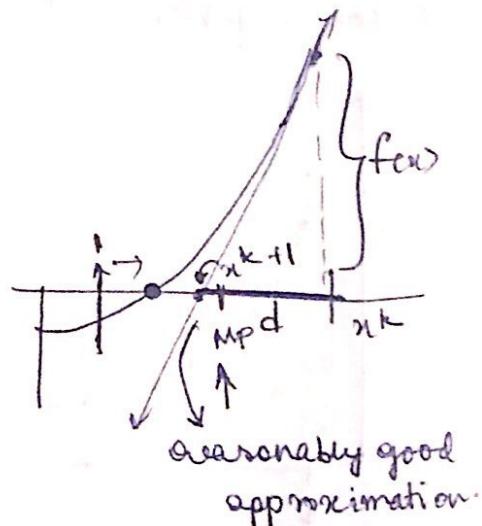
## Newton's Method $\rightarrow$

$$x^{k+1} \leftarrow x^k - \frac{f(x)}{f'(x)}$$

$$\tan \theta = f'(x) = \frac{f(x)}{d}$$

$$d = \frac{f(x)}{f'(x)}$$

$$x^{k+1} \leftarrow x^k - d \Rightarrow x^{k+1} \leftarrow x^k - \frac{f(x)}{f'(x)}$$



## Limitations of Newton's Method

→ we need to find the differentiation of  $f(x)$

Iterative process.

$$\left[ \begin{array}{l} x^1 \leftarrow x^0 - \frac{f(x)}{f'(x)} \\ \downarrow \\ x^{k+1} \leftarrow x^k - \frac{f(x^k)}{f'(x^k)} \end{array} \right]$$

$$f(x) = x^2 - x - 2.$$

start with  $x^0 = 1$ ;

$$x^1 \leftarrow 1 - \frac{[-2]}{1} = 3$$

$$x^2 \leftarrow 3 - \frac{4}{5} = 2.2.$$

↓  
slowly converges to  $x^i = 2$ .

Secant Method ⇒

$$\left\{ \begin{array}{l} x^{k+1} = x^k - \frac{f(x^k)}{f'(x^k)} \rightarrow \text{instead of differentiating;} \\ f'(x^k) = \frac{f(x^k) - f(x^{k-1})}{x^k - x^{k-1}} \\ = x^k - \frac{f(x^k) \cdot (x^k - x^{k-1})}{f(x^k) - f(x^{k-1})} \end{array} \right.$$

Advantage over Newton's method.

→ No need to compute the derivative,  
just evaluating will be enough

Disadvantage ⇒ approximate of derivative is computed  
→ might take more time to converge.

Computation - becomes easier

↓  
small amount of memory is needed.  
(reevaluation is not needed).

$$p = g(p)$$

$$x^{k+1} \leftarrow g(x^k)$$

Locally converging

$$\min |g(x)|$$



$$\text{solve: } \frac{g(x) = 0}{\downarrow}$$

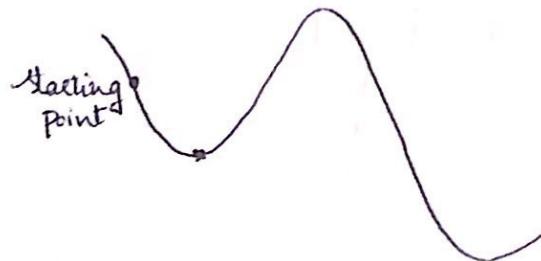
m var, n equations.

