



* Principle Component Analysis

lets say,

	x	y	z	
1	1	1	1	\vec{x}_1
0.5	0	0	0	\vec{x}_2
0.25	1	1	1	\vec{x}_3
0.35	2.5	1.5	1.5	\vec{x}_4
.
.
.
.	.	.	.	\vec{x}_m

Here x & y are highly correlated thus if any one is ignored it wouldn't make much difference in result.

Here if in following operation when $y_i = \bar{y} + \hat{y}_i$ then after updating all y_i if we calculate \bar{y}_i then if it is $\bar{y}_i = 0$ then it is called 0 mean.

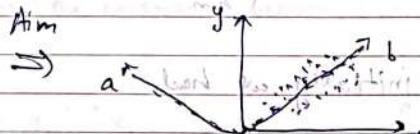
$$P_{yz} = \sum_{i=1}^n (y_i - \bar{y})(z_i - \bar{z})$$

$$(x,y,z) \rightarrow \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2} \times \sqrt{\sum_{i=1}^n (z_i - \bar{z})^2}$$

Notes



- In this first figure the data points have notable variance across both axis



- Aim is to find basis in the dimension such that dimension can be reduced

→ Here in 2nd figure a & b are independent vectors that can form entire plane. But unlike x,y variance of points is significantly affected by b and it can be safely ignored. Thus we reduce dimensions from newly formed vectors

→ Initially make data or mean & unit variance by Z-Score normalization

Now, representing x_i using basis P

if $x_i = \alpha_{i,1}p_1 + \alpha_{i,2}p_2 + \dots + \alpha_{i,n}p_n$

$$\text{then } x_i = \alpha_{i,1}p_1 + \alpha_{i,2}p_2 + \dots + \alpha_{i,n}p_n$$

Using orthonormal basis

for orthonormal basis

orthogonal projection of x_i onto

$$\Rightarrow \alpha_{i,j} = x_i^T p_j$$

if no coordinate space basis

use \hat{x}_i standard vec. p_j

$$x_i = \alpha_{i,1}p_1 + \alpha_{i,2}p_2 + \dots + \alpha_{i,n}p_n \quad \left\{ \begin{array}{l} \text{Dimension} = (1 \times n) \\ \text{row} \end{array} \right.$$

now \hat{x}_i has n entries

$$\begin{bmatrix} 1 \\ \vdots \\ n \end{bmatrix} \begin{bmatrix} \alpha_{i,1} & \alpha_{i,2} & \dots & \alpha_{i,n} \end{bmatrix}$$

Now for m vectors

$$\begin{pmatrix} \hat{x}_1 \\ \hat{x}_2 \\ \vdots \\ \hat{x}_m \end{pmatrix} \xrightarrow{\text{Dimension}} (m \times n)$$

initially we had

$$\begin{array}{ccc} \xrightarrow{\text{Dimension}} & \xrightarrow{\text{no cov.}} & \xrightarrow{\text{new cov.}} \\ \begin{pmatrix} \hat{x}_1 \\ \hat{x}_2 \\ \vdots \\ \hat{x}_m \end{pmatrix} & \xrightarrow{\text{Dimension}} & \begin{pmatrix} \alpha_{1,1} & \alpha_{1,2} & \dots & \alpha_{1,n} \\ \alpha_{2,1} & \alpha_{2,2} & \dots & \alpha_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_{m,1} & \alpha_{m,2} & \dots & \alpha_{m,n} \end{pmatrix} \end{array}$$

Now here we've got new dimension. Here

we want to prove that these new dimensions have low covariance i.e. they are not interrelated

Theorem

(1)

If X is a matrix such that its columns have zero mean and if $\tilde{X} = X P$ then the columns of \tilde{X} will also have zero mean.

$$\mathbf{q} \mathbf{x}^T (\mathbf{x}^T \mathbf{q}) \perp = \mathbf{q} \mathbf{x}^T (\mathbf{q} \mathbf{x}) \perp = \mathbf{x}^T \mathbf{q} \perp = 0.$$

Proof →

$$\mathbf{1}^T \mathbf{x} = 0 \text{ since columns are zero mean}$$

$$\text{now } \tilde{\mathbf{x}} = \mathbf{x} \mathbf{P} \Rightarrow \mathbf{x}^T \mathbf{y} + \mathbf{1}^T \mathbf{x}^T \mathbf{q} = \mathbf{1}^T \mathbf{x} \mathbf{P}$$

$$= 0 \mathbf{P}$$

(2)

Theorem

$\mathbf{x}^T \mathbf{x}$ is symmetric because $\mathbf{x}^T \mathbf{x} \perp$ is zero

Proof → $(\mathbf{x}^T \mathbf{x})^T = (\mathbf{x}^T (\mathbf{x}^T \mathbf{x}))^T = \mathbf{x}^T \mathbf{x} \perp$ means it is true

→ If X is a matrix whose columns are zero mean then $\frac{1}{m} \mathbf{x}^T \mathbf{x}$ is the covariance matrix. This means that each entry C_{ij} stores the covariance between column i & j of X .

Explanation

$$C_{ij} = \frac{1}{m} \sum_{k=1}^m (x_{ki} - \mu_i)(x_{kj} - \mu_j) \rightarrow \text{covariance}$$

$$= \frac{1}{m} \sum_{k=1}^m x_{ki} x_{kj} \quad (\because \mu_i = \mu_j = 0)$$

$$= \frac{1}{m} \mathbf{x}_i^T \mathbf{x}_j \text{ or resp. entry } C_{ij} \text{ after } \mathbf{x}^T \mathbf{x} = \mathbf{q}^T \mathbf{q}$$

Similarly matrix is formed

$$\frac{1}{m} (\mathbf{x}^T \mathbf{x})_{ij}$$

$\frac{1}{m} \hat{X}^T \hat{X}$ → Covariance matrix of transformed data

$\Sigma = \frac{1}{m} X^T X$ → Covariance matrix of original data

DELUXE
PAGE NO.:
DATE:

CloudTech

Now,

$$\hat{X} = X P$$

and as seen earlier $\frac{1}{m} \hat{X}^T \hat{X}$ is covariance matrix

$$\therefore \frac{1}{m} \hat{X}^T \hat{X} = \frac{1}{m} (X P)^T X P = \frac{1}{m} (P^T)^T X P$$

cancel out P and P^T to get $\Sigma = P^T \Sigma P$

$$Q = P^T \left(\frac{1}{m} X^T X \right) P \quad Q = \Sigma$$

$$\frac{1}{m} \hat{X}^T \hat{X} = P^T \Sigma P$$

Here, in $\frac{1}{m} \hat{X}^T \hat{X}$ is covariance matrix where we

want all i, j where $i \neq j$ "0" since if $i=j$ it is covariance of i & j column if $i=j$ it is variance

$\begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix}$ is Diagonal matrix

Now, as seen earlier

$\frac{1}{m} \hat{X}^T \hat{X}$ is symmetric

Theorem 2

$$\text{Thus } P^T = P^{-1}$$

(P is collection of eigen vectors of $\frac{1}{m} \hat{X}^T \hat{X}$)

and $P^T \Sigma P = \text{Diagonalization of } \Sigma \text{ gives Diagonal matrix}$



The method of transforming data to new basis with

- (1) low covariance
- (2) high variance

- proved

- not proved

is called PCA

→

$$n_i = \sum_{j=1}^n \alpha_{ij} p_j$$

Now, if we approximate by taking only k dimensions out of n then we can reduce dimensions

$$\hat{n}_i = \sum_{j=1}^k \alpha_{ij} p_k$$

So we select this k such that

$$S_e^2 = \sum_{i=1}^m (n_i - \hat{n}_i)^T (n_i - \hat{n}_i)$$

$$= \sum_{i=1}^m \left(\sum_{j=1}^n \alpha_{ij} p_j - \sum_{j=1}^k \alpha_{ij} p_j \right)^2$$

all dimensions k dimensions

$$= \sum_{i=1}^m \left(\sum_{j=k+1}^n \alpha_{ij} p_j \right)^2 = \sum_{i=1}^m \left(\sum_{j=k+1}^n \alpha_{ij} p_j \right)^T \left(\sum_{j=k+1}^n \alpha_{ij} p_j \right)$$

$$= \sum_{i=1}^m \sum_{j=k+1}^n \alpha_{ij} p_j^T p_j \alpha_{ij}$$

$$= \sum_{i=1}^m \sum_{j=k+1}^n \underbrace{\alpha_{ij} p_j}_{\text{vector}}^T \underbrace{\alpha_{ij}}_{\text{scalar}}$$

Given α_{ij} and α_{ij}^2 are probabilities of moving from i to j in last time slot.

$$= \sum_{i=1}^m \sum_{j=k+1}^n \alpha_{ij}^2 \quad (\because p_j^T p_j = 1, p_i^T p_j = 0 \forall i \neq j)$$

$$= \sum_{i=1}^m \sum_{j=k+1}^n (u_i^T p_j)^2$$

$$= \sum_{i=1}^m \sum_{j=k+1}^n (p_j^T u_i) (u_i^T p_j)$$

Now we want $\sum_{j=k+1}^m u_i^T p_j (\sum_{j=k+1}^m u_i^T p_j)$ p_{ji} will be zero if $i \neq j$

Now

lets say

$$\sum_{i=1}^m \begin{bmatrix} u_{i1} \\ u_{i2} \\ u_{i3} \end{bmatrix} \begin{bmatrix} p_{j1} \\ p_{j2} \\ p_{j3} \end{bmatrix}^T = \sum_{i=1}^m \begin{bmatrix} u_{i1} \\ u_{i2} \\ u_{i3} \end{bmatrix} \begin{bmatrix} u_{j1} & u_{j2} & u_{j3} \end{bmatrix}$$

$$= \begin{bmatrix} (u_{11}^2 + u_{12}^2 + u_{13}^2) & u_{11}u_{12} & u_{11}u_{13} \\ u_{12}u_{11} & (u_{12}^2 + u_{13}^2) & u_{12}u_{13} \\ u_{13}u_{11} & u_{13}u_{12} & (u_{13}^2 + u_{11}^2) \end{bmatrix}^T \begin{bmatrix} u_{21}^2 & u_{21}u_{22} & u_{21}u_{23} \\ u_{22}u_{21} & (u_{22}^2 + u_{23}^2) & u_{22}u_{23} \\ u_{23}u_{21} & u_{23}u_{22} & (u_{23}^2 + u_{21}^2) \end{bmatrix}$$

$$(u_{11}^2 + u_{12}^2 + u_{13}^2) (u_{21}^2 + u_{22}^2 + u_{23}^2) + \dots = (u_{11}^2 + u_{12}^2 + u_{13}^2) \dots$$

$$(u_{11}^2 + u_{12}^2 + u_{13}^2) \left[u_{11}^2 + u_{12}^2 + \dots + u_{m1}^2 + u_{11}u_{12} + u_{11}u_{13} + \dots + u_{12}u_{13} \right] = u_{11}^2 + u_{12}^2 + u_{13}^2$$

The Matrix $\sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T$ is the covariance matrix

$$\begin{pmatrix} \mathbf{x}_{11} & \mathbf{x}_{12} \\ \mathbf{x}_{21} & \mathbf{x}_{22} \\ \mathbf{x}_{31} & \mathbf{x}_{32} \end{pmatrix} + \begin{pmatrix} \mathbf{x}_{11}^2 & \mathbf{x}_{11} \mathbf{x}_{12} \\ \mathbf{x}_{21} \mathbf{x}_{11} & \mathbf{x}_{22}^2 \\ \mathbf{x}_{31} \mathbf{x}_{11} & \mathbf{x}_{32} \mathbf{x}_{11} \end{pmatrix} + \begin{pmatrix} \mathbf{x}_{11} \mathbf{x}_{21} & \mathbf{x}_{11} \mathbf{x}_{31} \\ \mathbf{x}_{21} \mathbf{x}_{11} & \mathbf{x}_{22} \mathbf{x}_{21} \\ \mathbf{x}_{31} \mathbf{x}_{11} & \mathbf{x}_{32} \mathbf{x}_{11} \end{pmatrix} + \dots = \mathbf{C}$$

column column column column

Now

the $\sum_{i=1}^m p_j^T m c p_j$ is Norm of $\left[\frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T \right] \mathbf{c} = \frac{\mathbf{c}^T \mathbf{c}}{m}$

p_j is eigenvectors of covariance matrix $\frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T$

$$c = \sum_{j=k+1}^n p_j^T m c p_j$$

(p_1, p_2, \dots, p_k are eigenvectors of covariance matrix) signified $A^{-1} \mathbf{c}$

$\min_{j=k+1}^n (p_j^T m c p_j)$ subject to $p_j^T p_j = 1$

Perp P_{k+1}, \dots, P_n

orthogonal eigenvectors of covariance matrix \rightarrow theorem

Deep learning book (7)

p_j are smallest eigen vectors

→ How to prove high variance - ② condition as mentioned earlier

mean μ non-zero is zero \rightarrow zero mean \rightarrow variance

$$\mathbf{x}_i = \mathbf{x}_i - \mu$$

variance

$$\frac{(\hat{\mathbf{x}}_i - \mu_i)^T (\hat{\mathbf{x}}_i - \mu_i)}{m} = \frac{\hat{\mathbf{x}}_i^T (\hat{\mathbf{x}}_i - \mu_i)}{m} \quad (\because \mu_i = 0)$$

zero mean

$$\begin{aligned} \text{where } \pi_i &= \frac{1}{m} p_i^T (X^T X p_i) \quad \rightarrow A_{ii} = \lambda_i \\ &= \frac{1}{m} p_i^T (\Lambda_i p_i) \\ &= \frac{\lambda_i}{m} \underbrace{p_i^T p_i}_1 \\ &= \frac{\lambda_i}{m} \quad \} \text{ eigen values} \end{aligned}$$

If we remove the small eigen values then we are left with large eigen values & large eigen value means high variance.

