

①

COL 865

Deep Learning

August 14, 2017

Last class:-

(a) Backpropagation

(b) Regularization (means to avoid overfitting and/or improve generalization error).

Some Linear Algebra:-~~Matrix~~ 12Let $A \in \mathbb{R}^{n \times n}$ symmetric.(a) An eigenvector $v \in \mathbb{R}^n$ of A is defined as:Right
eigenvector

$$A v = \lambda v$$

eigenvalue

A square matrix $A \in \mathbb{R}^{n \times n}$ has n such eigenvalues.

(b) A (real & symmetric) matrix can be decomposed as:-

$$A = Q \Lambda Q^T$$

where:-

 Q :- Orthogonal matrix $\rightarrow Q_{:i} \rightarrow i^{th}$ column of Q Λ :-

Diagonal matrix

composed of columns of A eigenvectors of A st $\Lambda_{ii} = i^{th}$ eigenvalueof A (Corresponding to corresponding column of Q).

$$Q_{:i} \cdot Q_{:j} = 0 \quad i \neq j$$

$$\text{Similarly } Q_{:i} \cdot Q_{:j} = 0 \quad i \neq j$$

Rows & columns of Q are \perp to each other

$$\text{Further, } \|Q_{:i}\|^2 = 1$$

$$\|Q_{:j}\|^2 = 1$$

Then:-

Let $\begin{bmatrix} v^{(1)} \\ v^{(2)} \\ \vdots \end{bmatrix}$ be an eigenvector of A .
 $v^{(2)}$ be an eigenvector of A .

columns of A

Property
 $AT = T^{-1}$

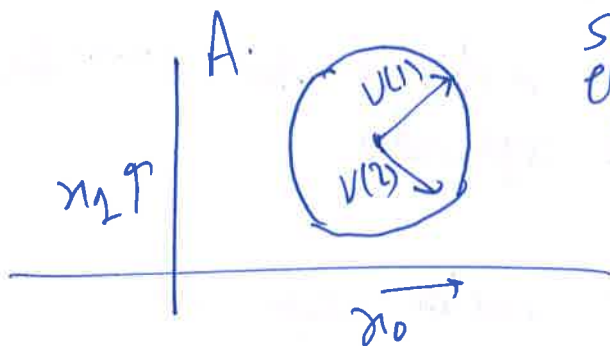
$$\Rightarrow AV^{(1)} = \lambda V^{(1)}$$

$$\Rightarrow AV^{(2)} = \lambda V^{(2)}$$

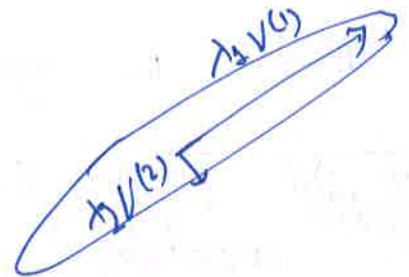
But $v^{(1)}, v^{(2)} \dots v^{(n)}$
 span the entire space

$$AV^{(n)} = \lambda V^{(n)}$$

Q - orthonormal



scaling the entire space
 \Rightarrow



Now, let us revisit L2-Norm:-

~~Assumption~~ - No

Question:- $R(\theta) = \frac{1}{2} W^T W \quad \hat{= (L2-Norm)^2}$

Why to ignore bias parameters. Not very useful to regularize

(W requires large amount of data)

$$y = W^T x + b \quad Z = W^T x$$

$$Z = W^T x + b$$

Does not capture correlation among variables interact

Not difficult to learn (less data required)

only H-1 (output)

②

Model No bias: - Only for exposition.

① Quadratic approximation to a function: -

approximation $\Rightarrow f'(w) = f(w^*) + f'(w^*)(w-w^*) + \frac{f''(w^*)}{2}(w-w^*)^2$ $x \approx x^*$
 $f: \mathbb{R} \rightarrow \mathbb{R}$

Now, suppose: - $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$f'(w) = f(w^*) + \nabla_w f(w^*)^T (w-w^*) + \frac{1}{2} (w-w^*)^T H (w-w^*)$$

Hessian evaluated at w^*

~~$\Rightarrow J'(w) =$~~

* Consider $J(w)$:- error function

Let (w^*) be a minima of $J(w)$

(1) Let w be a point close to w^* .

