

* Principle Component Analysis

lets say,

x	y	z
1	1	1
0.5	0	0
0.25	1	1
0.35	2.5	1.5
0.45	1	2
0.57	2	2.1
0.62	1.1	1
0.75	0.75	0.75
0.72	0.86	0.87

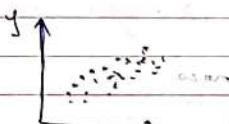
Here x & y are highly correlated thus if any one is ignored it wouldn't make much difference in result.

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.

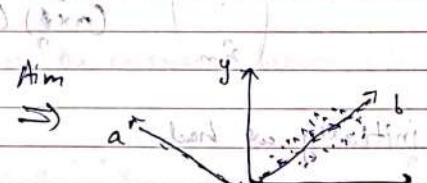
$$P_{yz} = \sum_{i=1}^n (y_i - \bar{y})(z_i - \bar{z})$$

$$(x \text{ axis}) \text{ variance} = \sum_{i=1}^n (y_i - \bar{y})^2, \text{ } (z \text{ axis}) \text{ variance} = \sum_{i=1}^n (z_i - \bar{z})^2$$

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

Ishan Modi

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

Now, representing \mathbf{x}_i using basis \mathbf{P}

and notation $\alpha_{ij} = \mathbf{x}_{i1} p_1 + \mathbf{x}_{i2} p_2 + \dots + \mathbf{x}_{in} p_n$

$$\text{Therefore } \mathbf{x}_{ij} = \alpha_{i1} p_1 + \alpha_{i2} p_2 + \dots + \alpha_{in} p_n$$

Using orthonormal basis

for orthonormal basis

orthogonal projection of \mathbf{x}_i onto

$$\Rightarrow \alpha_{ij} = \mathbf{x}_{i1}^T p_j$$

if no coordinate space basis

use \mathbf{x}_i^T standard basis p_j

$$\Rightarrow \alpha_{ij} = \mathbf{x}_{i1}^T p_j \quad \left\{ \begin{array}{l} \text{Dimension} = (1 \times n) \\ \text{row} \end{array} \right.$$

now \mathbf{x}_i^T has n entries

$$\left[\begin{matrix} 1 \\ \vdots \\ n \end{matrix} \right]$$

$$\left[\begin{matrix} 1 & \dots & 1 \end{matrix} \right] \left[\begin{matrix} 1 & \dots & 1 \end{matrix} \right]^T$$

Now for m vectors

$$\left(\begin{matrix} \hat{\mathbf{x}}_1 \\ \hat{\mathbf{x}}_2 \\ \vdots \\ \hat{\mathbf{x}}_m \end{matrix} \right) \xrightarrow{\text{Dimensions}} \left(\begin{matrix} m \\ 1 \\ \vdots \\ 1 \end{matrix} \right) \left(\begin{matrix} 1 & \dots & 1 \end{matrix} \right)^T \xrightarrow{\text{Dimensions}} (m \times n)$$

initially we had

$$\left(\begin{matrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_m \end{matrix} \right) \xrightarrow{\text{Dimensions}} \left(\begin{matrix} m \\ n \\ \vdots \\ n \end{matrix} \right) \xrightarrow{\text{Dimensions of rows}} \left(\begin{matrix} m \\ n \\ \vdots \\ n \end{matrix} \right)$$

$$\xrightarrow{\text{Dimensions of columns}} \left(\begin{matrix} m & n \\ 1 & \dots & 1 \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{matrix} \right) \xrightarrow{\text{Dimensions}} (m \times n)$$

but now we are getting m as new dimension

* Now here we got new dimension. Here we want to prove that these new dimensions have low covariance i.e. they are not interrelated

char. Mod. of these new dimensions makes it simple

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.

$$\tilde{X}^T (\tilde{X} \tilde{X}^T)_{\perp} = \tilde{X}^T (X P)^T (X P)_{\perp} = X^T X_{\perp} = 0.$$

Proof →

$$1^T \tilde{X} = 0 \text{ since columns are zero mean}$$

$$\text{now } \tilde{X} = X P \Rightarrow X^T \tilde{X} + 1^T \tilde{X}^T = 1^T X P \\ = 0 P = 0$$

(2)

Theorem

$X^T X$ is symmetric because $(X^T X)^T = X^T X$ as well

Proof → $(X^T X)^T = (X^T (X^T X))^T = 0 X^T X + I$ means it is diagonal matrix so it is symmetric if it is to be symmetric

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

Explanation

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

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

$$= \frac{1}{m} \sum_{k=1}^m x_{ki}^T x_{kj} \text{ or resp. entry } C_{ij} \text{ after } X^T X = 0$$

Similarly matrix is formed

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

Ishan Modi

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

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

DELUXE
PAGE NO.:
DATE:

CloudTech

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

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

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

mean shift and rotation be merged = P^T

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

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

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

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

Now, as seen earlier

$\frac{1}{m} \hat{X}^T \hat{X}$ is symmetric — Theorem 2

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

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

and $P^T \Sigma P =$ Diagonalization of gives diagonal matrix

Ishan Modi



Scanned with OKEN Scanner

The method of transforming data to new basis with

(1) low covariance

- proved

(2) high variance

- 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

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

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

all dimensions k dimensions

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

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

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

$$= \sum_{i=1}^m \sum_{j=k+1}^n \sum_{l=k+1, l \neq j}^n \alpha_{il} p_l^T \alpha_{il}$$

ShankModi

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

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

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

Now we can write $\sum_{j=k+1}^n u_i^T p_j (\sum_{j=k+1}^n u_i^T p_j) = p_{jk}$ which is non-zero if $i = k$

