① 

$$\frac{z}{u} = Au$$

$$\frac{\partial u}{\partial u} = A$$

$$v_j = \{A\}_{jk} u_k$$
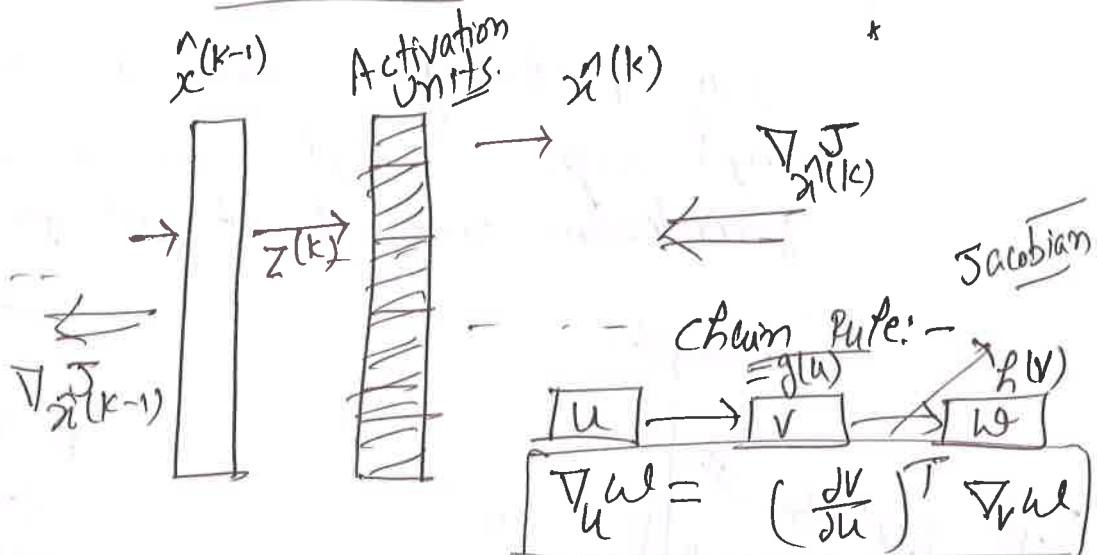
Deep Learning
Chain Rule

$$\left(\frac{\partial u}{\partial u}\right)_{jk} = \frac{\partial v_j}{\partial u_k}$$

$\hat{x}^{(k-1)}$  Activation units.  $x^{(k)}$

$z^{(k)}$

$\nabla_{\hat{x}^{(k)}} J$

$\nabla_{\hat{x}^{(k-1)}}$

Jacobian

Chain Rule:-

$u \xrightarrow{=g(u)} v \xrightarrow{f(v)} w$

$$\nabla_u w = \left(\frac{\partial v}{\partial u}\right)^T \nabla_v w$$

**Forward Computation:-**

Jacobian $\boxed{z^{(k)}} = W^{(k)T} \hat{x}^{(k)} + b^{(k)}$

$$\hat{x}^{(k)} = g^{(k)}\left(z^{(k)}\right)_i$$

Notes:-
(A)

$g(z^{(k)})$

$\boxed{z^{(k)}} - \boxed{\hat{x}^{(k)}} - \boxed{J}$

$\nabla_{\hat{x}^{(k)}}$

$\hat{W}^{(k)}$

$$\begin{bmatrix} \hat{W}^{(k)}_{:1} & - - & \hat{W}^{(k)}_{:n^{(k)}} \\ & & \\ & & \\ & & \end{bmatrix}$$

$$z^{(k)}_j = W^{(k)T}_{:j} \hat{x}^{(k-1)} + b^{(k)}$$

$\boxed{W^{(k)}}$ Linear Transform $z^{(k)} \longrightarrow J$

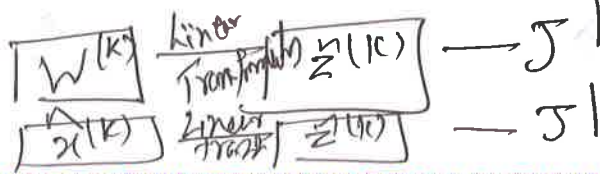$\boxed{\hat{x}^{(k)}}$ linear transf $z^{(k)} \longrightarrow J$