→ PCA Example (refer notes/week 6/1/PCA example.png)

→ Singular Value Decomposition (SVD)

Eigen values & vectors work only for square symmetric matrix

But when matrix are rectangular then we use SVD

Let's say A - rectangular matrix of size $n \times m$
matrix A transforms a vector from n dimension to m dimension

$$\text{i.e. } A v' = v''$$

\downarrow

$n \times n$

\downarrow

$m \times 1$

$(n \times 3) \times (m \times 3)$

Now (let us assume) that $\{v_i\}$ are such that

($n \times 1$) v_i is linearly independent & non-zero.

$\{v_i\}$ is basis of vector space \mathbb{R}^n

($m \times 1$) u_i is basis of vector space \mathbb{R}^m

then,

$$Av_i = \sigma_i u_i$$

$Av_i = \sigma_i u_i$ can be possible

$$[v_i = v_i] \Rightarrow v_i = v_i \sigma_i u_i$$

Now v_i is a vector such that $v_i = \sum_{i=1}^k \alpha_i v_i$

$$[v_i = v_i] \Rightarrow v_i = v_i \sum_{i=1}^k \alpha_i u_i$$

\therefore

$$Av_i = \sum_{i=1}^k \alpha_i A v_i \quad v_i \in \text{column space of } A \Rightarrow (n \times 1) \text{ vectors}$$

$$\Rightarrow \sum_{i=1}^k \alpha_i \sigma_i u_i \Rightarrow u_i \in (m \times 1) \text{ vectors}$$

$$(v_i = v_i) \Rightarrow (v_i = v_i) = A^T A$$

$$[v_i = v_i] \Rightarrow [v_i = v_i]$$

Here k is used because there are k linearly independent vectors in rowspace & columnspace of A for these k vectors as input $A^T A$ is non-zero

$$A^T A = A^T A \quad \text{columns}$$

Since $v_i \perp u_i$ common term k is used in summation

$n - \text{dim}$

$m - \text{dim}$

→ Thus we can write

$$A_{m \times n} \times V_{n \times k} = \begin{bmatrix} v_1 & v_2 & \dots & v_k \end{bmatrix} = \sum_{i=1}^k u_i \sigma_i u_i^T$$

\downarrow $n \times 1$ vectors \downarrow $m \times 1$ vectors \downarrow diagonal matrix

\downarrow $n \times k$ matrix

Now there are $(n \times k)$ & $(m \times k)$ vectors to solve
 but we need $n \times n$ & $m \times m$ vectors
 so by (Gram-Schmidt orthogonalization) we
 can find $(n-k)$ & $(m-k)$ orthogonal vectors respectively

$$\therefore A_{m \times n} v_{n \times n} = U_{m \times m} \Sigma_{m \times n}$$

Thus $U^T A V = \Sigma$ [$U^{-1} = U^T$]

$$A = U \Sigma V^T$$
 [$V^{-1} = V^T$]

→ Now we find $U \leq V$ & $A = U \Sigma V^T$

let say U, V, U^T, V^T exist

$$\begin{aligned} \therefore A^T A &= (U^T \Sigma V)^T (U^T \Sigma V) \\ &= V \Sigma^T U^T U \Sigma V^T \end{aligned}$$

Similarly $A^T A^T = U^T \Sigma^2 U$ { U, V are eigenvectors of $A^T A$ & $A^T A^T$ respectively }

$$\begin{aligned} A_{m \times n} &= \left[\begin{array}{c|c|c} 1 & 1 & 1 \\ \hline u_1 & u_2 & \dots & u_k \\ \hline 1 & 1 & 1 \end{array} \right] = \left[\begin{array}{c|c|c} \sigma_1 & & \\ \hline \vdots & \ddots & \\ \hline & & \sigma_k \end{array} \right] \times \left[\begin{array}{c|c|c} v_1 & & \\ \hline \vdots & \ddots & \\ \hline & & v_k \end{array} \right] \\ &= \sum_{i=1}^k \sigma_i u_i v_i^T \end{aligned}$$

complexity in storage for $A = mxn$
improved $= (m+n+1) k$

SVD

$\min_{\tilde{A}} \| A - \tilde{A} \|^2$ is given by

$A = U_{:,k} \Sigma_{l,k} V^T$ }
(first k columns of U & first k rows of V)

rank k have k non-zero singular values of A .

Ranking is $\sigma_1 \geq \sigma_2 \geq \dots$

columns of U spanning $C(A)$ with σ_1 as largest singular value of A with σ_2 as second largest singular value.

rank k grade of A

rank k maximum low rank subset of A

$$D + D^T = P$$

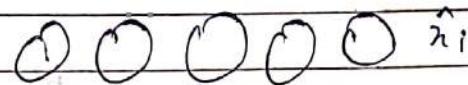
rank k maximum k rank subset of A

$$(D + D^T) \text{ diag} = P$$

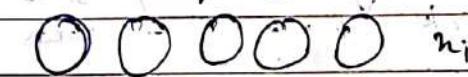
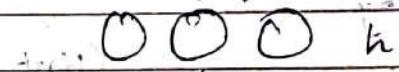
rank k largest k non-zero singular values of A

WEEK 7

* Auto Encoders



$$h = g(wn_i + b)$$



$$\hat{n}_i = f(w^*h + c)$$

Aim here is to reconstruct \hat{n}_i from h such that error

$$\|n_i - \hat{n}_i\|^2$$
 is minimized.

→ if $\dim(h) \geq \dim(n_i) \Rightarrow$ Over complete autoencoders
if $\dim(h) < \dim(n_i) \Rightarrow$ Under complete autoencoders
↓ as in above case

→ If the n_i nodes have real numbers then

$$\hat{n}_i = w^*h + c$$

If n_i nodes have binary numbers 0, 1 then

$$\hat{n}_i = \text{sigmoid}(w^*h + c)$$

f-in
at different
cases

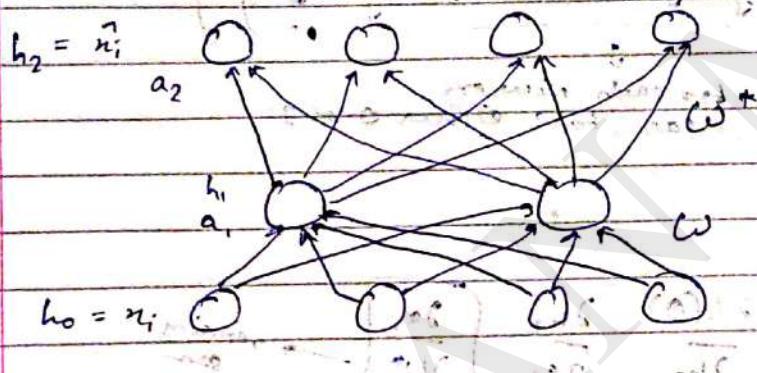
→ let's say loss function is sum of squared error

$$\text{ie } \min_{w, w^*, c, b} \sum_{m=1}^M \sum_{j=1}^n (\hat{x}_{ij} - x_{ij})^2$$

$$\Rightarrow \min_{w, w^*, c, b} \sum_{m=1}^M \sum_{i=1}^n (\hat{x}_i - x_i)^2$$

→ Backprop (for Real x_i)

$$L(\theta) = \frac{1}{2} \sum_{i=1}^n (\hat{x}_i - x_i)^2$$



Now earlier we did backprop with cross entropy loss

$$\frac{\partial L(\theta)}{\partial w^*} = \dots \rightarrow \frac{\partial L(\theta)}{\partial h_2} \left[\frac{\partial h_2}{\partial a_2} \frac{\partial a_2}{\partial w^*} \right] \xrightarrow{\text{remains same}} \downarrow \text{changes}$$

$$\frac{\partial L(\theta)}{\partial w} = \frac{\partial L(\theta)}{\partial h_2} \left[\frac{\partial h_2}{\partial a_2} \frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1} \frac{\partial a_1}{\partial w} \right] \xrightarrow{\text{remains same}}$$

Now

$$\frac{\partial L(\theta)}{\partial h_2} = \frac{\partial L(\theta)}{\partial \hat{x}_i} = \sum_j (\hat{x}_{ij} - x_{ij})^2 \xrightarrow{\text{scaled}}$$

$$= 2(\hat{x}_i - x_i) \xrightarrow{\text{refer machine learning book ①}}$$

→ Back Propagation (for binary a_{ij})

0 0 0 0

Cost = $\sum (0 - a_{ij})^2$

0 0 0 0

1 0 1 1 } we will use binary cross entropy loss

(in last row) (cross-entropy)

no. of entries : no. of neurons/inputs

$$\min \left(\sum_{i=1}^m \sum_{j=1}^4 \left[x_{ij} \log \hat{a}_{ij} + (1-x_{ij}) \log (1-\hat{a}_{ij}) \right] \right)$$

+ for each neuron
can take either 0 or 1

Now

$$\frac{\partial L(\theta)}{\partial w^*} = \frac{\partial L(\theta)}{\partial h_2} \frac{\partial h_2}{\partial a_2} \boxed{\frac{\partial a_2}{\partial w^*}} \rightarrow \text{remain same}$$

and previous steps & final operation like we did in back propagation

$$\frac{\partial L(\theta)}{\partial w^*} = \frac{\partial L(\theta)}{\partial h_2} \frac{\partial h_2}{\partial a_2} \boxed{\frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1} \frac{\partial a_1}{\partial w^*}} \rightarrow \text{remain same}$$

Now, loss for 1 data entry

$$\rightarrow \sum_{j=1}^4 [x_{ij} \log \hat{a}_{ij} + (1-x_{ij}) \log (1-\hat{a}_{ij})]$$

$$\frac{\partial L(\theta)}{\partial a_{ij}} = \frac{\partial L(\theta)}{\partial \hat{a}_{ij}} \cdot \frac{\partial \hat{a}_{ij}}{\partial a_{ij}} = -x_{ij} + 1 - x_{ij}$$

$$\frac{\partial L(\theta)}{\partial a_{ij}} = \frac{\partial \hat{a}_{ij}}{\partial a_{ij}} = \sigma(a_{ij})(1 - \sigma(a_{ij}))$$

* Link between PCA and Autoencoders

→ Conditions of noise in X for reconstruction held

- (1) Linear decoder $\rightarrow f$ is a linear function (given)
 - (2) Linear encoder $\rightarrow g$ is a linear function (prove) in order to minimize loss
 - (3) Loss function \rightarrow sum of squared errors (given) linear encoder is required
 - (4) Generalization \rightarrow $x_{ij} = \frac{1}{\sqrt{m}} \left(x_{ij} + \frac{1}{m} \sum_{k=1}^m x_{kj} \right)$
- Normalization zero mean

$$\rightarrow x_{ij} = \frac{1}{\sqrt{m}} (x_{ij} - \text{mean})$$

$$\therefore X = \frac{1}{\sqrt{m}} X'$$

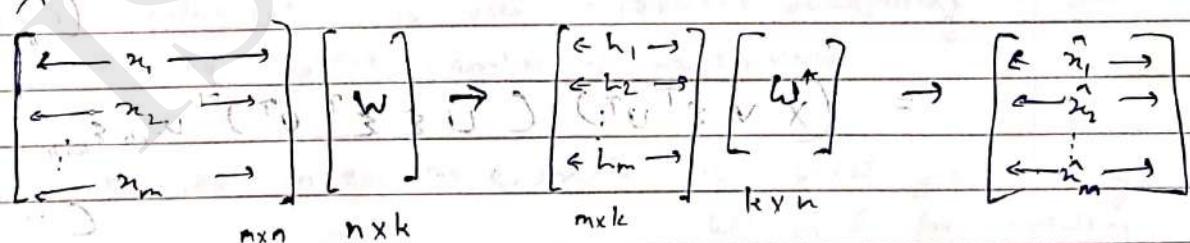
$$X^T X = \frac{1}{m} (X')^T (X')$$

prove it for m rows \rightarrow covariance matrix

X is full rank matrix \rightarrow $X^T X$ is invertible

→ Now loss function is

$$\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \hat{x}_{ij})^2$$



$$\min_{W^*} \|X - HW^*\|_F^2 = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \hat{x}_{ij})^2$$