Note: - (1) $\nabla_w J(w^*) = 0$

(2) H is (symmetric) true semi-definite

(at w^*): - since function should be concave at w^*

$$\Rightarrow J'(w) = J(w^*) + \nabla_w J(w^*)^T (w-w^*) + \frac{1}{2} (w-w^*)^T H (w-w^*)$$

$-2(w) = \frac{1}{2} w^T w$

~~We want to minimize~~

Further:-

$$\tilde{J}(w) = J(w) + \frac{1}{2} w^T w$$

$$\Rightarrow \tilde{J}'(w) = J'(w) + \frac{1}{2} w^T w$$

Quadratic Approx.

Now, we want to minimize $\tilde{J}(\omega)$
 \Rightarrow minimize $\tilde{J}'(\omega)$

$$\Rightarrow \nabla_{\omega} \tilde{J}'(\omega) = 0$$

$$\Rightarrow \nabla_{\omega} \tilde{J}(\omega) + \frac{\alpha}{2} \omega = 0$$

$$\Rightarrow \nabla_{\omega} \tilde{J}(\omega) = -\frac{\alpha}{2} \omega$$

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

$$H(\tilde{\omega} - \omega^*) + \alpha \tilde{\omega} = 0$$

$$\Rightarrow H \tilde{\omega} + \alpha I \tilde{\omega} = H \omega^*$$

$$\Rightarrow (H + \alpha I) \tilde{\omega} = H \omega^*$$

$$\Rightarrow \tilde{\omega} = (H + \alpha I)^{-1} H \omega^*$$

Now:- H is symmetric.

$$H = Q \Lambda Q^T \rightarrow \text{Orthogonal}$$

\Rightarrow if H is diagonal -

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

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

$$= \underbrace{[Q (\Lambda + \alpha I)^{-1} \Lambda Q^T]}_{H'} \omega^*$$

with eigenvector same as H & eigenvalues $\frac{\lambda_i}{\lambda_i + \alpha}$ if λ_i is an eigenvalue of H .

Now re-scales the w space by $\frac{\lambda_i}{\lambda_i + 2}$ along each dimension corresponding to eigenvector of H

$\Rightarrow w^{\text{new}}$ is w^* re-scaled on the re-scaled space

Value of w^* is shrunk by a factor of $\frac{\lambda_i}{\lambda_i + 2}$ in each direction. \rightarrow constant

if λ_i is large less shrinking.

λ_i is large: - High

second order derivative.

High curvature \Rightarrow less shrinkage.

$J(w)$ changes very fast in direction corresponding to λ_i . can't afford to have ~~large regularization~~ much regularization.

L2 - Regularization -

Consider quadratic approximation as before.

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

$$-2(w) = \frac{1}{2} \|w\|_2^2$$

(regularized)

$$\tilde{J}^q(w) = J^q(w) + \|w\|_2^2$$

$$\nabla_w \tilde{J}^q(w) = \nabla_w J^q(w) + \nabla_w \|w\|_2^2$$

$$\nabla_w \tilde{J}^q(w) = (\tilde{w}_j - w_j^*) H_{jj} + \alpha \text{sign}(\tilde{w}_j)$$