**Gradient Computation:-**

① $$\nabla_{z^{(k)}} J = \left(\frac{\partial \hat{x}^{(k)}}{\partial z^{(k)}}\right)^T \cdot \nabla_{\hat{x}^{(k)}} J$$

Diagonal matrix

$$= \cdot g'\left(z^{(k)}\right)$$

◎ $\nabla_{\hat{x}^{(k)}} J$

$$\frac{\partial J}{\partial W^{(k)}_{ij}} = \hat{x}^{(k+1)}_i \cdot \frac{\partial J}{\partial z^{(k)}_j}$$

② $$\nabla_{\hat{W}^{(k)}_{:j}} J = \nabla_{\hat{W}^{(k)}_{:j}} J$$

$$\frac{\partial J}{\partial \hat{W}^{(k)}_{ij}}$$

$z^{(k)}_j = (W^{(k)T})_{ji} \hat{x}^{(k-1)}_i + b^{(k)}_j$

$$\left(\nabla_{\hat{W}^{(k)}} J\right)_{ij} = W^{(k)} \hat{x}^{(k-1)}_i \left(\nabla_{z^{(k)}} J\right)_{ji}$$

$$\nabla_{b^{(k)}} J = 1 \cdot \nabla_{z^{(k)}} J$$

$$\nabla_{\hat{x}^{(k-1)}} J = W^{(k)T} \nabla_{z^{(k)}} J$$

Topic 2:-

## Regularization:-

Key & Challenge:- How to keep the model simple enough so that we d' have a good generalization error (& not just good training error).

Example:-

$$y = a_n x^n + a_{n-1} x^{n-1} + \cdots a_0 x^0$$

what is & the largest $n$, $a_{n-1}$

Polynomial Regression:-

A high degree polynomial $(n-1)$ fits $n$ points exactly.

Fitting a hyperplane in $n$-dimensions:-
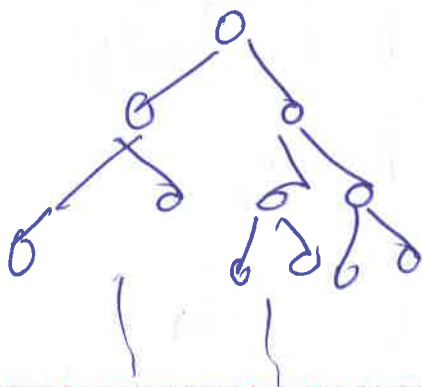
Regularizer:- ① Penalty for having higher degree $d$

② SVMs:-

$$\min_{w,b} \left[ \frac{1}{2} w^T w + c \sum_{i=1}^{m} \xi_i \right] \to \text{Loss}$$

$\Downarrow$
Regularizer    22-Norm

③ Decision Trees:-

A penalty on number of nodes in the tree

② For neural networks:-

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha \underbrace{\Omega(\theta)}_{\text{Regularizer}}$$

Regularized
Objective
Loss
Function

⇓
LOSS

$\alpha = 0$  No regularization

$$\Omega(\theta) = \frac{1}{2} W^T W \xrightarrow{\hspace{1cm}} \text{weight vector} \quad (\text{excluding bias term})$$

Note:- typically do not impose penalty on bias term $b$ (Does not capture correlation among

Norm of a vector - $W$:-

$$\boxed{P > 0}$$

p-norm:- $\qquad \left[ \sum_{j=1}^{n} \| W_j \|^P \right]^{1/p}$

Different kinds of effects of on parameters.

L2- norm:-        $P = 2$        $\| W \|_2^2$:-
                                   L2-norm square

L1- norm:-        $\sum_{j=2}^{n} |W_j|$

L0:- norm:-       # of non-zero components of $W$.