Now from SVD theorem it is suggested that

Best approximation of X is given by $HW^* \approx \hat{X}$

$$\hat{X} = U \Sigma V^T$$

$$H = U_{:,k} \Sigma_{k,k} V_{:,k}^T$$

Now we can say,

$$H = U_{:,k} \Sigma_{k,k}$$

$$W^* = V_{:,k}^T$$

→ Now given all statement we try to prove
 H is linear encoding of X .

$$\Rightarrow H = U_{:,k} \Sigma_{k,k} = (x x^T) (x x^T)^{-1} U_{:,k} \Sigma_{k,k}$$

$$= (x v \Sigma^T u^T) (u \Sigma^T v^T u^T)^{-1} U_{:,k} \Sigma_{k,k}$$

$$(\because x = u \Sigma v^T)$$

$$= (x v \Sigma^T u^T) (v \Sigma^T v^T u^T)^{-1} U_{:,k} \Sigma_{k,k}$$

$$(\because v^T v = I)$$

$$= x v \Sigma^T u^T (v \Sigma^T)^{-1} u^T U_{:,k} \Sigma_{k,k} \quad (\because (ABC)^{-1} = C^{-1} B^{-1} A^{-1})$$

$$= x v \Sigma^T (v \Sigma^T)^{-1} u^T U_{:,k} \Sigma_{k,k}$$

$$= x v \Sigma^T \Sigma^{-1} u^T U_{:,k} \Sigma_{k,k} \quad (\because (AB)^{-1} = B^{-1} A^{-1})$$

$(\text{condition} = \mathbf{x} \mathbf{v} \mathbf{\Sigma}^{-1} \mathbf{I}_{k,k} \mathbf{\Sigma}_{k,k} \text{ and so on till } (\mathbf{v} \mathbf{U}^T \mathbf{U})_{k,k} = \mathbf{I}_{k,k})$

$$= \mathbf{x} \mathbf{v} \mathbf{I}_{k,k} \quad (\mathbf{\Sigma}^{-1} \mathbf{I}_{k,k} = \mathbf{\Sigma}_{k,k})$$

$$= \mathbf{x} \mathbf{v}_{\cdot, k}$$

linear combination of \mathbf{x} & $\mathbf{v}_{\cdot, k}$
↓
is weight matrix \mathbf{w}

* Regularization in Autoencoders

(Overfitting can happen in both undercomplete & overcomplete Autoencoders)

Two types

$$(1) \quad (\text{Loss Function} + \frac{\lambda \|\mathbf{w}\|^2}{2})$$

} updating the weights using L2 norm
flatten all \mathbf{w} matrix & put them in a vector \mathbf{w}

W.L.O.G. $\frac{\partial L}{\partial \mathbf{w}} = \text{nderivative of loss functions} + \lambda \mathbf{w}$

(2) Weights tying

Since regularization is used to prevent overfitting of data
we need to reduce number of parameters

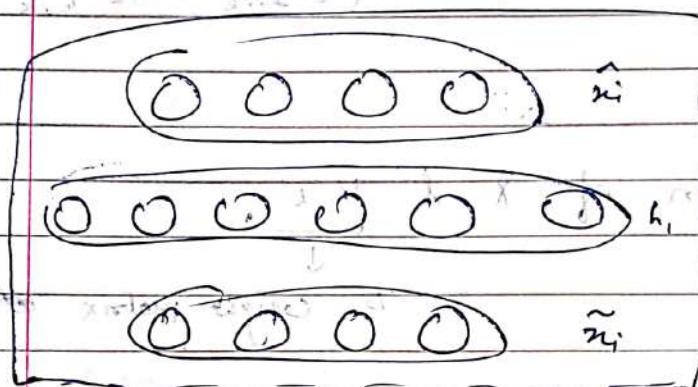
∴ As seen in earlier example \mathbf{w} is used for encoding
and \mathbf{w}^* is used for decoding

In this method $\mathbf{w}^* = \mathbf{w}^T$. Learning is as follows

At each instance only \mathbf{w} is updated.

The loss is for most frequent word which is used to find the word for filling at other positions.

* Denoising Autoencoders (Type of Regularization)



Autoencoder

\hat{z}_i is obtained by applying a function to z_i (denoising)

(original input z_i is corrupted and quantized)

Here, z_i is original input which is corrupted using a function & all ones are converted to zeros & vice versa in \hat{z}_i is formed from z_i + quantized noise (1)

Now,

\hat{z}_i is obtained at loss function. here will be

$$\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n (\hat{z}_{ij} - z_{ij})^2$$

but it is not (loss function) (2)

thus since it is overcomplete autoencoder

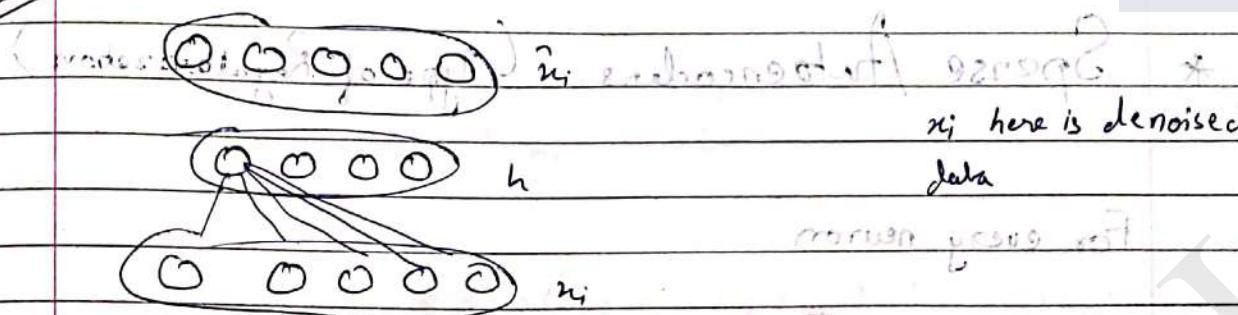
thus \hat{z}_i won't go in \hat{z}_i straight away because

\hat{z}_i is compared with z_i , so better fit in

thus it is given some loss function

By reconstructing over corrupted \hat{z}_i it also learns some important features of z_i and will be able to predict if some data is missing.

Experiment → shows that denoising the inputs helps learn better feature representations from input.



$$\text{Each neuron in } h_j = \sigma(w_j^T n_i)$$

Now each neuron fires when $(w_j^T n_i)$ is maximum.

Maximum occurs when it is zero because at this point sigmoid reaches zero and want to fire is -1. If this

is large enough then pass at each of them to get zero

$$\max \{ w_j^T n_i \} \text{ such that } \|x_i\|^2 = n_i^T n_i = 1$$

at this point n_i starts and when returning n_i is normalized

n_i is normalized

$w_j^T n_i$ dot product is maximum when $\theta = 0^\circ$ same direction

$$n_i = \frac{w_j}{\sqrt{w_j^T w_j}}$$

Each input here helps respective neuron to fire

$$w_1, \sqrt{w_1^T w_1}, \sqrt{w_2^T w_2}, \dots, \sqrt{w_n^T w_n}$$

n_1, n_2, \dots, n_n

variation of n_i on one input respects a $(0, 1) G$

This experiment is done after training set of all inputs

* Sparse Autoencoders (Type of Regularization)

For every neuron

$$\hat{p}_e = \frac{1}{m} \sum_{i=1}^m g(\omega^T, u_i) \quad \text{is calculated}$$

J

This gives us (at) an average number the neuron will fire and If \hat{p}_e is less than it is a sparse neuron & will be zero most of time. Since we have taken average it will be close to zero and hence sparse.

p is sparsity parameter which we choose close to zero lets say 0.005

Now we want $\hat{p}_e = p$ (at an average every neuron should be sparse)

∴ we use

$$\Omega(\theta) = \sum_{e=1}^k p \log \frac{\hat{p}_e}{p} + (1-p) \log \frac{1-p}{1-\hat{p}_e}$$

$$L(\theta) = L'(\theta) + \Omega(\theta)$$

will be minimum if $p = \hat{p}_e$

→ Back Propagation

$\frac{\partial L(\theta)}{\partial \omega}$ * charges thus we need to calculate

$$\frac{\partial \Omega(\theta)}{\partial \omega}$$

$$\frac{\partial \ell(\theta)}{\partial w} = + \frac{\partial \ell(\theta)}{\partial p} \times \frac{\partial p}{\partial w}$$

Now,

$$\frac{\partial \ell(\theta)}{\partial p} = \frac{\partial \log p}{\partial p} = \log p + (1-p) \log(1-p) -$$

$$-(1-p) \log(1-p)$$

$$\frac{\partial p}{\partial w}$$

$$= \left(\begin{array}{c} -p \\ - (1-p) \end{array} \right) \quad \begin{matrix} \text{vector with each} \\ \text{element representing} \\ \text{a neuron} \end{matrix}$$

(softmax is now converted to gradient form)

Now

$$\rightarrow \frac{\partial p}{\partial w_{xj}} = \frac{\partial}{\partial w_{xj}} \left(\frac{1}{m} \sum_{i=1}^m g(w_{xi}, n_{ij} + b_x) \right)$$

scalar

$\frac{\partial w_{xj}}{\partial w_{xj}}$

$$= \frac{1}{m} \sum_{i=1}^m \left[\underbrace{g'(w_{xi}, n_{ij} + b_x)}_{\text{scalar}} \times \underbrace{n_{ij}}_{\text{scalar}} \right]$$

$$\rightarrow \frac{\partial p}{\partial w_{xi}} = \frac{\partial}{\partial w_{xi}} \left[\frac{1}{m} \sum_{i=1}^m \left[g(w_{xi}, n_{ij} + b_x) \right] \right]$$

vector

$\frac{\partial w_{xi}}{\partial w_{xi}}$

$$= \frac{1}{m} \sum_{i=1}^m \left[\underbrace{g'(w_{xi}, n_{ij} + b_x)}_{\text{scalar}} \times \underbrace{n_{ij}}_{\text{vector}} \right]$$

$$\rightarrow \frac{\partial p}{\partial w} = \frac{\partial}{\partial w} \left(\frac{1}{m} \sum_{i=1}^m \left(g(w_{xi}, n_{ij} + b_x) \right) \right)$$

matrix

$\frac{\partial w}{\partial w}$

$$= \frac{1}{m} \sum_{i=1}^m \left[\underbrace{g'(w_i + b)}_{\text{vector}} \times \underbrace{n_i^T}_{\text{vector}} \right]$$

→ Final Derivative

$$\frac{\partial Q(\theta)}{\partial w} = \frac{1}{m} \sum_{i=1}^m \left(\underbrace{\begin{pmatrix} -p & (1-p) \\ \hat{p}_x & (1-\hat{p}_x) \end{pmatrix}}_{\text{vector from } \textcircled{1}} \odot \underbrace{g'(w_i + b)}_{\text{vector}} \right) \times \underbrace{n_i^T}_{\text{vector}}$$

* Contractive Autoencoder

we use Jacobian of h (vector) over x (vector)
refer the notes of machine learning book 2.

$$\therefore J = \begin{bmatrix} \frac{\partial h_1}{\partial x_1} & \frac{\partial h_1}{\partial x_2} & \dots & \frac{\partial h_1}{\partial x_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_k}{\partial x_1} & \frac{\partial h_k}{\partial x_2} & \dots & \frac{\partial h_k}{\partial x_m} \end{bmatrix}$$

x has dimension n (input layer vector or raw)

h has dimension k (hidden layer)

$$Q(\theta) = \sum_{j=1}^n \sum_{k=1}^l \left(\frac{\partial h_k}{\partial x_j} \right)^2$$

[using $\frac{\partial h_k}{\partial x_j} = \frac{\partial L}{\partial h_k} \cdot \frac{\partial h_k}{\partial z_j}$]

$$\min (L(\theta) + \alpha(\theta)) \rightarrow \text{aim}$$

$$\therefore \alpha(\theta) = 0$$

(residual sum of squares means don't capture variations in data)

$L(\theta) = \text{capture important variations in data}$

$L(\theta) = (\theta_i - f(\theta))^2 \quad (\text{changes with respect to } \theta_i)$

$\alpha(\theta) = \text{don't capture variations in data}$

(residual sum of squares doesn't change with respect to θ_i)

Trade off - capture only very important variations in data

$L(\theta) + \alpha(\theta) = L(\theta) + \alpha(\theta)$

$L(\theta) + \alpha(\theta) = L(\theta) + \alpha(\theta)$

Limitations of Principal Component Method

Residual sum of squares are just two

No short residual - big *

method { capture first few components
and ignore others }
can be optimally selected

Integers

Short

length

length

WEEK 8

* Bias... (Simple models have high bias)

→ Formula

$$\text{Bias}(\hat{f}(x)) = E[\hat{f}(x)] - f(x)$$

* Variance... (Complex Models have high variance)