Now

lets say

$$\sum_{i=1}^m \begin{bmatrix} u_{11} \\ u_{12} \\ u_{13} \end{bmatrix} \begin{bmatrix} p_{11} & p_{12} & p_{13} \\ p_{21} & p_{22} & p_{23} \\ p_{31} & p_{32} & p_{33} \end{bmatrix} = \sum_{i=1}^m \begin{bmatrix} u_{mi} \\ u_{mi} \\ u_{mi} \end{bmatrix} \begin{bmatrix} p_{11} & p_{12} & p_{13} \\ p_{21} & p_{22} & p_{23} \\ p_{31} & p_{32} & p_{33} \end{bmatrix}$$

$$= \begin{bmatrix} (u_{11}^2 + u_{12}^2 + u_{13}^2) p_{11} & (u_{11}^2 + u_{12}^2 + u_{13}^2) p_{12} & (u_{11}^2 + u_{12}^2 + u_{13}^2) p_{13} \\ (u_{11}^2 + u_{12}^2 + u_{13}^2) p_{21} & (u_{11}^2 + u_{12}^2 + u_{13}^2) p_{22} & (u_{11}^2 + u_{12}^2 + u_{13}^2) p_{23} \\ (u_{11}^2 + u_{12}^2 + u_{13}^2) p_{31} & (u_{11}^2 + u_{12}^2 + u_{13}^2) p_{32} & (u_{11}^2 + u_{12}^2 + u_{13}^2) p_{33} \end{bmatrix}$$

$$(u_{11}^2 + u_{12}^2 + u_{13}^2) p_{11} + \dots + (u_{11}^2 + u_{12}^2 + u_{13}^2) p_{33}$$

$$(u_{11}^2 + u_{12}^2 + u_{13}^2 + \dots + u_{mi}^2) p_{11} + (u_{11}^2 + u_{12}^2 + u_{13}^2 + \dots + u_{mi}^2) p_{21} + \dots + (u_{11}^2 + u_{12}^2 + u_{13}^2 + \dots + u_{mi}^2) p_{31}$$

Ishan Modi

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

$$\begin{matrix} & \mathbf{x}_{11} & \mathbf{x}_{12} \\ \mathbf{x}_{21} & + & \mathbf{x}_{22} \\ & \mathbf{x}_{31} & \mathbf{x}_{32} \end{matrix} \quad \begin{matrix} & \mathbf{x}_{11}^2 = \mathbf{x}_{11} \mathbf{x}_{11}^T \\ & + \mathbf{x}_{21} \mathbf{x}_{21}^T \\ & + \mathbf{x}_{31} \mathbf{x}_{31}^T \end{matrix} \quad \begin{matrix} & \mathbf{x}_{11} \\ & \mathbf{x}_{21} \\ & \mathbf{x}_{31} \end{matrix} \quad \begin{matrix} & \mathbf{x}_{11} \\ & \mathbf{x}_{21} \\ & \mathbf{x}_{31} \end{matrix}$$

column column column column

Now

the $\sum_{i=1}^m p_i^T m c p_i$ is Norm [$\frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T \in \frac{\mathbf{x}^T \mathbf{x}}{m} = C$]

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

(p_1, p_2, \dots, p_n are eigenvectors of m)

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

P_1, P_2, \dots, P_n

statements above ref. from steven's book & amlab notes $\xrightarrow{\text{theorem}}$ Deep learning book $\xrightarrow{\text{proof}}$

p_j are smallest eigen vectors

Q13. Q13.

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

mean μ is zero \Rightarrow condition ② variance

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

variance

$$(\hat{\mathbf{x}}_i - \mu_i)^T (\hat{\mathbf{x}}_i - \mu_i) \xrightarrow{m} \hat{\mathbf{x}}_i^T \hat{\mathbf{x}}_i \xrightarrow{m} \frac{1}{m} \hat{\mathbf{x}}_i^T \hat{\mathbf{x}}_i \quad (\because \mu_i = 0)$$

zero mean

Ishan Modi

$$\begin{aligned}
 & \text{where } \lambda_i = \frac{1}{m} p_i^T (X^T X) p_i \quad \xrightarrow{\text{let } p_i = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}} \lambda_i = \frac{1}{m} p_i^T (\lambda_i p_i) \\
 &= \frac{1}{m} \lambda_i p_i^T p_i \\
 &= \frac{\lambda_i}{m} \underbrace{p_i^T p_i}_{\text{sum of squares}} \quad \xrightarrow{\text{normalized vector}} \frac{\lambda_i}{m} \underbrace{\|p_i\|_2^2}_{\text{sum of squares}} \\
 &= \frac{\lambda_i}{m} \quad \{ \text{eigen values} \}
 \end{aligned}$$

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

→ 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 \underline{A} - operation of size $n \times m$
 matrix A transforms a vector from n dimension to m dimension

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

v'

v''

\downarrow
 nxn

\downarrow
 mx1

$(n \times n)(n \times 1) \xrightarrow{\text{operation}} (n \times 1)$

Ishan Modi

Now (let us assume) that $\{v_i\}$ are such that $(n \times 1)$ vectors whose sum is zero, then we find

$\{v_i\}$ is basis of vector space \mathbb{R}^n and also $\{v_i\}$ is basis of vector space \mathbb{R}^m and $n < m$

then,

$$Av = \sum_{i=1}^k \sigma_i u_i = \text{non-zero}$$

$Av_i = \sigma_i u_i$ can be possible

$$[Av = v] \Rightarrow v = VU^T$$

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

$$[v = v] \Rightarrow v = VU^T$$

\therefore

$$Av = \sum_{i=1}^k \alpha_i Av_i \quad v \in \mathbb{R}^n \Rightarrow (n \times 1) \text{ vectors} \quad \leftarrow$$

$$\Rightarrow b \in \mathbb{R}^n \quad k \leq n$$

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

$$(VU^T)^T = U^T V = A^T A \quad \therefore$$

$$U^T V = V^T U = I_m$$

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

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

Since $v_i \perp u_i$ } common term k is used in summation
 $n - \text{dim}$ $m - \text{dim}$

→ Thus we can write

$$A_{m \times n} \times V_{n \times k} = \begin{bmatrix} U_{m \times k} \\ \vdots \\ U_{m \times k} \end{bmatrix} \times \begin{bmatrix} X_1 & X_2 & \cdots & X_k \end{bmatrix} = \sum_{i=1}^k U_{m \times k} X_i = D$$