---

Why minimize $\| W \|^2$ ?

$$J(\theta) =$$

# $L^2$ Regularization :- Simplification — No bias parameter

$$\tilde{J}(\theta; X, y)$$

$$= J(\theta; X, y) + \frac{\alpha}{2} w^T w$$

$$\nabla_w \tilde{J}(w; X, y) = \nabla_w J(\theta; X, y) + \alpha w$$

$\Rightarrow$ Gradient Update Rule.

$$w^{new} = w - \eta \cdot \nabla_w \tilde{J}(\theta; X, y)$$

$$= w - \eta [\nabla_w J(\theta; X, y) + \alpha w]$$

$$= w - \eta \alpha w - \eta J(\theta; X, y)$$

$$\overset{2}{=} w(1 - \eta \alpha w) - \eta J(\theta; X, y)$$

$*$ Equivalent to ~~applying~~ ~~no~~ multiplying weight by $1 - \eta$ a constant

$(1 - \eta \alpha w)$ before applying non-regularized gradient update

# Another insight :-

$$Let \ w^* = \underset{w}{argmin} \ J(w)$$

Approximate $J(\theta)$ using a quadratic approximation around $w^*$.

③

Taylor approximation: $\theta(w,b)$

~~J approx~~
~~J app~~

$$J'(\theta) = J(w^*) + \nabla_w J(w^*)$$

Since $w^*$ is minimum $\to 0$

$$- \cdot (w - w^*) \, ,$$

$$+ \frac{1}{2}(w - w^*)^T H (w - w^*)$$

$\Downarrow$

$$= J(w^*) + \frac{1}{2}(w - w^*)^T H(w^*) (w - w^*)$$

Hessian matrix (Are semi-definite).

$$H = \begin{bmatrix} & j & \\ i & \frac{\partial J(\theta)}{\partial \theta_j \partial \theta_i} & \\ & & \end{bmatrix}$$

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

Quadratic approximation to a function $f: R \to R$ around $\underline{x^*}$.

Now, let w add the weight decay term to $J'(w)$

$\implies \nabla_w J'(w) = (w - w^*) H(w^*) (w - w^*)$

Adding regularization term & taking derivative:—

$$\nabla_w \tilde{J}'(w) = \cancel{(w - w^*)} + H(w - w^*)$$

$$H(\omega^*)(w - w^*)$$

$$+ \alpha \cdot \nabla_w \tfrac{1}{2} w^T w$$

$$= H(\omega^*)(w - w^*)$$

$$+ \alpha I w$$

Equating it to zero, we get:-

$$\cancel{H(w^*)w} \quad Hw \cancel{- Hw^*} + \alpha I w = H w^*$$

$$\boxed{w(H + \alpha I) = H w^*}$$

$$\boxed{w = (H + I\alpha)^{-1} H w^*}$$

if $\alpha = 0$, $\quad \boxed{w = w^*}$ $\qquad \boxed{Q Q^T = I}$

$H$: real & symmetric $\xrightarrow{}$ orthonormal basis

$$H = Q \Lambda Q^T \qquad \text{(Eigenvalue decomposition)}$$
$$\hookrightarrow \text{Diagonal}$$