\Rightarrow Equating gradient to zero, we get

$$(\tilde{w} + w^*)^T H = \alpha \nabla_w \|w\|_2^2$$

Assume H is diagonal
Taking componentwise derivatives $H_{jj} > 0$

$$[\tilde{w}_j + w_j^*] H_{jj} = \alpha \cdot \text{sign}(\tilde{w}_j)$$

$$\tilde{w}_j + w_j^* = \frac{\alpha}{H_{jj}} \text{sign}(\tilde{w}_j)$$

then first-

$$|\tilde{w}_j| \leq |w_j^*|$$

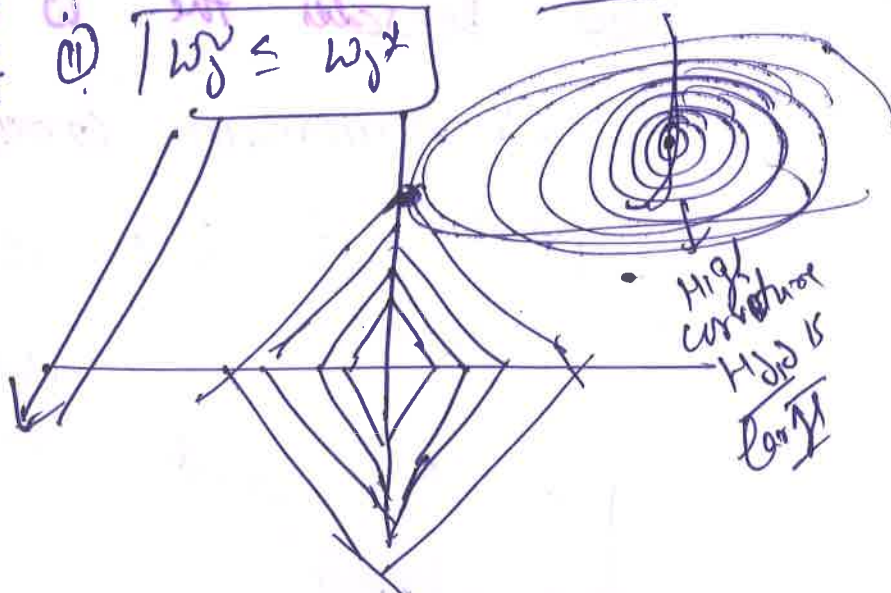
:- not

$$\textcircled{1} \text{sign}(\tilde{w}_j) = \text{sign}(w_j^*)$$

$$\textcircled{2} |\tilde{w}_j| \leq |w_j^*|$$

$$\text{if } w_j^* \geq 0 \\ \Rightarrow 0 \leq w_j \leq w_j^*$$

$$\text{if } w_j^* \leq 0 \\ \Rightarrow w_j^* \leq w_j \leq 0$$



Further $w_j^* \geq 0$

$$w_j = \text{sign}(w_j^*) \max \left\{ |w_j^*| - \frac{\alpha}{H_j}, 0 \right\} \quad \downarrow \text{induces sparsity}$$

$$w_j^* \geq 0 :-$$

$$\text{if } w_j^* > \frac{\alpha}{H_j} \Rightarrow \tilde{w}_j = w_j^* - \frac{\alpha}{H_j}$$

Non-zero

else $w_j = 0$

\Rightarrow L1 Norm induces sparsity.

as opposed to:-

$$\tilde{w}_j = w_j^*$$

Feature selection mechanism

LASSO:-

Linear Model + Least Squares

objective + L1-Penalty

Least Absolute Shrinkage & Selection Operator

④

Other variations:-

1.2 - 7.8 Read

Constrained Optimization:-

① $\boxed{-R(\theta) \leq k}$ Enforce Explicitly

$W^T W \leq k$ rather than having a soft penalty

Solve using method of Lagrangians (Dual)

or solving using gradient descent along with projection of

2) $\boxed{W^T W = k}$ with radius \sqrt{k}

argmin $J(\theta)$

$R(\theta) \leq k$

How to solve it?