$\left(\begin{array}{c} k \\ m \end{array} \right) \text{ vectors}$

$\left(\begin{array}{c} k \\ m \end{array} \right) \text{ vectors}$

diagonal matrix

Ishan Modi

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

$$\therefore A^T A = (U^T \Sigma V)^T (U^T \Sigma V)$$

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

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

$$\begin{bmatrix} A \\ m \times n \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ u_1 & u_2 & \dots & u_k \\ \downarrow & \downarrow & \ddots & \downarrow \end{bmatrix}_{m \times k} = \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_k & \\ & & & \Sigma_{k \times k} \end{bmatrix}_{k \times k} \begin{bmatrix} v_1 & & & \\ \leftarrow v_1 \rightarrow & \ddots & & \\ & & v_2 & \\ \leftarrow v_2 \rightarrow & & \ddots & \\ & & & v_k \end{bmatrix}_{k \times n}$$

$$= \sum_{i=1}^k \sigma_i u_i v_i^T$$

Ishan Modi

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

SVD

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

$$\tilde{A} = U_{:,k} \Sigma_{l,k} V^T$$

(first k columns of U & first k rows of V)

Example of PCA & SVD



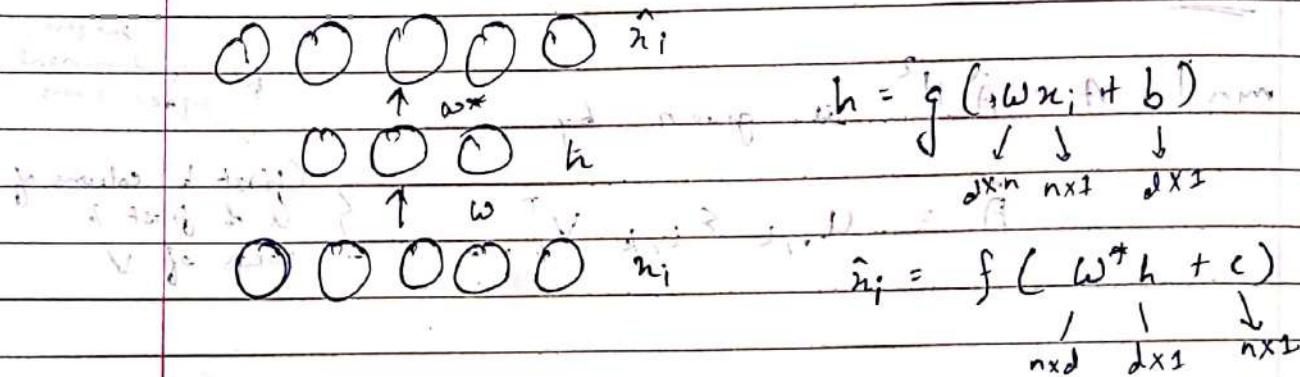
- These images are called eigenfaces and form a basis for representing any face in our database
- In other words, we can now represent a given image (face) as a linear combination of these eigen faces
- In practice, we just need to store p_1, p_2, \dots, p_k (one time storage)
- Then for each image i we just need to store the scalar values $\alpha_{i1}, \alpha_{i2}, \dots, \alpha_{ik}$
- This significantly reduces the storage cost without much loss in quality

$$\sum_{i=1}^{16} \alpha_{i1} p_i$$

Ishan Modi

WEEK 7

* Auto Encoders



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

$$\|x_i - \hat{x}_i\|^2$$
 is minimized.

→ if $\dim(h) \geq \dim(x_i) \Rightarrow$ Over complete auto encoders
if $\dim(h) < \dim(x_i) \Rightarrow$ Under complete auto encoders
↓ as in above case

→ If the x_i nodes have real numbers then

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

If x_i nodes have binary numbers 0, 1 then

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

f-in
of different
cases

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

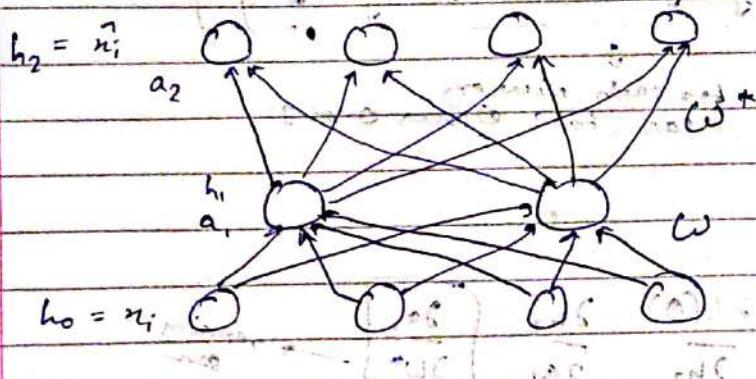
Ishan Modi

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

$$\Rightarrow \min_{w, w^*, c, b} \sum_{m=1}^M (\hat{x}_m - x_m)^T (\hat{x}_m - x_m)$$

→ Backprop (for Real x_i)

$$L(\theta) = (\hat{x}_i - x_i)^T (\hat{x}_i - x_i)$$



Now earlier we did backprop with cross entropy loss

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

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

Now

$$\frac{\partial L(\theta)}{\partial h_2} = \frac{\partial L(\theta)}{\partial \hat{x}_i} = \sum_j (\hat{x}_{ij} - x_{ij})^T (\hat{x}_{ij} - x_{ij}) \quad \begin{cases} \xrightarrow{\text{scaled}} \\ \xrightarrow{\text{vector}} \end{cases}$$

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

Ishan Modi

→ Back Propagation (for binary a_{ij})

0 0 0 0
Classified (0 - 00)

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

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 product term is like 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