$$w = (Q \Lambda Q^T) + \alpha I)^{-1} Q \Lambda Q^T w^*$$

$$= [Q \Lambda Q^T + Q \alpha I Q^T]^{-1} Q \Lambda Q^T w^*$$

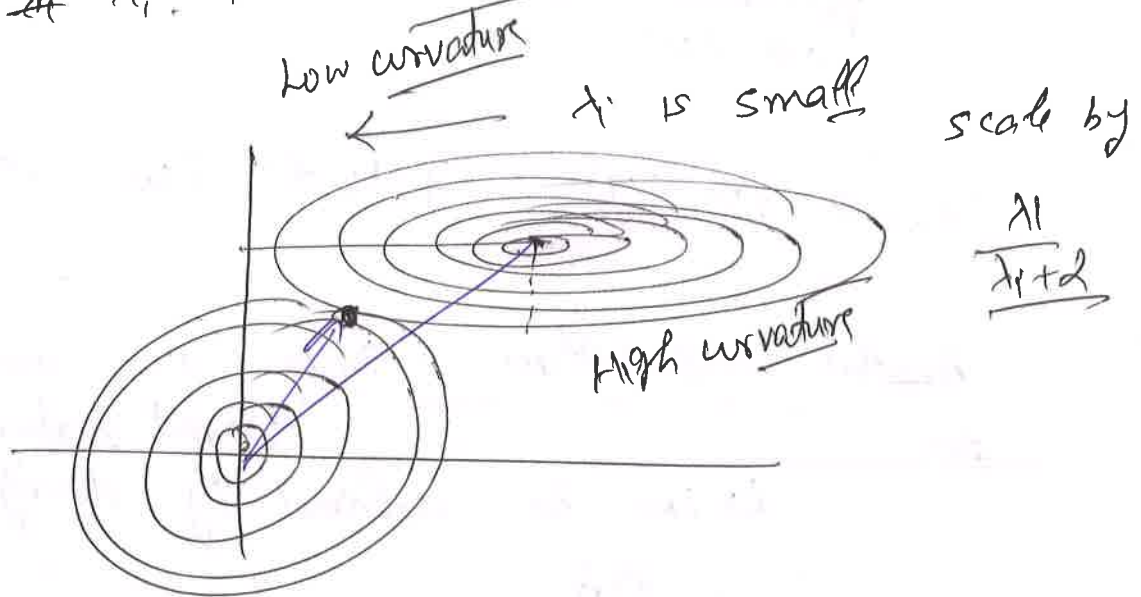$$= (Q(\Lambda + \alpha I) Q^T)^{-1} (Q \Lambda Q^T) w^*$$

$$= \boxed{Q (\Lambda + \alpha I)^{-1} \Lambda Q^T w^*}$$

Each eigenvalue $\lambda_i$ is scaled by $\frac{\lambda_i}{\lambda_i}$

scaling each direction (i-th eigenvetol) by $\frac{\lambda_i}{\lambda_i + \alpha}$

$\frac{\beta}{f} \cdot \frac{1}{4i}$   $\lambda_i$: - curvature.

Low curvature

$\lambda_i$ is small   scale by

$\frac{\lambda_i}{\lambda_i + \alpha}$

High curvature



## $L1$:- Regularization :-

$$\Lambda(\theta) = ||w||_1 = \sum_{i=1}^{A} |w_i|$$

$$\tilde{J}(w; x, y) = J(w; x, y) + \alpha ||w||_1$$

$$\nabla_w \tilde{J}(w; x, y) = \nabla_w J(w; x, y) + \alpha \, \text{sign}(w)$$

sub-gradient

Gradient update

old $\Big[$ ⇒ $w \leftarrow w - \eta \cdot \nabla_w J(w; x, y)$

new
(regularization) $\Big[$   $w \leftarrow w - \eta \left[ \nabla_w J(w; x, y) \pm v \right]$

Sign of $w$

vector of $1's$ & $-1's$
(or possibly zeros)
$v = \text{sign}(w)$

May not get a closed form

$$J'(w; X, y) = J(w^*, X, y) + (w - w^*)^T \underbrace{\nabla J(w^*; X, y)}_{\to 0}$$

$\Downarrow$ Quadratic approximation around $w^*$

$$+ (w - w^*)^T \tfrac{1}{2} H(w^*) (w - w^*)$$

$$\nabla_w J'(w; X, y) = H(w^*)(w - w^*)$$

Another **assumption** - Data is un-correlated
(input features)
$\hookrightarrow$ can be achieved by doing linear
PCA.