#### \* Fixed point method -

$h(x) \rightarrow$  locally converges if →

$$|h'(x^*)| < 1$$

$$\frac{x^{k+1} - x^k}{e^{k+1}} = \underbrace{h(x^k) - h(x^*)}_{\text{From Meanvalue theorem}}$$



— From Meanvalue theorem

$$e^{k+1} = h'(\cdot) [x^k - x^*]$$

$$e^{k+1} = \underbrace{h'(\cdot)}_{(h'(x) < 1)} \cdot e^k$$

→ on increasing iterations,  
clear decrease.

$$e^k \rightarrow e^{k+1} \downarrow$$

#### \* Newton's method -

$$x^{k+1} = x^k - \frac{f(x^k)}{f'(x^k)}$$

$$h(x) = x - f(x)/f'(x) \Rightarrow h'(x) = 1 - \frac{f \cdot f' - f \cdot f''}{f' \cdot f''}$$

$$= \frac{f \cdot f''}{f' \cdot f''}$$

$$\text{Ex: } x^3 - 2x - 5 = f(x)$$

Newton's iteration -  $x^0 = -1$

$$x^{k+1} \leftarrow x^k - \frac{x^3 - 2x - 5}{3x^2 - 2}$$

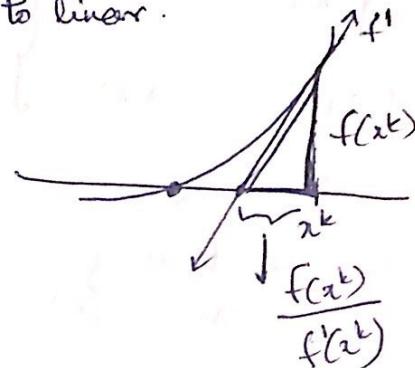
$$x^1 \leftarrow 3$$

$$x^2 \leftarrow 59/25$$

approximates to  $\approx 2.1$

$$x^{k+1} - x^k = -\frac{f(x^k)}{f'(x^k)}$$

limiting to linear.



Taylor series -

$$f(y) = f(x) + \frac{f'(x)}{1!}(y-x) + \frac{f''(x)}{2!}(y-x)^2 + \dots$$

we have to find  $x^{k+1}$  such that  $f(x^{k+1}) = 0$ .

$$y = x^{k+1}$$

$$f(x^{k+1}) = f(x^k) + \frac{f'(x^k)}{1!}(x^{k+1} - x^k) + \underbrace{\frac{f''(x^k)}{2!}(x^{k+1} - x^k)^2}_{\text{approximated as linear, so higher order terms = 0.}} + \dots$$

$$x^{k+1} = x^k - \frac{f(x^k)}{f'(x^k)}$$

approximated as linear, so higher order terms = 0.

Quadratic approximation  $\rightarrow$

$$f(x^{k+1}) = f(x^k) + \frac{f'(x^k)}{1!}(x^{k+1} - x^k) + \frac{f''(x^k)}{2!}(x^{k+1} - x^k)^2$$

$$0 = f(x^k) + \frac{f'(x^k)}{1!}(x^{k+1} - x^k) + \frac{f''(x^k)}{2!}(x^{k+1} - x^k)^2$$

Solve for  $x^{k+1}$ .

(but Hessian matrix  
(memory intensive))

$\Rightarrow$  gives a better solution than linear approximation

and even more higher order terms give better and better

solutions  $\rightarrow$  so why are we not doing it?

- computationally intensive

In higher order approximations  
 ↳ no. of iterations - will be slow  
 but each iteration - is computationally complex to solve.

Min  $g(x)$

$$g(x_1, x_2, \dots, x_n)$$

$$g: \mathbb{R}^N \rightarrow \mathbb{R}$$

(scalar function of  
n variables)

$$\left. \begin{array}{l} \\ \end{array} \right\} \rightarrow \nabla g = 0.$$

$$\begin{bmatrix} \frac{\partial g}{\partial x_1} \\ \frac{\partial g}{\partial x_2} \\ \vdots \\ \frac{\partial g}{\partial x_n} \end{bmatrix} = 0$$

system  
of  
equations.