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

$$\frac{\partial L(\theta)}{\partial a_{ij}} = \frac{\partial L(\theta)}{\partial h_2} \frac{\partial h_2}{\partial a_2} \frac{\partial a_2}{\partial a_{ij}} = -x_{ij} + 1 - x_{ij}$$

$$\frac{\partial L(\theta)}{\partial a_{ij}} = \frac{\partial h_2}{\partial a_{ij}} = \sigma(a_{ij})(1 - \sigma(a_{ij}))$$

Ishan / Nodi

* 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)
 - (4) Generalization \rightarrow $x_{ij} = \frac{1}{\sqrt{m}} \left(n_{ij} + \frac{1}{m} \sum_{k=1}^m x_{kj} \right)$
- Normalization zero mean

$$\rightarrow n_{ij} = \frac{1}{\sqrt{m}} (\text{zero mean})$$

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

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

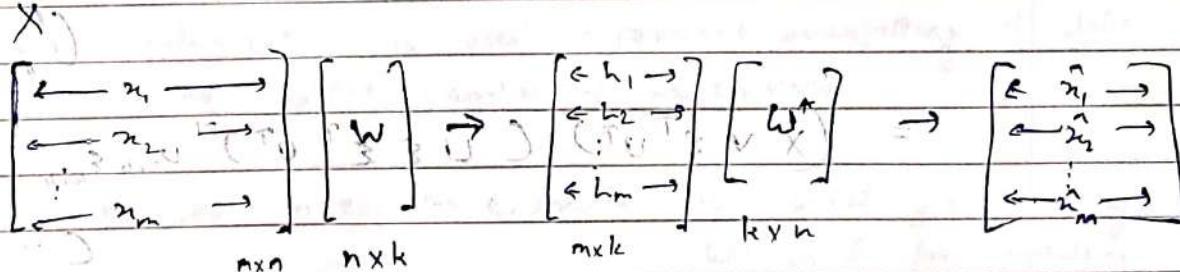
prove it for n_{ij} \rightarrow covariance matrix

X is probability constant $\forall i \in N$

→ Now loss function is

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

$$\min_{W^*} \|X - HW^*\|_F^2$$



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

Shan Modi

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 .

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

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

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

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

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

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

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

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

Ishan Modi

$$\text{Consider } \mathbf{x} = \mathbf{U} \mathbf{\Sigma}^{-1} \mathbf{I}_{\mathbf{z}, k} + \mathbf{\epsilon}_{\mathbf{z}, k} \text{ and so } \mathbf{x} = \mathbf{U} \mathbf{\Sigma}^{-1} \mathbf{I}_{\mathbf{z}, k} + \mathbf{\epsilon}_{\mathbf{z}, k}$$

$$= \mathbf{X} \mathbf{V} \mathbf{I}_{\mathbf{z}, k} \quad (\mathbf{\Sigma}^{-1} \mathbf{I}_{\mathbf{z}, k} = \mathbf{\epsilon}_{\mathbf{z}, k})$$

$$= \mathbf{X} \mathbf{V} \mathbf{I}_{\mathbf{z}, k}$$

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

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

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

$$\text{Ans: } \frac{\partial L}{\partial \mathbf{w}} = \text{nderivative of loss functions} + \lambda \mathbf{w}$$

(2) Weights tying

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

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

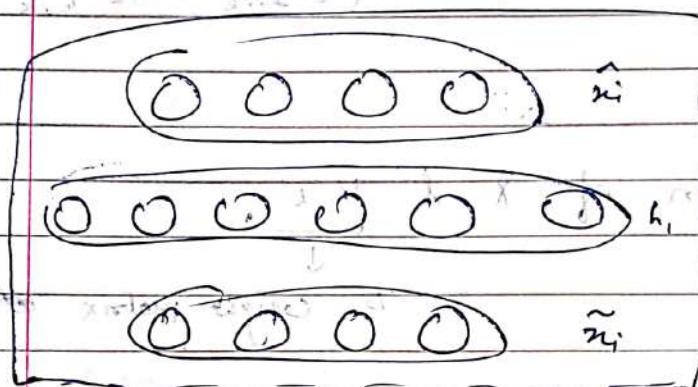
in this method $\mathbf{w}^* = \mathbf{w}^T$ learned by it

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

\mathbf{w} is a row vector containing all the values for a particular row of \mathbf{w}^* .
The row \mathbf{w} for a particular instance is used to calculate the error in that row for finding the old

Ishan Modi

* Denoising Autoencoders (Type of Regularization)



Autoencoder

z_i-hat is denoising

(original input z_i corrupted by noise)

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

Now,

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

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

loss function (2)

models to fit not $\left(\frac{1}{m} \sum_{i=1}^m (\tilde{z}_{ij} - z_{ij})^2 \right)$ in reconstruction error

thus, since it is overcomplete autoencoder
the \tilde{z}_i won't go in \tilde{z}_i straight away because

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

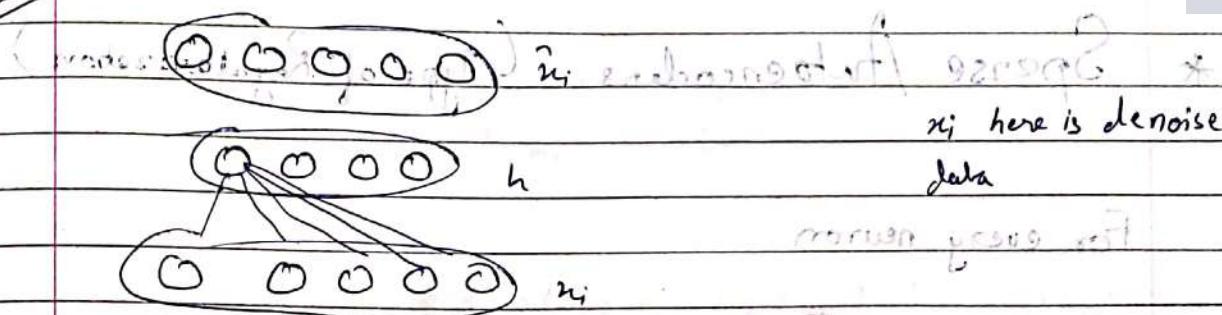
reconstruction is done when some loss is

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

Ishan Modi

Experiment

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



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

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

Maximum occurs when x_i is aligned because at this point sigmoid output would be equal to 1.0 (from $\sigma(z) = 1$).