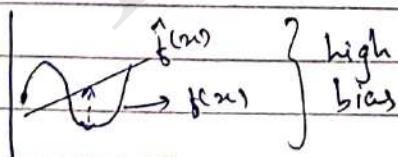
→ Formula

$$\text{Variance}(\hat{f}(x)) = E[(\hat{f}(x) - E[\hat{f}(x)])^2]$$

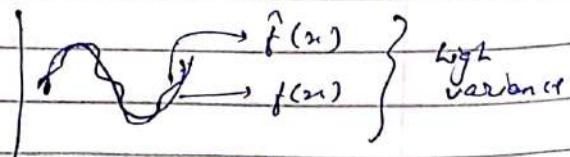
~~Precision~~ (~~High~~) ~~(Low)~~

Here unlike bias predicted are not compared to original but they are compared to one another

* Bias - Variance Trade off



Simple model

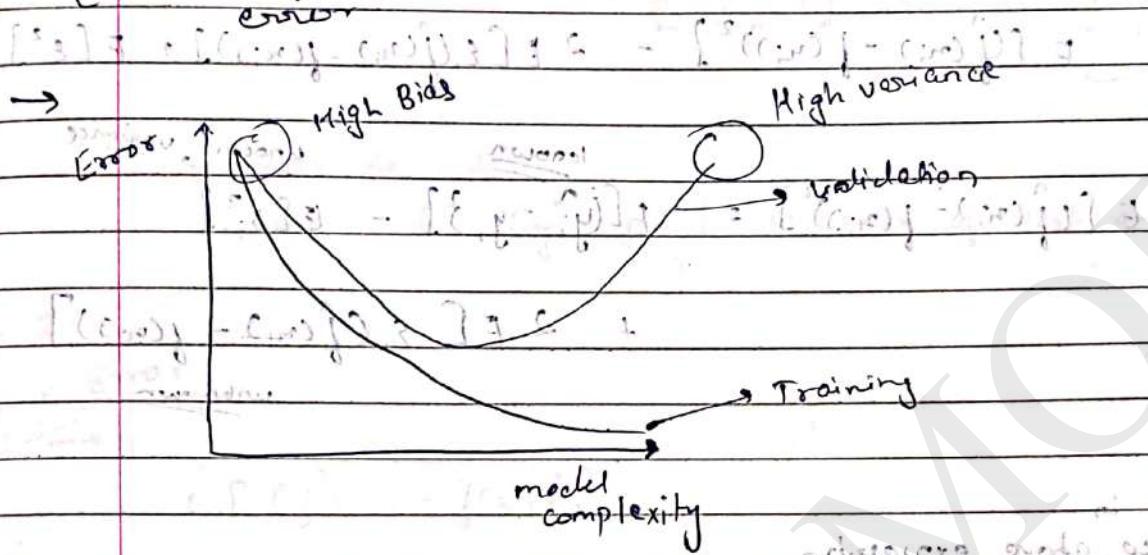


Complex model

It can be proved

$$E[(y - f(x))^2] = \text{Bias}^2 + \text{variance} + \sigma^2 \text{ (irreducible error)}$$

mean squared error



If there are outliers in training & in testing point :

$$\text{data for training} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

$$\text{data for test} = \frac{1}{m} \sum_{i=n+1}^{n+m} (y_i - \hat{f}(x_i))^2$$

→ Mathematically

we know training data (x_i, y_i) corresponding to it

∴ we can calculate \hat{y}_i

$$(\hat{y}_i - y_i)^2 = (\hat{f}(x_i) - y_i)^2$$

∴ $E[(\hat{y}_i - y_i)^2]$ can be empirically calculated

$$E[(\hat{y}_i - y_i)^2] = E[(\hat{f}(x_i) - f(x_i) + \epsilon_i)^2]$$

$$(Empirical) E[(\hat{y}_i - y_i)^2] = E[(\hat{f}(x_i) - f(x_i) + \epsilon_i)^2] = E[(\hat{f}(x_i) - f(x_i))^2 + 2\hat{f}(x_i)\epsilon_i + \epsilon_i^2] = E[(\hat{f}(x_i) - f(x_i))^2] + E[2\hat{f}(x_i)\epsilon_i] + E[\epsilon_i^2]$$

Here we assume that we know $f(x_i)$ but not $\hat{f}(x_i)$

$$\begin{aligned}
 & E[(y_i - \hat{y}_i)^2] = E[(f(x_i) - \hat{f}(x_i))^2] \\
 & = E[(\hat{f}(x_i) - f(x_i))^2 + 2\varepsilon_i(\hat{f}(x_i) - f(x_i)) + \varepsilon_i^2] \\
 & = E[(\hat{f}(x_i) - f(x_i))^2] + 2E[\varepsilon_i(\hat{f}(x_i) - f(x_i))] + E[\varepsilon_i^2] \\
 & \therefore E[(\hat{f}(x_i) - f(x_i))^2] = \underbrace{E[(\hat{y}_i - y_i)^2]}_{\text{known}} - \underbrace{E[\varepsilon_i^2]}_{\text{variance}} \\
 & \quad + 2E[\varepsilon_i(\hat{f}(x_i) - f(x_i))]
 \end{aligned}$$

Test Error

Now, above expression

- $E[(\hat{y}_i - y_i)^2]$ can be empirically written as

$$\frac{1}{m} \sum_{i=n+1}^{n+m} (\hat{y}_i - y_i)^2 \quad \left. \begin{array}{l} \text{for test data} \\ \text{this is not exact (E)} \\ \text{but it is value} \\ \text{very close to expectation} \end{array} \right\}$$

Similarly,

$$E[\varepsilon_i^2] \approx \frac{1}{m} \sum_{i=n+1}^{n+m} \varepsilon_i^2 \quad \left. \begin{array}{l} \text{for test data} \\ \text{(independent)} \end{array} \right\}$$

- $E[\varepsilon_i(\hat{f}(x_i) - f(x_i))]$ against covariance $(\varepsilon_i, (\hat{f}(x_i) - f(x_i)))$

$\hat{f}(x_i)$ for test data is calculated using parameters from training data which depend on ε_i of training data and not ε_i of test data

$$\therefore \text{cov}(x, y) = E[(x - \bar{x})(y - \bar{y})]$$

$$\begin{aligned}
 & \text{Variables are zero mean} \\
 & = E[xy]
 \end{aligned}$$

Now, since (ε_i) & $(\hat{f}(x_i) - f(x_i))$ are independent

$$(\varepsilon_i)(\hat{f}(x_i) - f(x_i)) = 0 \quad \text{if } \varepsilon_i \perp (\hat{f}(x_i) - f(x_i))$$

∴ cov = 0 (in simple words)

$$\therefore E[(\hat{f}(x_i) - f(x_i))^2] = \text{Train Error}$$

$$\text{Train Error} = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 + 3 \underbrace{\frac{1}{m} \sum_{i=1}^m \varepsilon_i^2}_{\text{small constant}}$$

\therefore Error depends

on $E[(\hat{y}_i - y_i)^2]$ in case of test error

$$= \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 + (\text{cov. } \varepsilon_i, \hat{f}(x_i)) \cdot \frac{3}{n} \sum_{i=1}^n \varepsilon_i^2$$

Train Error

\rightarrow Model is good or bad

$$E[(\hat{f}(x_i) - f(x_i))^2]$$

$$= \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 + - \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 + 2E[\varepsilon_i(\hat{f}(x_i) - f(x_i))]$$

Learning error depends on elements from error covariance ($\varepsilon_i, (\hat{f}(x_i) - f(x_i))$)
elements in cov. opn.

Here $E[\cdot]$ is independent on training dataset with ε_i

\therefore covariance $\neq 0$

Thus $E[(\hat{y}_i - y_i)^2]$ can't give correct estimate of error

\rightarrow Conclusion

Validation set gives better estimation of error than training set

→ Earlier

$E[\epsilon_i(f(x_i)) - f(x_i)]$ was unknown

Now,

using Stein's lemma

$$\frac{1}{n} \sum_{i=1}^n \epsilon_i(f(x_i)) - f(x_i) = \frac{\sigma^2 \sum_{i=1}^n \frac{\partial f(x_i)}{\partial y_i}}{n}$$

This value is high when $\frac{\partial f(x_i)}{\partial y_i}$

is high and it is high

$\frac{\partial f(x_i)}{\partial y_i}$ gives small change in observation causes large change in estimation

So this value is high in complex models

∴ True error = empirical train error

+ small constant term

+

α (model complexity)

→ This value minimize the fitted curve for model bias

$$\min_{\text{wrt } \theta} L_{\text{train}}(\theta) + \alpha R(\theta) = L(\theta)$$

concept of regularization

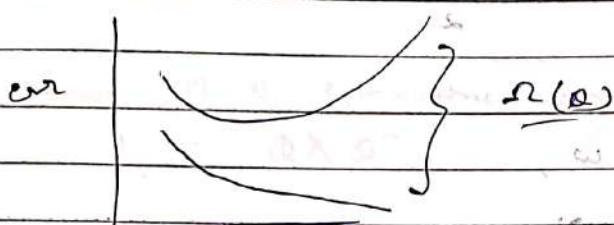
instead of train errors

Now we still don't know this quantity $\omega(\theta)$
thus we use

$$(H^T(\omega)) \propto (\omega) \text{ or } H^T(\omega) = (\omega) \Rightarrow H^T(\omega) = (\omega)$$

any function as $\omega(\theta)$ such that it is high for complex models & low for simple models.

$$(\omega - \omega_0) H^T(\omega - \omega_0) \geq 0$$



* L2 Regularization

$$L(\tilde{\omega}) = L(\omega) + \frac{1}{2} \alpha \|\omega\|^2$$

$$\nabla L(\omega) = \nabla L(\omega) + \alpha \omega \quad (\omega) \propto \omega$$

update

$$\omega_{t+1} = \omega_t - \eta (\nabla L(\omega_t) + \alpha \omega_t)$$

Now $\omega = (\omega)^T$ loss function of $\tilde{\omega}$ is

$$\text{let } \omega^* \text{ be such that } \nabla L(\omega^*) = 0$$

Loss
without
regularization

& $h = \omega - \omega^*$ where (ω) is a point
in neighbourhood of ω^*

Now quest is to find $\nabla L(\omega)$

from Taylor series in upto 2nd terms note we will 2nd derivative

$$L(\omega^* + h) = L(\omega^*) + (h)^T \nabla L(\omega^*) + \frac{1}{2} (h)^T H(h)$$

$$\begin{aligned} L(\omega) &= L(\omega^*) + \underbrace{(h - \omega^*)^T \nabla L(\omega^*)}_{0} - (\because \nabla L(\omega^*) = 0 \text{ from } ①) \\ &\quad + \frac{1}{2} (\omega - \omega^*)^T H (\omega - \omega^*) \end{aligned}$$

Differentiating wrt ω ,

$$\begin{aligned} \nabla L(\omega) &= \frac{\nabla L(\omega^*)}{0} + \cancel{\frac{1}{2} (\omega - \omega^*)^T H (\omega - \omega^*)} - (\because \nabla L(\omega^*) = 0 \text{ from } ①) \\ &\quad + (\omega - \omega^*)^T H \end{aligned}$$

$$\nabla L(\omega) = (\omega - \omega^*)^T H$$

Now,

$$\tilde{L}(\tilde{\omega}) = \nabla L(\omega) + \alpha \omega + (\omega)^T V \quad \text{regularized loss function } \tilde{L}(\omega)$$

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



Let $\tilde{\omega}$ be such that $\tilde{L}(\tilde{\omega}) = 0$ with \leftarrow

regularized loss

$$\therefore \text{from } H(\tilde{\omega}, \omega^*) + \alpha \tilde{\omega} = 0 \Leftrightarrow \tilde{\omega} = -H^{-1} \alpha$$

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

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

Note: $\alpha \rightarrow 0$ then $\tilde{\omega} = \omega^*$ (no regularization)

→ Assume H is symmetric matrix

$$H = Q \lambda Q^T \quad (Q Q^T = I)$$

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

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

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

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

$$= Q^T (Q \lambda Q^T)^{-1} Q \lambda Q^T \omega^*$$

$$= Q \boxed{(Q \lambda Q^T)^{-1} Q \lambda Q^T} \omega^*$$

Diagonal matrix D

$$\tilde{\omega} = Q D Q^T \text{ weight } \xrightarrow{\text{rotates } \omega^*}$$

rotates ω^*
 scales ω^*
 rotates ω^*

Here $\tilde{\omega}$ & ω^* are weight vectors

$$D = \begin{bmatrix} \frac{\lambda_1}{\lambda_1 + \alpha} & \dots & \frac{\lambda_n}{\lambda_n + \alpha} \\ \vdots & \ddots & \vdots \\ \frac{\lambda_1}{\lambda_1 + \alpha} & \dots & \frac{\lambda_n}{\lambda_n + \alpha} \end{bmatrix}$$

(Centered by α) $\omega = \frac{1}{\alpha} \begin{bmatrix} \lambda_1 & \dots & \lambda_n \end{bmatrix}^T - \frac{1}{\alpha} \mathbf{I}$

$$= (\mathbf{I} + \alpha \mathbf{D})^{-1} \quad \text{where } \mathbf{D} \text{ is diagonal with entries } \lambda_i$$