n equations

$$\Rightarrow \bar{f}(\bar{x}) = 0.$$

$$\text{where } \bar{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

optimising  $g(x)$

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \\ \vdots & & \\ \frac{\partial f_1}{\partial x_n} & \dots & \end{bmatrix} \quad \text{Jacobian}$$

equivalent

$$H = \begin{bmatrix} \frac{\partial^2 g}{\partial x_1 \partial x_1} & \frac{\partial^2 g}{\partial x_1 \partial x_2} & \dots \\ \vdots & \ddots & \end{bmatrix} \quad \text{Hessian (second derivative matrix)}$$

system of equations  $\Rightarrow \bar{f}(\bar{x}) = 0.$

$$\bar{x}^{k+1} \leftarrow \bar{x}^k + s.$$

$$\bar{f}(\bar{x} + s) = \bar{f}(\bar{x}) + J \cdot s$$

assume  $s=0$ , just like in Taylor's series

$$J \cdot s = -\bar{f}$$

Jacobian

change or  $\Delta x$

$$J.S = -\bar{f} \quad (\text{since } \bar{f}(\pi) = 0)$$

→ in terms of  $g$ ;  $\min g(x)$

$$\Rightarrow \boxed{\nabla^2 g \cdot s = -\nabla g}$$

$$\Rightarrow \boxed{x^{\text{new}} \leftarrow x^{\text{old}} - (\nabla^2 g)^{-1} \nabla g}$$

$$g(x) = (x-2)^2 - 5 \quad \boxed{\min = 2} \Rightarrow \eta^* = \frac{1}{2} \quad [\text{optimal learning rate}]$$

$$= x^2 - 4x - 1 \quad \Rightarrow \quad x^{\text{new}} \leftarrow x^{\text{old}} - \frac{1}{2}[2x - 4]$$

$$g'(x) = 2x - 4$$

$$\text{suppose; } x^0 = 5$$

$$g''(x) = 2.$$

$$\begin{cases} x^1 = 5 - \frac{1}{2}[10-4] = 2 = x^* \\ x^0 = 6 \\ x^1 = 6 - \frac{1}{2}[12-4] = 2 = x^* \end{cases}$$

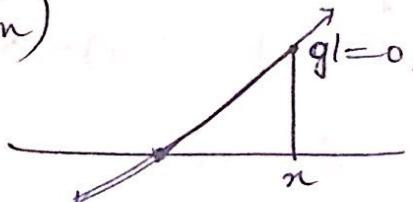
why?  $x$  gets cancelled out and we end up with opt, why?

→ we are getting the optimal [true solution] independent of the initialisation

$$\text{Problem} = \min g(x)$$

But, we are solving  $\Rightarrow \nabla g = 0$  (a linear function)

so in one go, we get the solution.



$$\left[ \eta = \frac{1}{\sqrt{2} g} \right]$$

learning rate [in gradient descent]

$$\Rightarrow x^{\text{new}} \leftarrow x^{\text{old}} - \eta \nabla g \quad [\text{in the above case, if } \eta^* = \frac{1}{2}, \text{ it worked perfectly}]$$

$$\eta = 2;$$

$$\begin{aligned} x^{\text{new}} &\leftarrow x^{\text{old}} - 2 \nabla g \\ &= x^k - 2[2x^k - 4] \end{aligned}$$

$$x^0 = 5$$

$$\begin{aligned} x^1 &\leftarrow 5 - [10-4] = -1 \\ x^2 &\leftarrow -1 - [-6] = 5 \end{aligned} \quad \text{oscillates}$$

$$x^0 = 5;$$

$$x^1 = 5 - 2[2.5-4] = -7$$

$$\begin{aligned} x^2 &= -7 - 2[2.5-4] = 29 \\ &\downarrow \text{Diverging.} \end{aligned}$$

Min  $g(x)$

$$\nabla g = 0$$

$$\begin{bmatrix} \frac{\partial g}{\partial x_1} \\ \frac{\partial g}{\partial x_n} \end{bmatrix} = 0 \quad // f(x) = 0.$$

$$H = \begin{bmatrix} \frac{\partial^2 g}{\partial x_1^2} & \frac{\partial^2 g}{\partial x_1 \partial x_n} \\ \vdots & \vdots \\ \frac{\partial^2 g}{\partial x_n \partial x_1} & \frac{\partial^2 g}{\partial x_n^2} \end{bmatrix}$$