It's easy to see this at each of the 3 neurons.

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

at which w_i stands and when x_i is normalized by $\|x_i\|$ or $x_i^T x_i = 1$.

$w_j^T x_i$ dot product is maximum when x_i is in same direction.

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

Each input here helps respective neuron to fire

$$w_1, \underbrace{w_2}_{\text{and so on}}, \underbrace{w_n}_{\text{in total } n \text{ weights}}$$

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

This experiment is done after training
by Sharad Modi

set of all inputs

* Sparse Autoencoders (Type of Regularization)

↳ for every neuron

For every neuron

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

↓

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 times. Since we have taken average it will be close to zero and hence sparse.

$\hat{p}_e = \frac{1}{m} \sum_{i=1}^m g(\omega^T, u_i)$ (at an average every neuron)

↳ $\hat{p}_e = 0.005$ (at an average every neuron)

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

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

↳ from the left side

$$\frac{\partial \Omega(\theta)}{\partial \omega} \quad \text{at minimum}$$

(minimum value)

Ishan Modi

$$\frac{\partial \Omega(\theta)}{\partial \omega} = + \frac{\partial \Omega(\theta)}{\partial p_i} \times \frac{\partial p_i}{\partial \omega}$$

Now,

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

$$\frac{\partial p_i}{\partial \omega}$$

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

(softmax is now considered to be constant and we

Now

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

scalar

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

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

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

vector

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

$$\rightarrow \frac{\partial p_i}{\partial w} = \frac{\partial}{\partial w} \left(\frac{1}{m} \sum_{i=1}^m g(w n_i + b) \right)$$

matrix
chain Mod

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

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

→ Final Derivative

$$\frac{\partial \mathcal{L}(0)}{\partial w} = \frac{1}{m} \sum_{i=1}^m \left(\underbrace{\begin{pmatrix} -p & (1-p) \\ \hat{p}_x & (1-\hat{p}_x) \end{pmatrix}}_{\text{vector from } \textcircled{1}} \odot \underbrace{g'(w_n; + b)}_{\text{vector}} \right) \times \underbrace{n_i^T}_{\text{vector}}$$

* 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)

$$(n, d + 1) \times (k, n) \in \mathbb{R}^{n \times k}$$

$$\mathcal{L}(0) = \sum_{j=1}^n \sum_{i=1}^k \left(\frac{\partial h_k}{\partial x_j} \right)^2$$

(Contractive loss)

Ishan Modi

Page 6



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

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

$\frac{\partial L(\theta)}{\partial \theta} = 0$ means don't capture variations in data

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

$L(\theta) = \sum_i (y_i - \hat{y}_i)^2$ (y_i changes with respect to x_i)

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

(mean value and slope doesn't change with respect to x_i)

Trade off - can't capture only very important variations in data

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

($L(\theta) - \alpha(\theta)$) = $L(\theta)$ (minimizing $L(\theta)$)

Limitations of linear regression are less flexible and often with

restrictions on the parameters are given to fit

Linear model

No short generalization

Only few training data

Only few training data

Model

$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

Model

Input

Output

Ishan Modi

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)

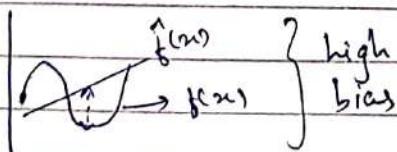
→ Formula

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

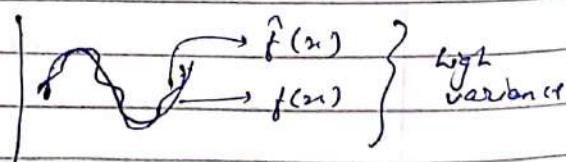
~~Precision~~ (~~Var~~) ~~(~~Efficiency~~)~~

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

* Bias - Variance Trade off



Simple model



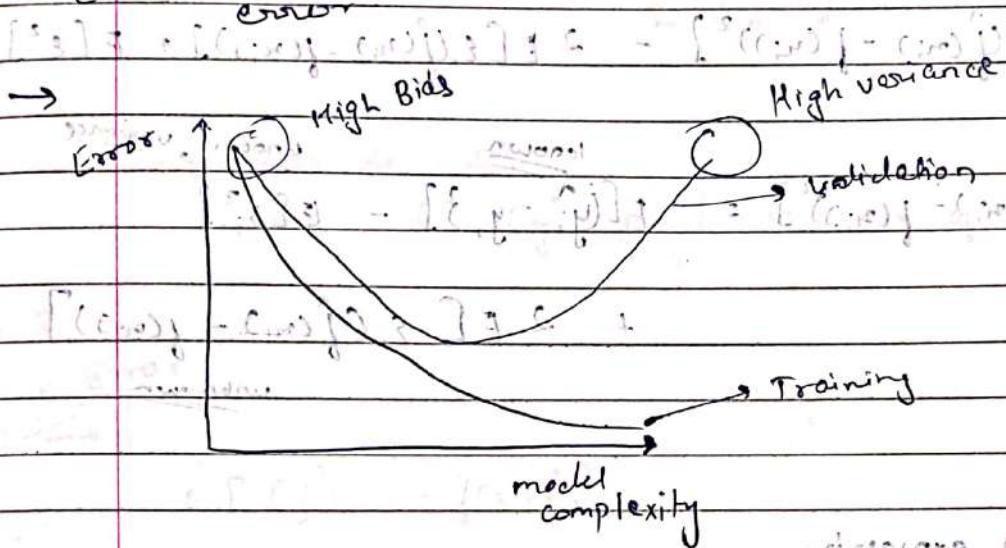
Complex model

shanModi

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 & on testing point :