$(\mathbf{I} = \text{Identity matrix}) \quad \text{rank } \mathbf{D} = n \rightarrow$

$$\begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & & \lambda_1 \\ & \frac{1}{\lambda_2 + \alpha} & \dots & & \lambda_2 \\ & & \ddots & & \vdots \\ & & & \frac{1}{\lambda_n + \alpha} & \lambda_n \end{bmatrix}$$

$+ \alpha^{-1} \text{rank } \mathbf{D}^{-1} (\text{rank } \mathbf{D} = n) = \frac{1}{\alpha} \text{rank } \mathbf{D}$

$+ \alpha^{-1} \text{rank } \mathbf{D}^{-1} (\text{rank } \mathbf{D} = n) = \alpha^{-1} n$

Note

$$\text{if } \lambda_i \gg \alpha, \frac{\lambda_i}{\lambda_i + \alpha} \approx 1$$

$$\text{if } \lambda_i \ll \alpha, \frac{\lambda_i}{\lambda_i + \alpha} \approx 0$$

A similar thought

* Dataset Augmentation

Increasing the amount of data

- Rotating images
- Shifting horizontally / vertically
- Blurring, adding noise etc

* Parameter Sharing & Tying is used in CNN and in Encoder/Decoder

$\theta = \theta_{\text{enc}} \oplus (\theta - \theta_{\text{enc}})$ ensures that

* Adding Noise to Data (Inputs) / (Encoding)

$\epsilon \sim N(0, \sigma^2)$ Gaussian distribution
(Zero mean, unit variance)

$\tilde{x}_i = x_i + \epsilon_i$ or (zero mean, multivariate)

$$\hat{y} = \sum_{i=1}^n w_i x_i$$

$$= \hat{y} + \sum_{i=1}^n w_i \epsilon_i$$

$$\begin{aligned} \hat{y} &= \sum_{i=1}^n w_i \tilde{x}_i = \sum_{i=1}^n w_i x_i + \sum_{i=1}^n w_i \epsilon_i \\ &= \hat{y} + \sum_{i=1}^n w_i \epsilon_i \end{aligned}$$

Now, we are interested in $E[(\hat{y} - y)^2]$

$$\therefore E[(\hat{y} - y)^2] = E\left[\left(\hat{y} + \sum_{i=1}^n w_i \epsilon_i - y\right)^2\right]$$

$$= E\left[\left((\hat{y} - y) + \left(\sum_{i=1}^n w_i \epsilon_i\right)\right)^2\right]$$

$$= E[(\hat{y} - y)^2] + E\left[2(\hat{y} - y) \sum_{i=1}^n w_i \epsilon_i\right]$$

$$= E[(\hat{y} - y)^2] + E\left[\left(\sum_{i=1}^n w_i \epsilon_i\right)^2\right]$$

Here in second term ~~we add noise to~~ ~~original values~~

$(\hat{y} - y) \& \varepsilon_i$ are independent
original values when noise was not added

thus covariance $(\hat{y} - y) \varepsilon_i = 0$

Covariance (cyclic) total of small probabilities

A in third term

$$(\text{covariance})^2 \Rightarrow (\hat{y} - y)^2$$

(since there is no bias)

$$(\varepsilon_i \varepsilon_j)^2 \approx 0 \quad \text{independent}$$

$$(\varepsilon_i^2) \neq 0$$

$$\therefore (\sum w_i \varepsilon_i)^2 = \sum w_i^2 \varepsilon_i^2$$

$$= E[(\hat{y} - y)^2] + 0 + E\left[\sum_{i=1}^n w_i^2 \varepsilon_i^2\right]$$

$$\text{Here } E\left[\sum_{i=1}^n w_i^2 \varepsilon_i^2\right] = \left(\sum_{i=1}^n w_i^2\right) E[\varepsilon_i^2]$$

since w_i is not random variable

$$= E[(\hat{y} - y)^2] + \sigma^2 \sum_{i=1}^n w_i^2$$

↓ same as L₂

regularization

This is used in overcomplete autoencoders for encoding, in order to reduce the complexity by number of dimensions increase.

Input

* Adding Noise to Data (Output) (softmax & softplus)

At the output we have

$$\sum_{i=1}^n p_i \log q_i \quad (\omega - \omega)^T H = (\omega)^T \nabla$$

$\downarrow \quad \downarrow$
original predicted

we add noise to original so that training error doesn't drive to 0

$$p_i = \{0, 0, 1, 0, \dots\} \quad \text{Assuming there are } q \text{ elements}$$
$$p_i = \left\{ \frac{1+\epsilon}{q}, \frac{\epsilon}{q}, \frac{1-\epsilon}{q}, \frac{\epsilon}{q}, \dots \right\} \quad \begin{matrix} \text{if } \epsilon = \text{noise} \\ \text{the correction} \\ \text{is present} \end{matrix}$$

$$q_i = \text{remains as it is } (1 + \epsilon)$$

* Early Stopping

Patience parameter p

→ At every epochs check train error & validation error and if for previous p epoch if validation error has increased or remained same stop

→ In gradient Descent ω is propagated plus ϵ and

$$\omega_{t+1} = \omega_t + \eta \nabla \omega_t (t+1) - \epsilon] \quad \epsilon = \epsilon_1$$

$$= \omega_0 + \eta \sum_{i=1}^t \nabla \omega_i \quad \epsilon = \epsilon_1$$

maximum of them

$$\omega_{t+1} \geq \omega_0 + \eta t \epsilon$$

Thus t controls how far ω_t is from ω^* .

→ As seen earlier

$$\nabla L(\omega) = H(\omega - \omega^*)$$

$$\omega_t = \omega_{t-1} + \eta \nabla L(\omega_{t-1})$$

$$= \omega_{t-1} + \eta H(\omega_{t-1} - \omega^*)$$

$$\approx (I + \eta H)\omega_{t-1} - \eta H\omega^*$$

$$\therefore \omega_t = (I + \eta H)\omega_{t-1} - \eta H\omega^*$$

- Using EVD, $H = Q\Lambda Q^T$

$$\omega_t = (I + \eta Q\Lambda Q^T)\omega_{t-1} - \eta Q\Lambda Q^T\omega^*$$

$$\omega_t = Q [Q^T(I - (I - \epsilon\Lambda)^{-1})Q^T\omega^*] \quad \text{①}$$

- This \Rightarrow early stopping is similar to L_2 regularization.

$$\tilde{\omega} = Q [I - (\lambda + \alpha I)^{-1}\alpha] Q^T \omega^* \quad \text{②}$$

- $\omega_t = \tilde{\omega}$ for following condition

$$(I - \epsilon\Lambda)^{-1} = (\lambda + \alpha I)^{-1}\alpha$$

* Ensemble Methods

→ Methods like Bagging (ML book 2) dataset is broken down into parts and models are prepared for each training data and error is evaluated for each part.

Evaluate

Avg of all errors $\bar{\epsilon}$ is given (Here the number of models is k)

$$\frac{1}{k} \sum_i \epsilon_i$$

Expectation over the inputs

$$\text{Expected mean squared error} = E \left[\left(\frac{1}{k} \sum_{i=1}^k \epsilon_i \right)^2 \right]$$

equal to variance of average of errors

$$\text{Variance} = \text{Var} \left[\frac{1}{k} \sum_{i=1}^k \epsilon_i \right] = \frac{1}{k^2} \sum_{i=1}^k \sum_{j=1, j \neq i}^k E(\epsilon_i \epsilon_j)$$

$$= \frac{1}{k^2} \left[\sum_{i=1}^k E(\epsilon_i^2) + k(k-1) \sum_{i \neq j} E(\epsilon_i \epsilon_j) \right]$$

$$= \frac{1}{k^2} \left[\sum_i E(\epsilon_i^2) + \sum_{i \neq j} E(\epsilon_i \epsilon_j) \right]$$

$$= \frac{1}{k^2} \left[kV + (k-1)C \right]$$

$$= V + \left(\frac{k-1}{k} \right) C$$

independence

independence

covariance

Note

If errors are independent $C = 0$

method is useful (s and m) passed still gradient descent because V becomes zero after many steps this will happen in case when back propagation

If errors are correlated $C = V$

$$\text{mse} = V$$

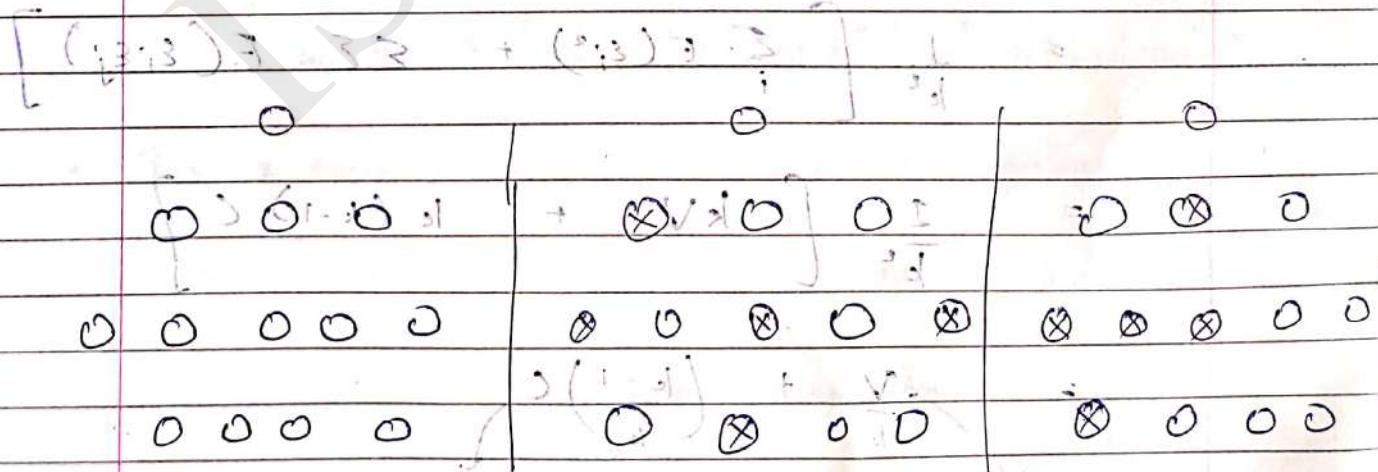
standard

Thus using averaging helps reduce errors in forward pass

* Dropout

Let W be the weight matrix which needs to be updated

- Dropout means to eliminate neurons in layers with probability p form completely new model train for one epoch on that model. update W .
- Now eliminate from original to get new model I use updated W and train for 1 epoch



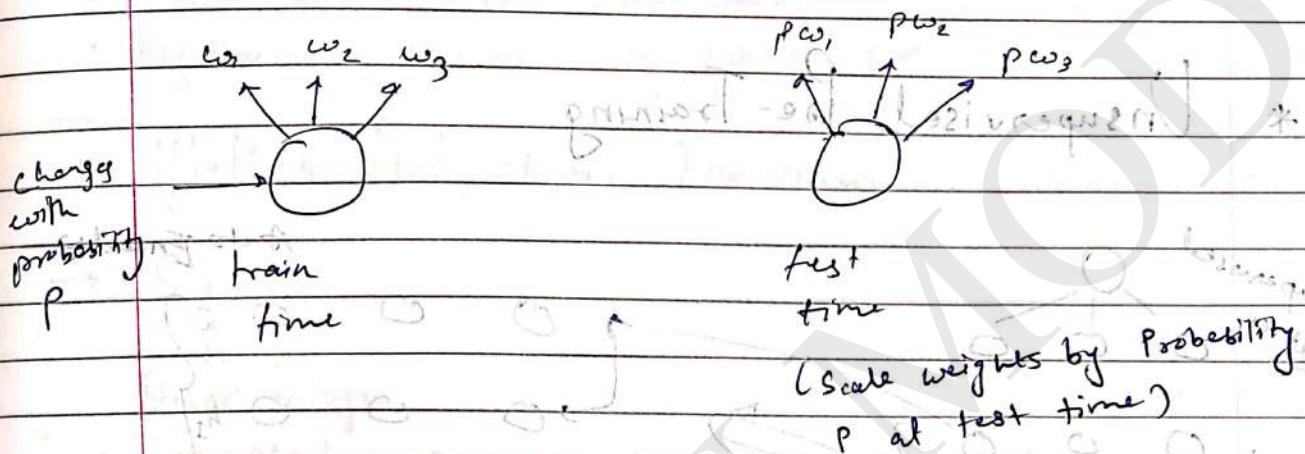
Original

Dropout 1

Dropout 2

This will make all neurons independent from one another because they can disappear any time.