Solve  $\bar{f}(\bar{x}) = \bar{0}$

$$\bar{f}(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_n(x) \end{bmatrix}$$

$$\underbrace{J}_{\text{Jacobian}} = \begin{bmatrix} \frac{\partial \bar{f}}{\partial x_1} & \frac{\partial \bar{f}}{\partial x_2} & \dots & \frac{\partial \bar{f}}{\partial x_n} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial f_1}{\partial x_1} \\ \vdots \\ \frac{\partial f_n}{\partial x_n} \end{bmatrix} \quad \begin{bmatrix} \frac{\partial f_1}{\partial x_n} \\ \vdots \\ \frac{\partial f_n}{\partial x_n} \end{bmatrix}$$

→ H and J are the same matrices  
because  $f = \nabla g$ ;

$$J_f = \nabla^2 g \rightarrow \text{Hessian of } g.$$

Problem: Solve  $f(x) = 0$ ;

$$f(\bar{y}) = f(\bar{x}) + [\bar{y} - \bar{x}]^T \nabla f(\bar{x})$$

$$\bar{y} - \bar{x} = s$$

Find  $y$  st.  $f(y) = 0$ .

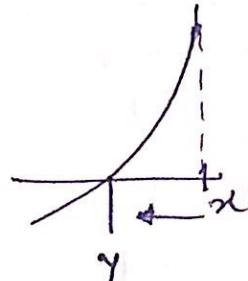
$$\bar{f} + J_s = 0 \quad , x \text{ is a vector.}$$

$$Js = -\bar{f}$$

$$s = -J^{-1}\bar{f}$$

$$y = x + s \\ = x - J^{-1}\bar{f}$$

$$\Rightarrow x^{k+1} \leftarrow x^k - [J^2]^{-1} \nabla g$$



Iterative Algorithm:-

Initialize  $x^0$

for  $k=0$  to  $\infty$   
Solve  $s = -J^{-1}\bar{f}$

$$x^{k+1} \leftarrow x^k + s = x^k - J^{-1}\bar{f}$$

$\nabla^2$  or  $[\nabla^2 g]^{-1}$  computation is expensive

$$x^{k+1} \leftarrow x^k - [\nabla^2 g]^{-1} \nabla g \quad \begin{array}{l} \text{right step size [but complex]} \\ \text{computation} \end{array}$$

↓

$$x^{k+1} \leftarrow x^k - \eta \nabla g \quad \begin{array}{l} \text{Gradient Descent.} \\ (\text{putting a scalar in the place of } \nabla^2 g) \end{array}$$

matrix sizes are usually huge in practice [in deep learning]

→ though second order approximation is useful, it is not practically feasible [gradient descent is much easier]

Ex:  $g = (x-2)^2 - 5 = x^2 - 4x - 1$

$$\nabla g = 2x - 4.$$

ad hoc learning rate

$$\nabla^2 g = 2.$$

instead, if we use gradient descent

$$x^{k+1} \leftarrow x^k - \eta \nabla g$$

$$x^0 = 5 \rightarrow \eta = 1$$

$$x^1 = 5 - 6 = -1$$

$$x^2 = -1 - (-6) = 5$$

→ oscillating.