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

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

\rightarrow Mathematically

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

\therefore we know 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

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

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

$$\therefore y_i = f(x_i) + \epsilon_i$$

Ishan Modic we assume that we know $f(x_i)$ but not $f(x_i)$

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

Test Error

Now, above expression

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

$$\left(\frac{1}{m} \sum_{i=n+1}^{n+m} (y_i - \hat{y}_i)^2 \right) \text{ for test data}$$

this is not exact (E)
but it is value

Similarly,

$$E[\varepsilon_i^2] \approx \left(\frac{1}{m} \sum_{i=1}^n \varepsilon_i^2 \right)$$

$$y_i = f(x_i) + \varepsilon_i \text{ for test data}$$

($\hat{f}(x_i)$ 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)

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

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

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

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

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

Thus, Model covariance = 0 (in simple words)

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

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

\therefore Error depends on Train Error

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

$$= \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 + (\text{train error})^2, \quad i=1 \quad n$$

Train Error

\rightarrow Model is good or bad

$$E[(\hat{y}_i - y_i)^2]$$

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

Learning error depends on training items using covariance ($\varepsilon_i, (\hat{y}_i - y_i)$)
independent in different epochs.

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

Ishan Modi

accuracy metric for boosters



→ Earlier

$E \sum_i \epsilon_i (f(\mathbf{x}_i) - f(\mathbf{y}_i))$ was unknown

Now,

using Stein's lemma

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

This value is high when $\frac{\partial f(\mathbf{x}_i)}{\partial y_i}$

is high and it is high

$\frac{\partial f(\mathbf{x}_i)}{\partial y_i}$ 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 we minimize the fitted error for model

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

concept of regularization

Ishan Modi

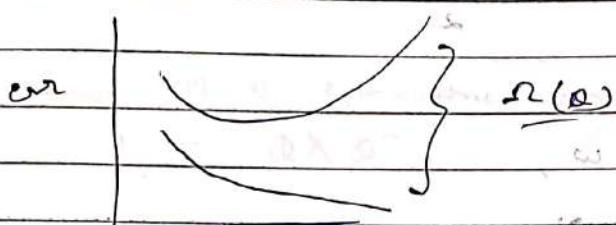
instead of train errors

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

$$(H^T(\omega))^{-1} = (\omega) \Rightarrow H^T(\omega) = (\omega)^{-1}$$

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

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



* L2 Regularization

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

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

$$\nabla \tilde{L}(\omega) = \{\nabla L(\omega) + \alpha \omega\}^T$$

update

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

Now $\omega = \tilde{\omega}^*$ (Let's find $\tilde{\omega}^*$)

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

Loss
without
regularization

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

Ishan Modi

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

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

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

\Rightarrow now do it 9

$$L(\omega) = \underbrace{L(\omega^*)}_{0} + \underbrace{(h - \omega^*)^T \nabla L(\omega^*)}_{0} - (\because \nabla L(\omega^*) = 0 \text{ from } ①)$$

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

differentiating wrt ω ,

$$\nabla L(\omega) = \underbrace{\nabla L(\omega^*)}_{0} + \underbrace{\cancel{h^T \nabla L(\omega^*)}}_{0} - (\because \nabla L(\omega^*) = 0 \text{ from } ①)$$

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

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

$$\|H\omega\|^2 + (\omega)$$

Now,

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

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

$$+ (\omega)^T V$$



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

① $\Rightarrow (\omega)^T V = 0$ \Rightarrow last term is 0 \Rightarrow $\omega = \tilde{\omega}$
 regularized loss

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

∴ to understand m

Ishan Modi

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

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

Note:

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

→ Assume H is symmetric matrix ($\alpha \neq 0$)

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

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

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

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

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

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

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

Diagonal matrix D

$$\tilde{\omega} = Q D Q^T \text{ (orthogonal transformation)}$$

↓ ↓ ↓
rotates scales rotates
 ω^* ω^* $\tilde{\omega}$

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

The second problem is now

Ishan Modi

$$\mathbf{D} = \begin{bmatrix} \frac{\lambda_1}{\lambda_1 + \alpha} \mathbf{H}^T \mathbf{H} & \frac{\lambda_1}{\lambda_1 + \alpha} \mathbf{H}^T \mathbf{y} \\ \frac{\lambda_2}{\lambda_2 + \alpha} \mathbf{H}^T \mathbf{H} & \frac{\lambda_2}{\lambda_2 + \alpha} \mathbf{H}^T \mathbf{y} \\ \vdots & \vdots \\ \frac{\lambda_n}{\lambda_n + \alpha} \mathbf{H}^T \mathbf{H} & \frac{\lambda_n}{\lambda_n + \alpha} \mathbf{H}^T \mathbf{y} \end{bmatrix}$$

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

$$\begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & & \lambda_1 \\ & \frac{1}{\lambda_2 + \alpha} \mathbf{H}^T \mathbf{H} & & & \lambda_2 \\ & & \ddots & & \vdots \\ & & & \frac{1}{\lambda_n + \alpha} \mathbf{H}^T \mathbf{H} & \lambda_n \end{bmatrix} = \mathbf{D}^{-1}$$

Note

$$\mathbf{w}, \mathbf{D} \in \mathbb{R}^{n \times (n + \alpha)}$$

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

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

A similar logic

* Dataset Augmentation

Increasing the amount of data

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

Shan Modi

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

$\theta = \theta_{\text{enc}} \cup \theta_{\text{dec}}$ (parameters used in encoder and decoder)

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

$$[\epsilon \sim \mathcal{N}(0, \sigma^2)] \quad \text{Gaussian distribution}$$

$(\mu + \sigma \epsilon) \leftarrow (\text{zero mean}, \text{unit variance})$

$\tilde{x}_i = x_i + \epsilon_i \quad \theta = \theta_{\text{enc}}$ (or (zero mean, multivariate))

$$\hat{y} = \sum_{i=1}^n w_i x_i \quad \theta = \theta_{\text{enc}}$$

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

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

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

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

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

of expectation of products of independent random variables is zero (i.e., $E[AB] = E[A]E[B]$)