Under Constrained Problems

X :- design matrix

$X^T X$:-

if k need to invert

$X^T X$

Morse-Penrose pseudoinverse

Invert:-

$(X^T X + \alpha I)$

$(X^T X + \alpha I)^{-1} X$

Linear Regression

$W =$

$(X^T X)^{-1} X^T y$

$(X^T X + \alpha I)^{-1} X^T y$

equivalent to L2 regularization

Linear Regression with quadratic loss

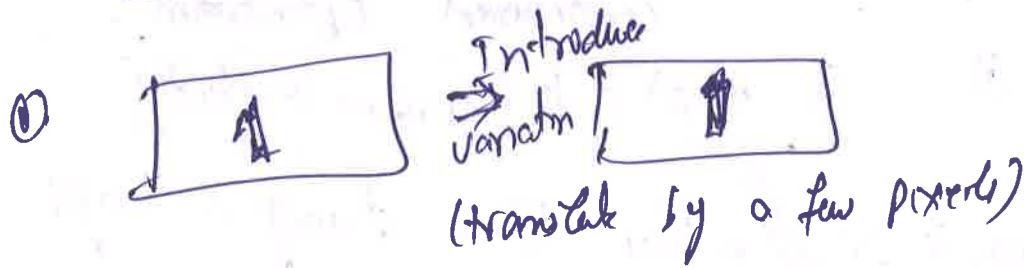
$\tilde{J}(W) = J(W) + \frac{\alpha}{2} W^T W$

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

$W = (X^T X + \alpha I)^{-1} X^T y$

$J(W) = \frac{1}{2} (X^T W - y)^T (X^T W - y)$

Dataset Augmentation



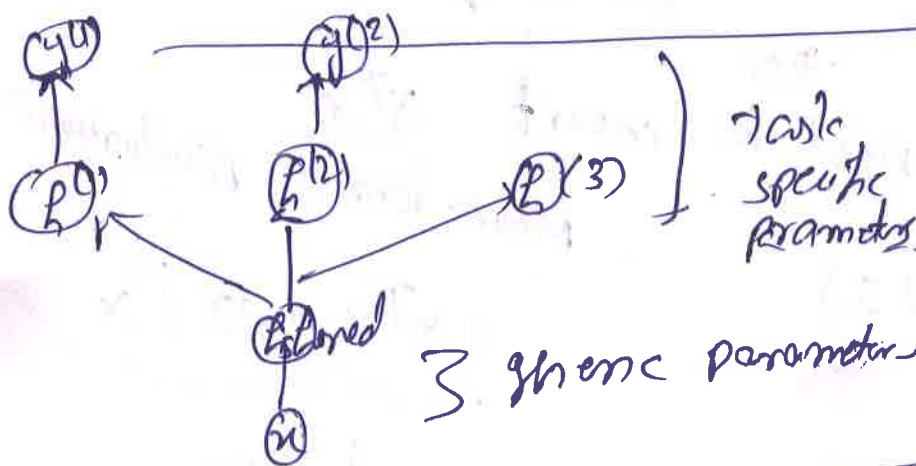
② explicitly introduce noise (random noise)
 ↳ can apply to encoder units as well

Noise Introduction in weights :-

$$W \rightarrow W + \epsilon W \quad \text{for each example, we randomly perturbed weights}$$

$$\epsilon W \sim N(\epsilon; 0, \eta I)$$

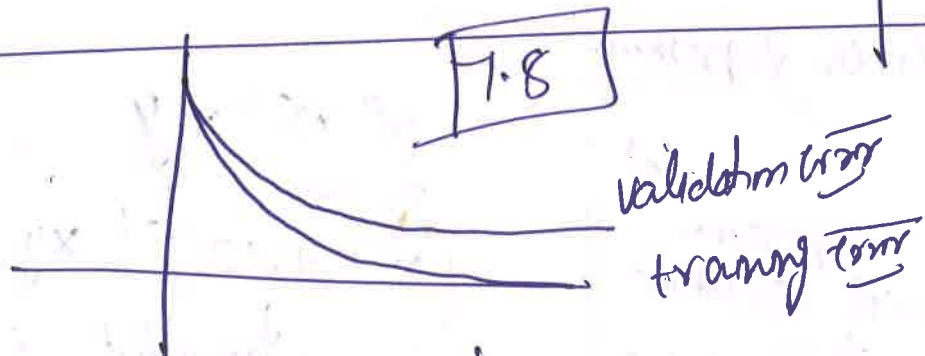
↓
stability



Multitask Learning from

① Ensembles
 ② Parameter Tuning

Early Stopping :-



$$\alpha = \frac{1}{\tau \epsilon}$$

$$\alpha = \frac{1}{\tau \eta} \Rightarrow \text{learning rate}$$

iteration