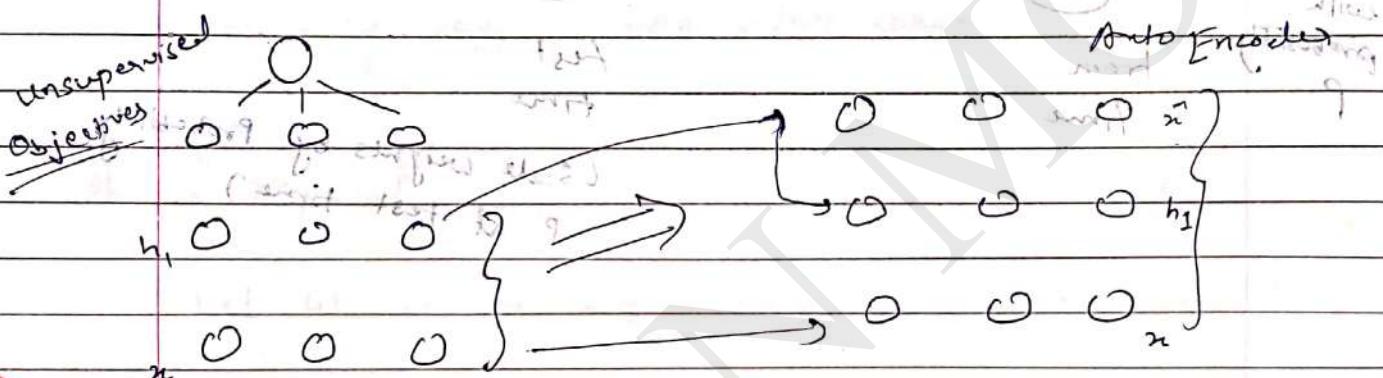
These neurons will be robust (due to fading).



WEEK 9

→ Gradient at a particular layer depends on input at that layer (refer backprop deep learning ①)

* Unsupervised Pre-Training



Thus, between each layer there is an encoder which gives abstraction of previous layer

Thus we are trying to minimize

$$\text{Loss} = \sum_m \sum_{i=1}^n \sum_{j=1}^n (n_{ij} - \hat{n}_{ij})^2$$

After this is done for entire networks the weights will be initialized such that each layer is abstraction of previous.

Instead of random weights we used weight from unsupervised pre-training.

~~supervised
objectives~~

$$L(\theta) = \frac{1}{m} \sum_{i=1}^m (y_i - f(x_i))^2$$

Normal
Backprop

- Helps converge Faster
- Optimization Problem (Train Data)
- Regularization Problem (Test Data)

* Better Activation Functions

→ Sigmoid

Disadvantages

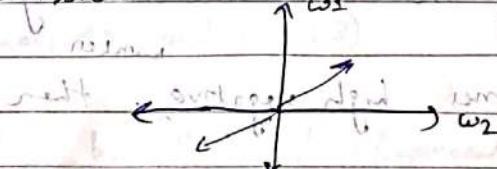
- Saturated neurons cause gradient to vanish

Here, in this case if value of x is very high or very low then sigmoid(x) saturates that is, its value becomes 1 and 0 respectively.

thus derivative $\sigma(x) \times (1 - \sigma(x))$ which is used in gradient becomes zero. Thus gradient become zero.

- Not zero centered (from 0 to 1)

the movement of gradient becomes restricted and thus convergence is slow due to either all the gradients at a layer either both are +ve or both are -ve



- Computationally expensive ($\exp(x)$)

→ Tanh.

Disadvantage

- Gradient vanishes at saturation after optimal point.
- Computationally expensive due to multiple additions.

Solve

- zero-centered.

→ ReLU

$$f(x) = \max(0, x)$$

$$f(x) = \max(0, x+1) - \max(0, x-1) \quad \left\{ \begin{array}{l} \text{equivalent} \\ \text{to sigmoid} \end{array} \right.$$

or first pass is the ReLU of the ReLU of input

2nd pass is the ReLU of the ReLU of the ReLU of input

Mostly used in neurons having sigmoidal activation

Only problem (dead neuron problem)

$$\frac{\partial f(x)}{\partial x} = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases}$$

derivatives

In practice a large fraction of ReLU units can die

if the learning rate is set too high

number

If bias becomes high negative, then

$$w_1x_1 + w_2x_2 + b < 0$$

(negative) vanishing phenomenon

$\therefore \frac{\partial f(x)}{\partial x} = 0$ if weights are not updated

→ Leaky ReLU (gradient, non-saturation) $f(x) = \max(0, \alpha x, x)$ (solution to dead neurons)

→ Parametric ReLU

$$f(x) = \max(\alpha x, x)$$

lower sticking parameter not higher parameters don't work
tough neurons

→ Exponential Linear Unit

more id. non-saturation and higher values enough work

$$f(x) = x \text{ if } x > 0 \\ (\text{middle position between 0 and infinity}) \\ = \alpha e^x - 1 \text{ if } x \leq 0$$

→ Maxout Network (many way of doing non-saturation)

$$\max(C + w_1^T x + b_1, \dots, w_m^T x + b_m) \quad \text{--- (1)}$$

variant of ReLU $\max(0, \alpha x)$ at --- (2) parallelism

$$\left. \begin{array}{l} w_1 = 0, w_2 = 0 \text{ and} \\ w_2 = 1, b_2 = 0 \end{array} \right\} \text{Representing (1) as (2)}$$

$$\max(\alpha x, x) \quad \text{--- (3)}$$

$$\left. \begin{array}{l} w_1 = \alpha, b_1 = 0 \\ w_2 = 1, b_2 = 0 \end{array} \right\} \text{Representing (1) as (3)}$$

* Better Weight Initialization Strategies

→ Weights are initialized by zero.

$$a_{11} = w_{11}n_1 + w_{12}n_2$$

$$a_{12} = w_{21}n_1 + w_{22}n_2$$

$$a_{11} = a_{12} \quad \& \quad h_{11} = h_{12}$$

Thus both weights will get the same update and remain equal.

→ Same happens when weights are initialized by same values.

(Known as Symmetric Breaking Problem)

→ Initializing weights to very small values.

① Vanishing gradient problem occurs if more layers are present.

→ Initializing weights to large values.

② Saturation occurs and thus the vanishing gradient problem persists.
variance is low when plotted

Principled Way

$$\rightarrow S_{11} = \sum_{i=1}^n w_{11}n_i$$

$$V(XY) = \underbrace{(\bar{x})^2 V(\gamma) + (\bar{\gamma})^2 V(x) + V(\gamma) V(x)}_{\text{Formula}} + \text{DELUXE}$$

PAGE NO.:
DATE:



$$V(S_{1,i}) = V(\sum_{i=1}^n w_i x_i) = \sum_{i=1}^n V(w_i x_i)$$

Now $(\sum w_i x_i)$ has weight w_i

$$\text{Expected variance} = \sum_{i=1}^n E[(E[w_i x_i])^2] - V(x_i)$$

$$= \text{const} + \text{const} (\mathbb{E}[x_i])^2 + \text{const} V(x_i)$$

$$+ V(x_i) V(w_i x_i)$$

$$(\sum) \text{now} = (\omega) \text{now}$$

lets say x_i & w_i

both wt are 0 mean 1

standard deviation

$$V(x_i) = V(n), V(w_i) = V(\omega)$$

$$= \sum_{i=1}^n V(x_i) V(w_i)$$

$$(S)_{\text{now}} = (n)_{\text{now}}$$

$$= n V(\omega) V(n)$$

$$\therefore V(S_{1,i}) = (n V(\omega)) (V(n)) \quad \text{--- (1)}$$

Now,

$$V(S_{2,i}) = \sum_{i=1}^n V(S_{1,i}) V(w_{2,i})$$

$$= n V(S_{1,i}) V(w_2)$$

Now we have to substitute (1) in it
(not to find now)

$$V(S_{2,i}) \propto \left[\underbrace{[n V(\omega)]^2 V(n)}_{(1)} \right]$$

thus low values & high values

would lead to a variance of $S_{2,i}$
or very high variance of $S_{2,i}$



$$\text{Var}(s_{ki}) = [n \text{Var}(w)]^k \text{Var}(x_i) = (1.2)^{\text{avg}}$$

if

$n \text{Var}(w) = 1$ then $\text{var}(s_{ki})$ will remain normal

→ If $(z)_{avg} = 0$ and $(z)_{std} \sim N(0, 1)$ and we take weights such that

$$(z, w)_{\text{avg}} = 0$$

$$n \text{Var}(w) = n \text{Var}(z)$$

thus

\sqrt{n} is number of neurons in the layer
normalizing variable

$$\text{since } (w = \frac{z}{\sqrt{n}})$$

$$\text{Now } \left. \begin{array}{l} \text{var}(z) \\ \text{var}(w) = \frac{1}{n} \end{array} \right\} \text{var}(az) = a^2 \text{var}(z)$$

$$= n \times \frac{1}{n} \text{var}(z)$$

$$= (z)_{\text{avg}}^2 (w)_{\text{avg}} = (1.2)^{\text{avg}}$$

$$= \text{var}(z) = 1 \text{ since normal distribution}$$

Applicable for Sigmoid & Tanh (1.2)^{\text{avg}}

$$(1.2)^{\text{avg}} (1.2)^{\text{avg}}$$

→ For ReLU (since it is 0 for negative values that is approx half of time)

$$w = \frac{[z]_{\text{avg}}^2 [(w)_{\text{avg}}]}{\sqrt{\frac{n}{2}}} \times (1.2)^{\text{avg}}$$

positive avg & smaller weight

1.2 for minimum of best choice
1.2 for maximum of best choice



input vector = $\{x_1, x_2, x_3, \dots, x_n\}$

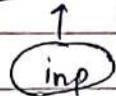
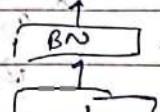
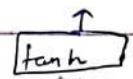
here we normalize
with respect to
features of
input

DELUXE
PAGE NO.:
DATE:



so if there are 100 samples
 $x_i = 100 \times 1$ vector similar
to others & we normalize x_i

* Batch Normalization



in between layers normalize values so that input to the next layer comes from the same distribution as the entire network

$$\rightarrow \hat{s}_{ik} = \frac{s_{ik} - E[s_{ik}]}{\sqrt{\text{var}(s_{ik})}} \quad 0 \text{ mean, unit variance}$$

for backprop. we need to differentiate the above expression also

\rightarrow If we don't want normalized values we can convert batch Norm to original

$$y^{(k)} = \gamma^{(k)} \hat{s}_{ik} + \beta^{(k)}$$

where, network learns γ & β which should be as follows

$$\gamma^{(k)} = \sqrt{\text{var}(s_{ik})} \quad \text{and} \quad \beta^{(k)} = E[s_{ik}]$$

\rightarrow Batch Norms helps network converge faster in some cases

For better understanding refer this [link](#) - [this](#)

deeplearning1 notes / week9 / batchnorm.tutorial.html

/ batchnorm.png



WEEK 10

* One-Hot Representation of Words (Sparse Representation)

Corpus → Collection of Sentences

Vocabulary - All unique words in corpus

→ Let's say vocabulary is

cat, dog, truck

One Hot = $[1, 0, 0]$, $[0, 1, 0]$, $[0, 0, 1]$

- It Occupies lot of space

- No similarity between words (e.g. Euclidean distance between cat and dog should be less because both are animals)

* Distributed Representation of Words

→ Write the individual unique words in corpus as rows and columns

→ Make a co-occurrence matrix

words Human, machine, interface, for, computer, applications

Corpus - Cat eats a dog

Dog eats a bone

Cat eats a mouse

vocab - Cat, eats, a, Dog, bone, mouse

(Count in (1,1) hours : 200)

Table 7 - Co-occurrence matrix $\Rightarrow k=1$ TM9 & (2,1).

context		Cat	eats	Dog	Bone	Mouse	TM9
words			(2,1)	Dog	Mouse	TM9	-
Cat		0	2	0	0	0	0
eats	1	2	0	1	0	0	0
Dog				0			
Bone				0			
Mouse					0		

• Here k is not for a current neighbour word.
we need to check 1 neighbour towards left and right.

→ (Cat) appeared two times around the word eat
in a window of 1 word around it.