$\eta = \frac{1}{2}$  optimal solution

$$x^0 = 5, \eta = \frac{1}{2}$$

$\eta = 1$  oscillating

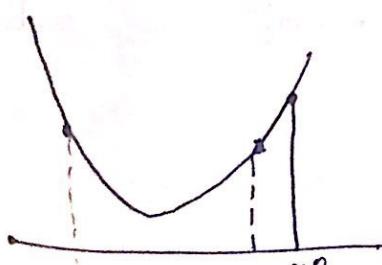
$$x^1 = 5 - 2(6) = -7$$

$\eta = 2$  diverging.

$$x^2 = -7 - 2(-18) = 20$$

→ diverging.

that is why we need the appropriate learning rate.



changing,  $\nabla g$ -direction of change.  
for larger  $\eta$

Least Square problem  $\Rightarrow$

$$\text{Min } q(\bar{x}) = [Ax - b]^T [Ax - b]$$

$$= \bar{x}^T A^T A \bar{x} - 2(A\bar{x})^T b + b^T b.$$

$$\nabla q = 2A^T A \bar{x} - 2A^T b.$$

$$\nabla^2 q = 2A^T A.$$

Initialise  $\bar{x}^0$ :

$$\bar{x}^1 \leftarrow \bar{x}^0 - (A^T A)^{-1} [2A^T A \bar{x}^0 - 2A^T b]$$

$$\leftarrow \bar{x}^0 - (A^T A)^{-1} \cdot (A^T A) \bar{x}^0 + (A^T A)^{-1} A^T b.$$

$$\leftarrow \bar{x}^0 - \bar{x}^0 + (A^T A)^{-1} A^T b.$$

$$\leftarrow (A^T A)^{-1} A^T b. \quad (\sim \text{independent of initialisation})$$

→ you get the solution in one single step.

[No need of iterative process]

Even in the case of nonlinear least squares' problem, we can use this  $\rightarrow$

$$Ax = b$$

$$a_i^T x = b_i \quad [\text{linear}]$$

Non linear case  $\Rightarrow$

$$f_i(x_i) = b_i \quad \xrightarrow{\text{then:}} \quad f_1(x) = a_{00}x_1^2 + a_{11}x_2^2 + a_{22}x_3^2 = b_1$$

$$f_2(x) = a_{00}x_1^2 + a_{21}x_2^2 \dots = b_2$$

$$\tilde{f} = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{bmatrix}$$

Tradeoff between cost per iteration and no. of iterations.

Gradient Descent vs Stochastic Gradient Descent  $\Rightarrow$

[GD]

[SGD]

$$g(\bar{w}) \rightarrow \sum_{i=1}^N \phi(x_i, w)$$

$$w^{k+1} \leftarrow w^k - \eta \nabla g \quad \xrightarrow{\text{an estimate of } \nabla g \text{ is computed on a much smaller sample set.}}$$

$$\boxed{\nabla \tilde{g} = \sum_{i=1}^M}$$

where  $M \ll N$

$\rightarrow$  updating  $w$  is not a big deal but computation of  $\nabla g$  is computationally intensive

$\rightarrow$  so we are reducing the cost per iteration by estimating gradient from  $M$  samples. [where  $M \ll N$ ]

↓  
GD implemented as SGD.

$\rightarrow$  if samples are good, then with a smaller set, ~~we can~~ the process converges.

Issues with Gradient Descent -

- \* GD and SGD
- \* Optimality
- \* Convergence [at what rate will it converge?]
- \* Non-differentiable
- \* Higher order terms (like Hessian for finding learning rates)
- \* Practical issues

Newton's method with backtracking =>

for  $k = 1, \dots$

$$x^{k+1} \leftarrow x^k + s$$

if  $g(x^{k+1}) > g(x^k)$  → To avoid divergence or oscillations.

$$s = s/2$$

GD for ML →

Weights are learned

w

$$w^* = \underset{w}{\operatorname{argmin}} f(w) \quad - \text{No explicit form of } w \text{ that we are minimizing}$$

f is usually an expectation over a distribution, D.

$$\text{samples } n \leftarrow \left\{ \sum_{i=1}^n \phi(w_i, x_i)^T y_i \right\} \rightarrow \text{loss function}$$

$(x_i, y_i)$   
↑      ↑  
data    label:  
 $x_i \in D$

n samples coming from a distribution

$$w^{k+1} \leftarrow w^k - \eta \nabla f$$

finding the gradient

$\nabla f$

stochastic gradient

$\nabla f$  is computed from a random subset of the samples

Subgradient

- function might not be differentiable at all points.