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

Ishan Modi

Here in second term ~~presented parameter~~

$(\hat{y} - y) \cdot \epsilon_i$ & ϵ_i are independent

original values when noise was not added

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

Covariance (co-variance) total of smooth weight

A in third term

$$\text{Covariance} \cdot (\sum w_i \epsilon_i)^2 \Rightarrow (a + b + c)^2$$

(since there is no cross)

$$(\epsilon_i \epsilon_j)^2 \approx 0 \quad ; \text{ independent}$$

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

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

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

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

since w_i is not random variable

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

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

↓ same as L regularization

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

Input
Ishan Modi

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

At the output we have

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

\downarrow
 original predicted
 (correct)

$$(\omega - \omega)^T \nabla + \epsilon \omega = \omega$$

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

$$p_i = \{0, 0, 1, 0, \dots\}$$

Assuming there are 9 elements

$$p_i = \left\{ \frac{1+\epsilon}{9}, \frac{\epsilon}{9}, 1 - \frac{\epsilon}{9}, \frac{1-\epsilon}{9}, \dots \right\}$$

ϵ = noise term

$$q_i = \text{remains as it is } p_i + \epsilon$$

* 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 w_{t+1} is propagated plus ϵ

$$w_{t+1} = w_t + \eta \nabla w_t (t+1) - \epsilon$$

$$= w_0 + \eta \sum_{i=1}^t \nabla w_i - \epsilon$$

maximum of them

$$w_{t+1} \geq w_0 + \eta t \epsilon$$

Chand Modi

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, $H = Q\Lambda Q^T$

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

If ω_{t-1} is now near ω^* , then ω_t is also near ω^* .

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

- 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

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

Ishan Modi

* 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 squared $\bar{\epsilon}_k^2$ (Here the number of models is k)

$$\frac{1}{k} \sum_i \bar{\epsilon}_i^2$$

Expectation over the inputs

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

equal to variance of average of errors

$$\text{where } \text{var} = \frac{1}{k^2} E \left[\sum_{i=1}^k \sum_{j=1}^k \bar{\epsilon}_i \bar{\epsilon}_j \right] + \frac{1}{k^2} \sum_{i=1}^k \sum_{j=1}^k \bar{\epsilon}_i^2$$

$$= \frac{1}{k^2} \left[\sum_i E(\bar{\epsilon}_i^2) + \sum_{i \neq j} E(\bar{\epsilon}_i \bar{\epsilon}_j) \right]$$

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

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

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

Ishan Modi covariance

NoteIf errors are independent $C = 0$

method is useful (s and m) passed still gradient to
mse = V when backpropogation
then this will be used in back propagation