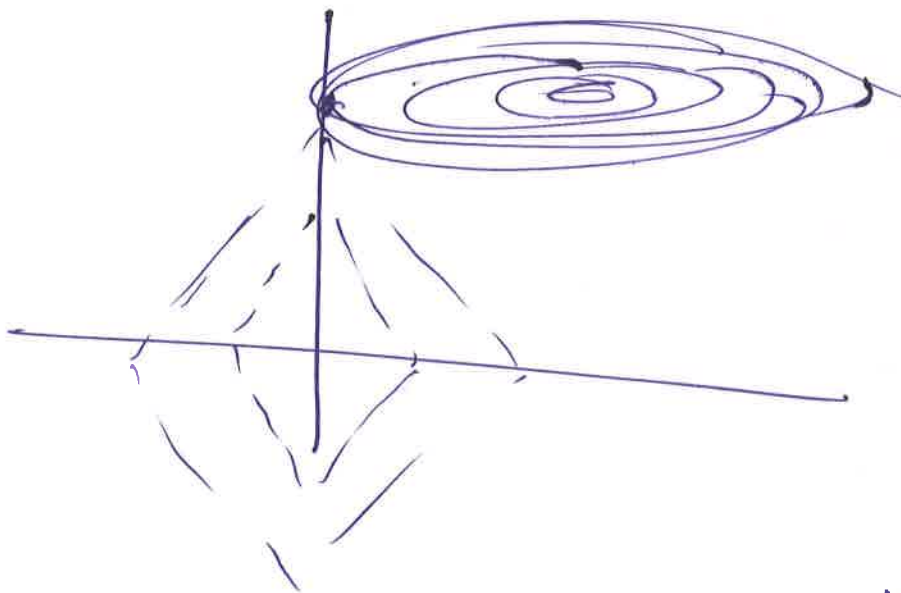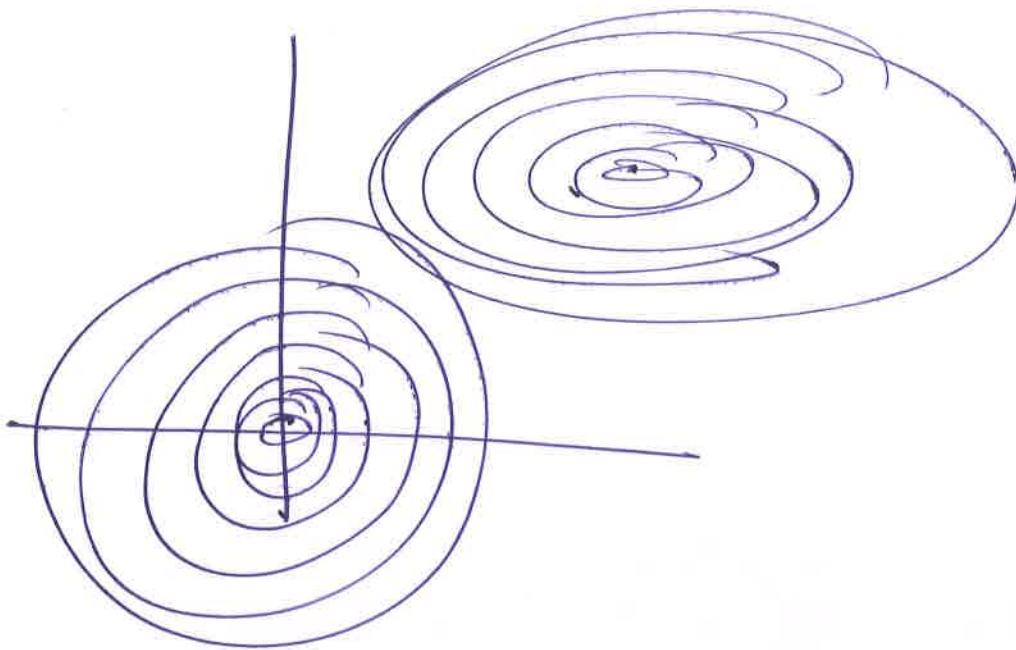$$H = \begin{bmatrix} H_{11} & & & & \\ & H_{22} & & & \\ & & H_{32} & & \\ & & & \ddots & \\ & & & & H_{n,n} \end{bmatrix}$$

$\Rightarrow$

$$\nabla_w \tilde{J}'(w; X, y) = H(w - w^*) + \alpha \cdot \text{sgn} \quad \alpha \sum_{i=1}^{n} |w_i|$$

$\Rightarrow$

$$\nabla_w \tilde{J}'(w; X, y) = H_{i,i}(w_i - w_i^*) + \alpha \xi_i \cdot \text{sgn}(w_i)$$

Equating to zero:

$w_i^* \geq 0$

$$w_i - w_i^* = \frac{\alpha \cdot \text{sign}(w)}{H_{i,i}}$$

$$w_i = \frac{w_i^* + \frac{\alpha}{H_{i,i}}}$$

$$w_i = w_i^* + \frac{\alpha}{H_{i,i}}$$

$$w_i = \text{sign}(w_i^*) \max\left\{ |w_i^*| - \frac{\alpha}{H_{i,i}}, 0 \right\}$$