[Ex: hinge loss]

computation - faster

- removes oscillations  
(because of the randomization introduced in the samples)

Instead of one gradient ( $\nabla f$ ), we compute a set of gradients  $\{\nabla y\}$

↳ finding the best out of all potential directions.

only for selected points

Minimising  $\rightarrow [y - xw]^T [y - xw] \rightarrow w^* = (\underbrace{x^T x}_{{d \times d} \text{ (usually high)}})^{-1} x^T y$

 $f(w) = \sum_{i=1}^n (y_i - x_i^T w)^2$ 

$w^{k+1} \leftarrow w^k - \eta \nabla f$

If this is done over a subset of samples  $\rightarrow$  SGD.

PCA:  $\max_w w^T \Sigma w \quad , \|w\| = 1$

$w \leftarrow$  eigen vectors of covariance matrix.

Update functions

$$\begin{aligned} \tilde{w} &\leftarrow w^k - \eta \nabla f \\ w^{k+1} &\leftarrow \frac{\tilde{w}}{\|\tilde{w}\|} \quad [\text{makes sure norm remains 1}] \end{aligned}$$

If  $n$  is very large, computed estimate will be closer to the real gradient

If  $n$  becomes very low  $\rightarrow$  it becomes noisy.

$$\begin{aligned} f(w+s) &= f(w^k) + s^T \nabla f \dots \quad \text{can be ignored} \\ &\approx f(w^k) - \eta \nabla f^T \nabla f \\ &\Rightarrow f(w^k) - \eta \quad [\text{kept small}] \rightarrow \text{not to jump too far} \\ f(w+s) &\leq f(w^k) \quad \Rightarrow f(w^{k+1}) \leq f(w^k) \end{aligned}$$

In every iteration, new value of the function is smaller than the current value.  
 $\sim$  converging

why do oscillations happen?  
 $\rightarrow$  Assumed that Taylor series approximation is accurate.

$f(w^{k+1})$   $\rightarrow$  estimation can be highly inaccurate.

$\rightarrow$  finding  $s$  becomes complex (when additional terms are used in the Taylor series)  
 [for quadratic - computing hessian inverse]  $\rightarrow$  though iterations become simpler.  
 $s = H^{-1} \nabla f$ .

## How to find $\eta$ ?

$\eta$  - has to be small  $\rightarrow$  good for convergence  
 but  
 but how small bad for computation time.

larger  $\eta \rightarrow$  risk of convergence and stability

Trade off between  $\Rightarrow$  computation/iter and #iterations.  
 stability and speed.

vary  $\eta$  over time

$\leftarrow \eta_t$  (initially make big jumps and when you come closer to the optima, reduce it slowly)

Instead of using a single scalar for all the dimensions -  $X$  not good as, so, we have to take a d dimensional  $\eta$  vector  
 Most popular.

Stable gradient descent  $\Rightarrow$  ADAGRAD  
 implementations

ADAGRAD still does  $\eta \nabla f$ , but  $\nabla f$  is modified as -

$$\text{Normalizing } \left\{ \nabla f^* = \frac{\nabla f^i}{\sqrt{g_i + \epsilon}} \right. \quad \begin{array}{l} \text{maintain a } g \text{ vector} \\ \text{and divide } \nabla f \text{ by } g \text{ vector.} \end{array}$$

this is same as to having diff gradients for diff dimensions.

$$g \rightarrow \text{cumulative gradient} ; \quad g_i = \sum_{i=1}^k (\nabla f^i)^2$$

what is the use of  $\epsilon \rightarrow$  when  $g^i \rightarrow 0$ , stability.

$$s = -H^{-1} \nabla f$$

$H_{d \times d}$  (computing inverse - is difficult)

$$s = -H^{-1} \nabla f$$

if  $H$  - diagonal —  $H^{-1}$  can be computed easily

"Hessian free" method — uses conjugate gradient