→ This is still very complex

→ Solution of QV2 problem

- Remove stopwords from column (context)
- Ignore very frequently occurring words

• Use threshold, $t = 100$

entry in above matrix

$$x_{ij} = \min(\text{count}(w_i, c_j), t)$$

word context

glossary entries

for words less than

threshold words less than

- Use $\text{PMI } l = \frac{c}{w}$ (high if $\text{count}(w, c)$ is low) $\rightarrow p(c|w)$

$$\begin{aligned} - \text{PMI}(w, c) &= \log \frac{p(c|w)}{p(c)} = \log \left(\frac{\text{count}(w, c)}{\text{count}(w)} / \frac{\text{count}(c)}{N} \right) \\ &= \log \left(\frac{\text{count}(w, c) * N}{\text{count}(c) * \text{count}(w)} \right) \end{aligned}$$

N is total no. of words

$$- \text{count}(w, c) = 0, \text{ PMI} = -\infty$$

(+) $\text{PMI}_0(w, c) = \text{PMI}(w, c)$, if $\text{count}(w, c) > 0$

else $\text{PMI}_0(w, c) = 0$, i.e., otherwise

(2) $\text{PPMI}(w, c) = \text{PMI}(w, c)$ if $\text{PMI}(w, c) > 0$
positive if $\text{count}(w, c) > 0$, otherwise 0

* Using SVD for dimensionality reduction

$$\hat{X}_{m \times n} = \underbrace{(U_{m \times k} \Sigma_{k \times k} V_{k \times n})}_{\text{approximate}} \text{discrete matrix}.$$

$$\hat{X} = u \sigma v^T$$

rank k approx
we construct it using k matrix

Now $(\hat{u}, (\hat{\sigma}, \hat{v}))$ rows = \hat{x}

$$\hat{X}\hat{X}^T = \text{dot product} = \text{cosine similarity}$$

if we ignore denominator

(Thus each element of $\hat{X}\hat{X}^T$ is cosine similarity between $(i^{\text{th}} \text{ row of } X) \text{ and } (j^{\text{th}} \text{ column of } X^T)$)

Now let initial PPMI table be X & transformed table be $\hat{X} \hat{X}^T$

$$\therefore X \leq \hat{X} \hat{X}^T \quad \{ \text{similarity between various}$$

+ terms, so smaller off diagonal rows & cols increases by it is able to learn relation

(Now, what is the role of Σ in this?)

$$X = U\Sigma V^T$$

$$XX^T = (U\Sigma V^T)(U\Sigma V^T)^T$$

$$(1) \text{ trace comparison } U\Sigma V^T \cdot V\Sigma^T U^T$$

$$= U\Sigma \Sigma^T U^T$$

$$(2) XX^T = \text{diag}(U\Sigma \Sigma^T U^T)$$

$$m \times n \xrightarrow{\text{dimensions}} m \times k \text{ & } k \times m$$

$$\text{if let } U\Sigma V^T = W_{\text{word}} \quad \left\{ \begin{array}{l} \text{dimension of } W_{\text{word}} = m \times k \\ \text{expressed in terms of basis } U \end{array} \right.$$

$$XX^T = W_{\text{word}} W_{\text{word}}^T$$

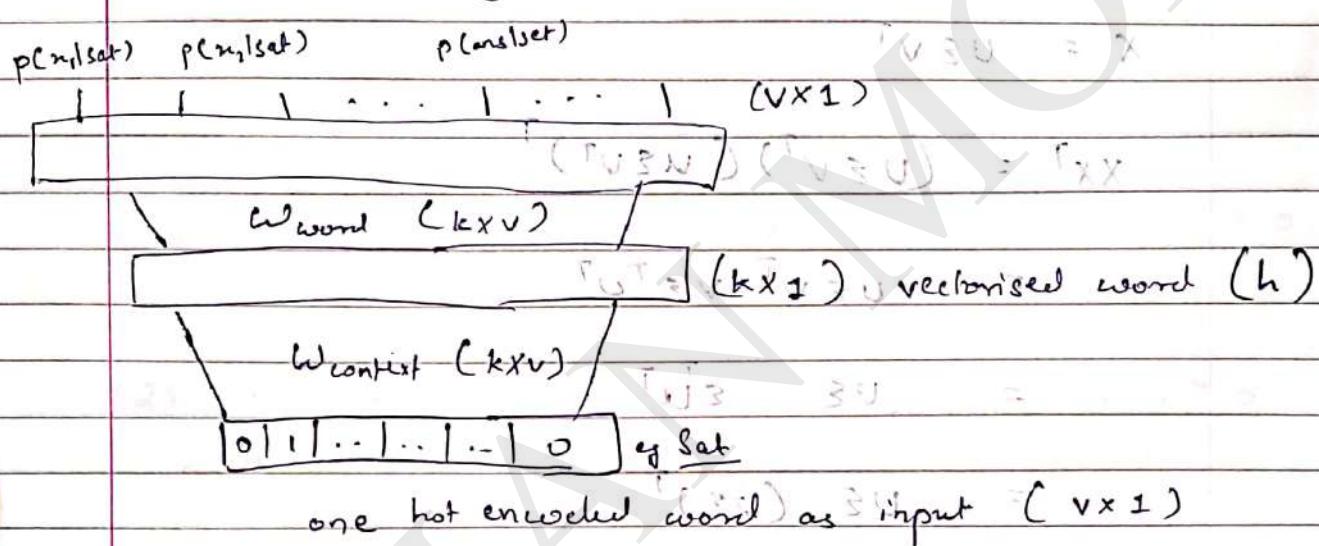
both have same cosine similarity (dot product) and ignore denominator

though dimension of W_{word} is very less than X

$$W_{\text{context}} = V_{n \times k \times p} \rightarrow y$$

part 2: if it is the reduction for columns or context words

* Continuous Bag of Words (given $n-1$ words predict n^{th} word)



So, here we have made a Neural Network.

→ w_{word} & w_{context} are weight matrix to be learnt between layers

→ Final activation is Softmax since we need probability

→ loss is Cross-Entropy loss.

Working

→ When we move from input x_i to x_{i+1}

(Word context is $x = h$)

g-

w_{context}

$$x = h$$

$$\begin{bmatrix} -1 & \{0.5\} & 2 \\ 3 & -1 & -2 \\ -2 & 1.7 & 3 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.5 \\ -1 \\ 1.7 \end{bmatrix} = (a)$$

∴ if i^{th} index of x is $\neq 1$ then

More $h = i^{\text{th}}$ column of w_{context} .

→ Now,

$$\text{softmax}(w_{\text{word}} \cdot h) = \text{output. } (\sim 1 \times 1)$$

Let's see any one entry from 1×1

$$p(\text{con} | \text{sat}) = \frac{e^{(w_{\text{word}} \cdot h)(i)}}{\sum_j e^{(w_{\text{word}} \cdot h)(j)}}$$

Now,

$$w_{\text{word}}^T \cdot h = \text{o/p}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 2 \\ 1 & 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 4 \\ 6 \\ 6 \end{bmatrix} \rightarrow 4$$

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} = 4 \quad | \quad j^{\text{th}} \text{ column of } w_{\text{word}} : h$$

Now, $h = i^{\text{th}}$ column of $\mathbf{W}_{\text{context}}$ is marked as

$(j^{\text{th}} \text{ col of } \mathbf{W}_{\text{word}}) \times (i^{\text{th}} \text{ col of } \mathbf{W}_{\text{context}})$



$$L(\theta) = -\log \hat{y}_w = -\log (w(c))$$

$$h = W_{\text{context}} \cdot n_c = u_c$$

now I is the i^{th} column of $\mathbf{W}_{\text{context}}$

$$\hat{y}_w = \exp(u_c \cdot v_w)$$

$$\sum_{w \in V} \exp(u_c \cdot v_w)$$

Now we want to learn v_w

$$L(\theta) = -\log \hat{y}_w$$

$$= -\log \left(\frac{\exp(u_c \cdot v_w)}{\sum_{w' \in V} \exp(u_c \cdot v_{w'})} \right)$$

$$= -\log(u_c \cdot v_w) + \log \sum_{w' \in V} \exp(u_c \cdot v_{w'})$$

$$\nabla_{v_w} = \frac{\partial}{\partial v_w} (u_c \cdot v_w + \log \sum_{w' \in V} \exp(u_c \cdot v_{w'}))$$

loss funct

wrt v_w

$$\sum_{w' \in V} \exp(u_c \cdot v_{w'})$$

$$= -u_c(1 - \hat{y}_w)$$

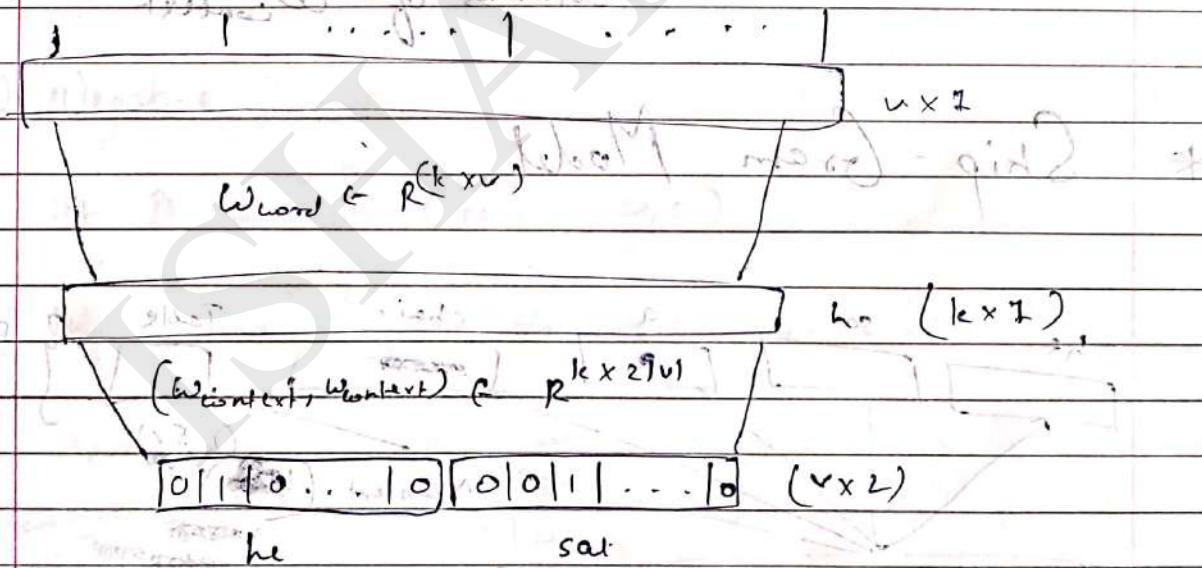
Update rule:

$$\Delta V_w = V_w - \eta \nabla V_w$$
$$= V_w + \eta u_c (1 - y_w)$$

Note - Current entry of V_w comes closer to entries in u_c

→ More than one word are given

multiple dialogue going on present
written in two sets p(m1, w1, set) p(m2, w2, set) based on
plausibility, relevance, context



Now,

number of words ($h \times 2$)

$\sum_{i=1}^{h-1} u_{ci}$

$h = \sum_{j=1}^{h-1} u_{cj} \rightarrow$ sum of columns in $W_{context}$ corresponding to one lot of words given

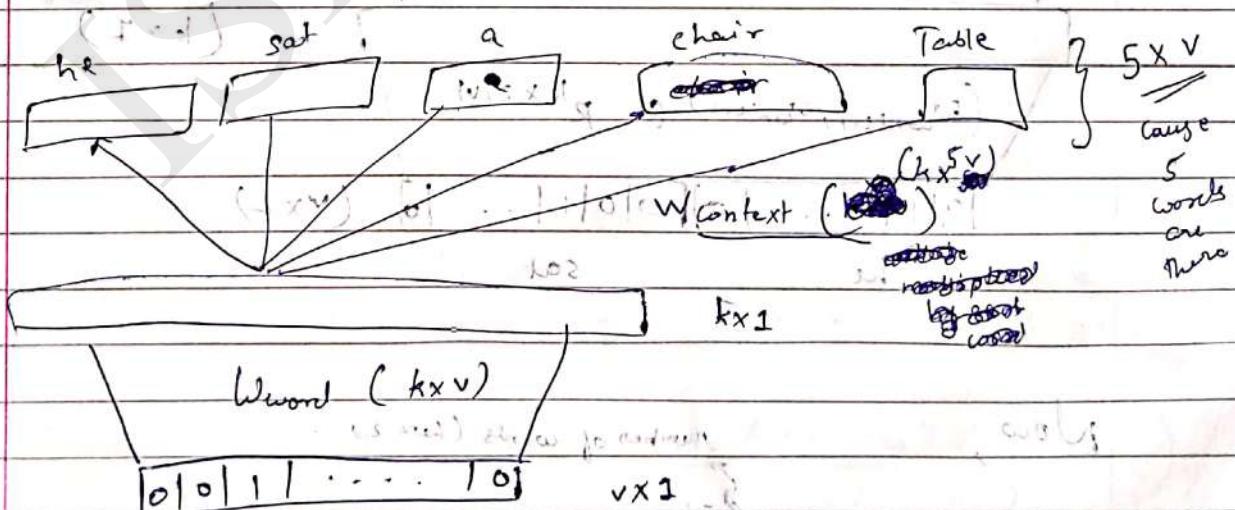
$$\begin{array}{c}
 A+B \\
 \overbrace{\left[\begin{array}{ccccc} -1 & \{0.5\} & 2 & -1 & 0.5 \\ 3 & -1 & -2 & 3 & -1 \\ -2 & 1.7 & 3 & -2 & 1.7 \end{array} \right]}^{\text{W context}} + \overbrace{\left[\begin{array}{ccccc} 2 & \{2\} & -2 & 3 & 0 \\ -1 & -2 & -1 & 3 & 1 \\ 1 & 3 & 0 & 0 & 0 \end{array} \right]}^{\text{W target}} = \left[\begin{array}{ccccc} 0 & \{0\} & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \end{array} \right]_{\text{set}}
 \end{array}$$

$$= \left[\begin{array}{c} 2.5 \\ -3 \\ 4.7 \end{array} \right] = \underline{A+B}$$

Now,

During Backprop, update all values of ω_{word} & only the participating columns of $\omega_{context}$

* Skip-Gram Model.



→ This model allows us to predict context words for a given word, words which can occur on the left or right of the given word.

$$\rightarrow L(\theta) = - \sum_{i=1}^{d+1} \log \hat{y}_{w_i}$$

If there was only 1 context word at output model is same as bag of words = 5 \times 5 dimension.

But here 5 words are there random. Thus loss function is addition of all loss.

Problem

- lot of computation at final layer

Solution

(1) Negative Sampling

let D be correct pair (w, c)

let D' be random pair incorrect pair (w, z)

$$u_c = c, v_w = w, u_z = z$$

correct context words $P(z|w, c)$

if sigmoid $\frac{1}{1 + e^{-x}}$

} dot product

$$v_w \quad u_c$$

now we will draw bayesian network of two words before soft c
 $p(z=1|w, c) = \sigma(u_c^T v_w)$ we draw diagram
 bayesian network

$$= \frac{1}{1 + e^{-u_c^T v_w}}$$

for all $(w, c) \in D$

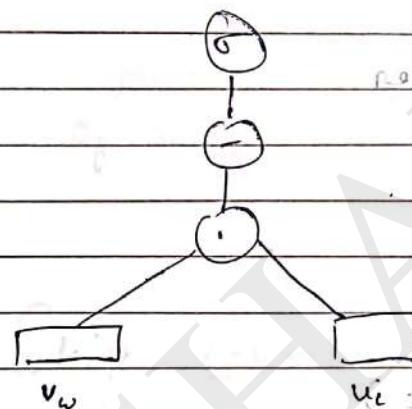
between thoughts to know instances & prior probability of it
 maximize $\prod p(z=1|w, c)$ and it works

is maximum over $(w, c) \in D$ result and choose (w, c) with max

prob. No. of mistakes

~~incorrect context words~~

Now,



now for $p(z=0|w, c) = 1 - \sigma(u_c^T v_w)$

$$= 1 - \sigma(u_c^T v_w)$$

$$= \frac{1}{1 + e^{u_c^T v_w}}$$

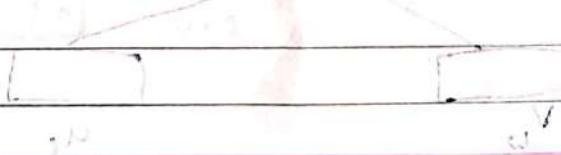
$$= \frac{1}{1 + \exp(u_c^T v_w)}$$

$$\frac{1}{1 + e^{u_c^T v_w}} \approx 1 - \sigma(u_c^T v_w)$$

$$\therefore p(z=0|w, c) = \sigma(-u_c^T v_w)$$

maximize $\prod p(z=0|w, c)$

$w, c \in D$



Total

$$\text{maximize}_{\theta} \prod_{(w, c) \in D} p(z=1 | w, c) \prod_{(w, r) \in D'} p(z=0 | w, r)$$

$$\text{maximize}_{\theta} \sum_{w, c \in D} \frac{1}{1 + e^{-u_c^T v_w}} + \sum_{w, r \in D'} \frac{1}{1 + e^{u_r^T v_w}}$$

$$= \text{maximize}_{\theta} \sum_{w, c \in D} \log \frac{1}{1 + e^{-u_c^T v_w}} + \sum_{w, r \in D'} \log \frac{1}{1 + e^{u_r^T v_w}}$$

$$\text{Final maximized } \sum_{w, c \in D} \log \sigma(u_c^T v_w) + \sum_{w, r \in D'} \log \sigma(u_r^T v_w)$$

bringing u_c & v_w

moving
 u_r & v_w
away

Number of samples in D' is k times more than D

→ Selection of (w, r)

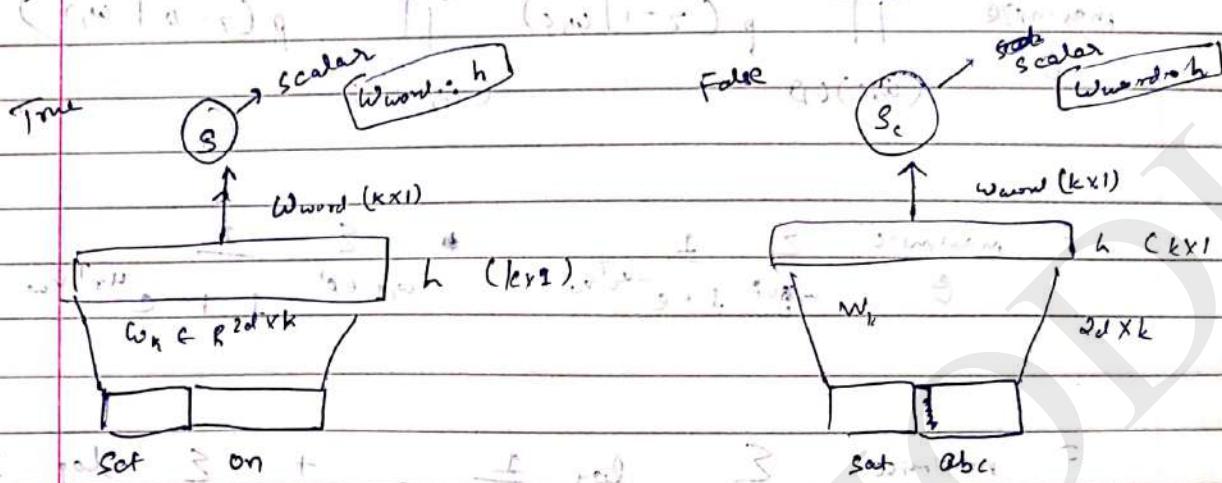
The random context word is drawn from a modified unigram distribution

$$\text{prob}(r) = \frac{\text{count}(r)^{\frac{3}{4}}}{(\text{count}(r) - 2) + 2} + \frac{2}{(\text{count}(r) - 2) + 2}$$

$$\text{count}(r) \sim \frac{\text{count}(r)}{N}$$

it is better than picking random because frequency of appearance is different.

(2) Contrastive Estimation



We want to maximize

(large difference)

$s - s_c + \text{margin}$ between them

Thus we maximize

$$s - (s_c + \text{margin}) \rightarrow \text{optimization}$$

Now, if $s > s_c + m$ don't do anything

$$\max \left[\min \left(0, s - (s_c + m) \right) \right] \quad \left\{ \begin{array}{l} \text{if } s < s_c + m \\ \text{loss function} \end{array} \right.$$

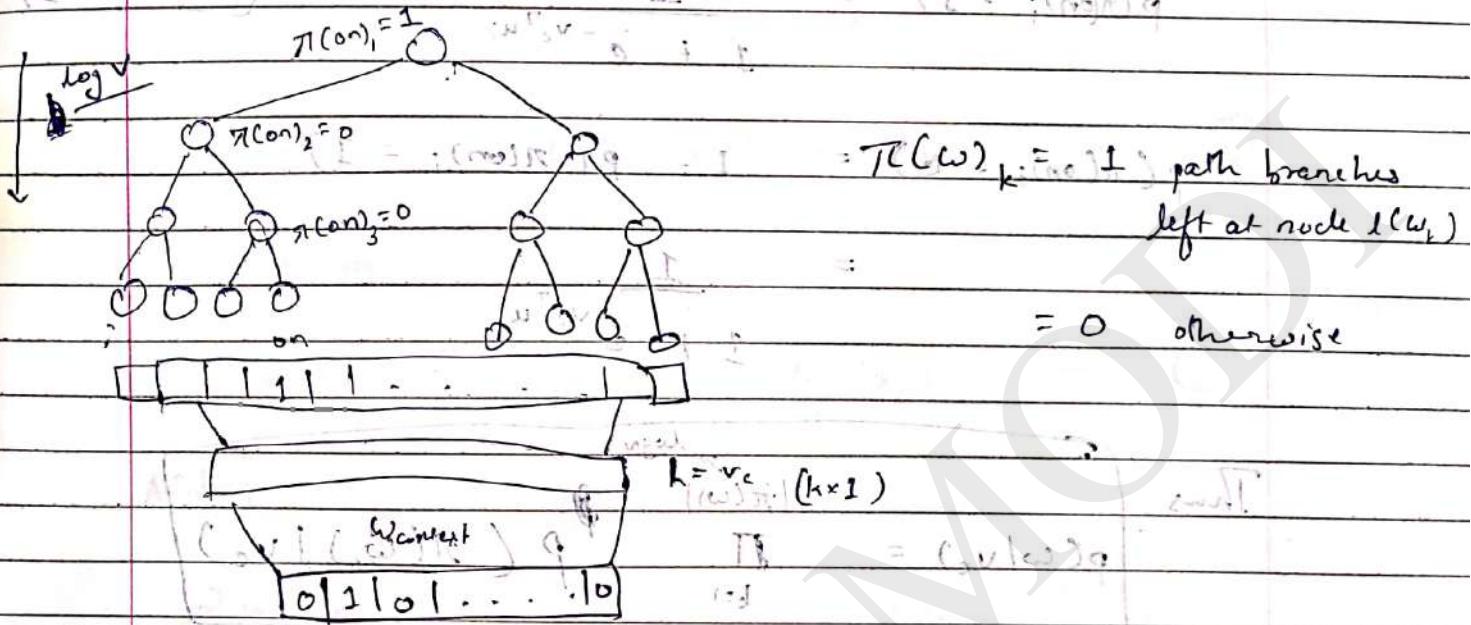
Here

if $s - (s_c + m) > 0$, loss = 0 (don't do anything)

if $s - (s_c + m) < 0$, loss = $(s - (s_c + m))$

Transferring learned model using multi-class at 71
transferring to same stage

(3) Hierarchical Softmax



number of nodes to be evaluated = log
value of V more branches

Now, number of nodes in the tree excluding leaf nodes are

$$u_i = v_i - 1$$

↓
no. of leaf nodes (first n nodes) (oval)

Thus the parameters of this network will be

w_{context} & u_1, u_2, \dots, u_{v-1}

every non-leaf node has a vector associated with it

→ Now,

$$\underline{i}^x = \underline{i}^y = (\dots)$$

$$p(\text{col } v_c) = \prod_k \pi(\pi(w_k) | v_c)$$

$$\begin{aligned} \text{g- } p(\text{con} | v_{\text{set}}) &= p(\pi(\text{con}_1 = 1 | v_{\text{set}})) \\ &\quad * p(\pi(\text{con}_2 = 0 | v_{\text{set}})) \\ &\quad * p(\pi(\text{con}_3 = 0 | v_{\text{set}})) \end{aligned}$$

Here

$$p(\pi(\text{on})_i = 1) = \frac{1}{1 + e^{-v_c^T u_i}}$$

Vectors u_i (ω)
are $k \times 1$

$$p(\pi(\text{on})_i = 0) = 1 - p(\pi(\text{on})_i = 1)$$

$$= \frac{1}{1 + e^{-v_c^T u_i}}$$

Thus,

$$p(\text{col} | v_c) = \prod_{k=1}^{\log n} P(\pi(\omega_k) | v_c)$$

total computations at activation
reduces from V to $\log V$

→ Constructing Binary Tree

* Glove Representations

Mixture of SVD based & Predictive based methods

→ Let,

x_{ij} be each entry in the cooccurrence matrix

$$p(j|i) = \underline{x_{ij}} = \underline{x_{ij}}$$

$$\sum_j x_{ij} (x_{ij})^T = (x_{ij}) q$$

$$(x_{ij}) q = (x_{ij})^T q$$

$$(x_{ij}) q = (x_{ij})^T q$$

Now, let's say

$$u_i^T v_j = \log(p(c_{ij})) \quad \text{--- equation 1}$$

$$u_i^T v_i = \log(x_{ii}) - \log x_i \quad \text{--- (1)}$$

Similarly for $p(c_{ilj})$

$$v_j^T u_i = \log(x_{ij}) - \log(x_{ji}) \quad \text{--- (2)}$$

Adding (1) & (2)

$$2u_i^T v_j = 2\log(x_{ij}) - \log x_i - \log x_j$$

$$u_i^T v_j = \log(x_{ij}) - \frac{1}{2}\log x_i - \frac{1}{2}\log x_j$$

$$\therefore u_i^T v_j + a + b = \log(x_{ij})$$

learnable parameters known

→ loss function

$$\min_{u_i, v_j, a, b} \sum \left(\underbrace{u_i^T v_j + a + b}_{\text{predicted value from model}} - \underbrace{\log x_{ij}}_{\text{value computed from query}} \right)$$

* Evaluating Word Representations

→ Semantic Relatedness
Comparing computer predictions with human perception

→ Synonym Detection

Term - derived
Candidates - imposed, believed, requested, correlated
 $\text{Similarity} = \arg\max_{v \in C} \cosine(v_{\text{term}}, v_c)$

→ Analogy

• Semantic Analogy

$v_{\text{sister}} = v_{\text{brother}} + v_{\text{grandson}}$ $\Rightarrow v_{\text{daughter}}$

• Syntactic Analogy

$v_{\text{works}} = v_{\text{worker}} + v_{\text{speaks}} + \dots + v_{\text{specialist}}$

* Relation Between SVD & Word2Vec

→ refer / notes / Creek10 / svdrelation.py