If errors are correlated $C = V$

$$\text{mse} = V$$

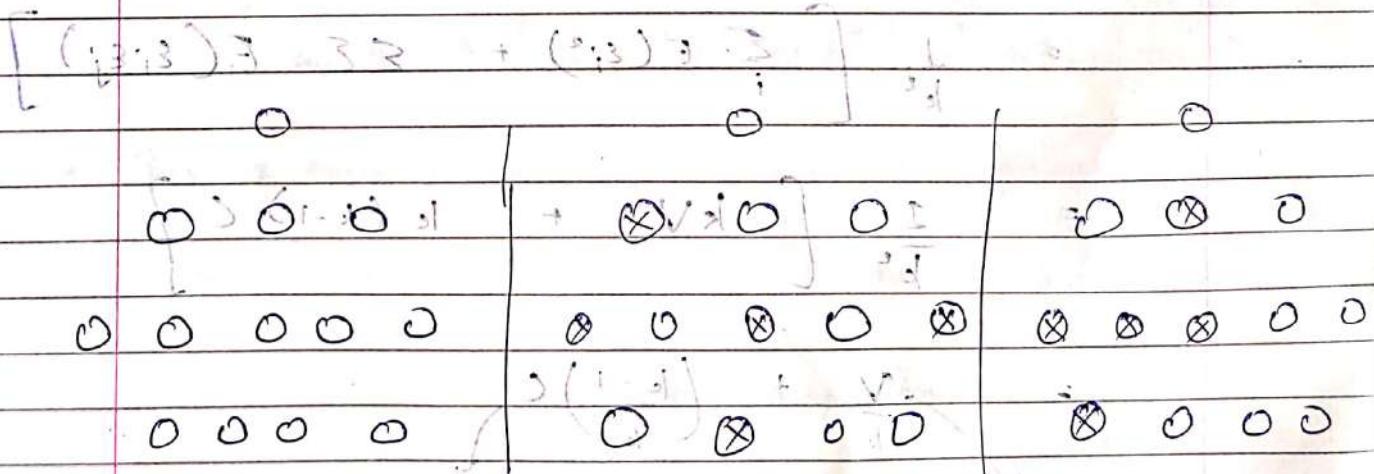
standard

Thus using averaging helps reduce errors in output

* Dropout

Let w be the weight matrix ~~matrix because now forget~~

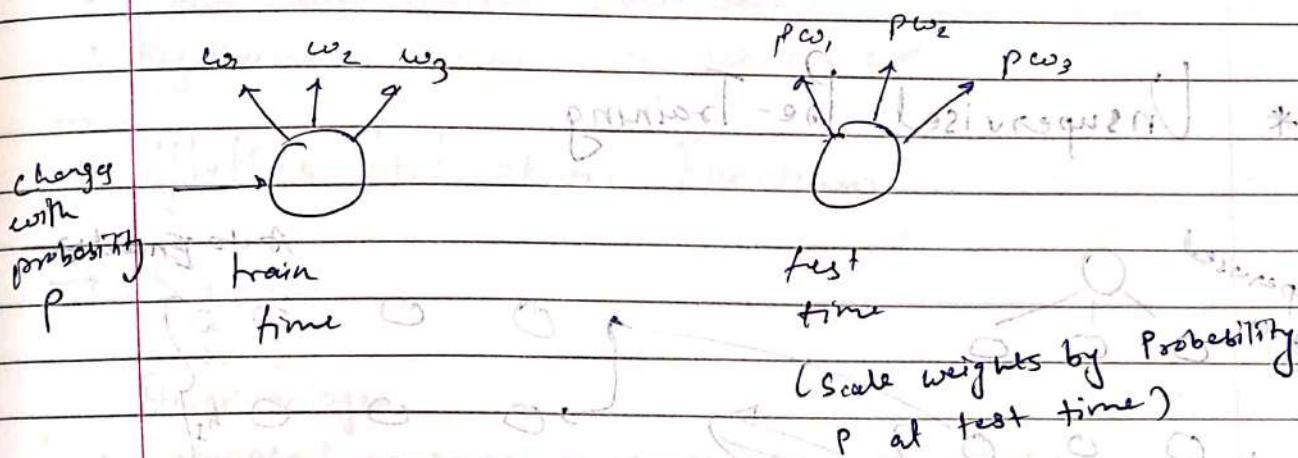
- Dropout means to eliminate neurons in layers with probability p form completely new model train for one epoch on that model. update w .
- Now eliminate from original to get new model I use updated w and train for 1 epoch



Ishan Modi

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

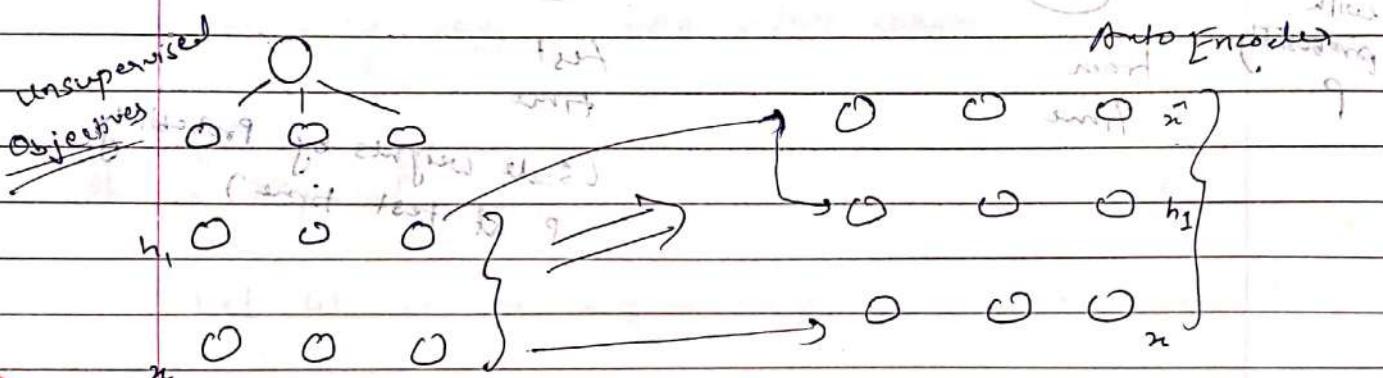
These neurons will be robust (due to fading).



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

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

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.

Ishan Modi

~~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

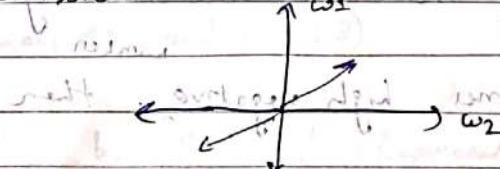
- Saturated neurons cause gradient to vanish

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 ($\exp(x)$)

Shankar Modi

→ Tanh.

Disadvantage

- Gradient vanishes at saturation after optimal point.
- Computationally expensive as it needs tanh function.

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 a ReLU layer if x_i is x_{-i}

2nd ordering techniques next will pass

Mostly used in CNNs due to increased speed

Only problem \rightarrow (dead neuron problem)

most neurons will (not) receive backprop in back

$$\text{Case } 2: \frac{\partial \text{ReLU}}{\partial x} = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases} \quad \left\{ \begin{array}{l} \text{derivatives} \\ \text{constant and 0} \end{array} \right.$$

$$\therefore \text{In practice a large fraction of ReLU units can die}$$

if the learning rate is set too high

If bias becomes high negative, then

$$w_1x_1 + w_2x_2 + b < 0$$

(dead) neurons phenomenon

$\therefore \frac{\partial \text{ReLU}}{\partial x} = 0$ & weights are not updated

Ishan Modi

→ Leaky ReLU (gradient, non-saturation) $f(x) = \max(0.01x, x)$ (solution to dead neurons)

$$f(x) = \max(0.01x, x)$$

→ Parametric ReLU

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

less sensitive to parameter changes but higher performance than ReLU

→ Exponential Linear Unit

$$f(x) = x \text{ if } x > 0$$

(modified gradient saturation for weights)

$$= \alpha e^x - 1 \text{ if } x \leq 0$$

→ Maxout Network (more power of representation)

$$\max(C_1 w_1^T x + b_1, C_2 w_2^T x + b_2) \quad \text{--- (1)}$$

variant of ReLU $\max(0, x)$ at $x = 0$ (parallel lines)

$$\begin{cases} w_1 = 1, b_1 = 0 \\ w_2 = 1, b_2 = 0 \end{cases} \quad \left. \begin{array}{l} \text{Representing (1) as 2 (2)} \\ \text{using multiple lines} \end{array} \right\}$$

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

$$\begin{cases} w_1 = \alpha, b_1 = 0 \\ w_2 = 1, b_2 = 0 \end{cases} \quad \left. \begin{array}{l} \text{Representing (1) as 2 (3)} \\ \text{using multiple lines} \end{array} \right\}$$

Ishan Modi

* 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 Variance is low when plotted.
Saturation occurs and thus the vanishing gradient problem persists.

Principled Way

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

Ishan Modi

$$v(XY) = \underbrace{(\bar{x})^2 v_{\text{ar}}(Y) + (\bar{y})^2 v_{\text{ar}}(X) + v_{\text{ar}}(Y) v_{\text{ar}}(X)}_{\text{Formula}}$$

PAGE NO.:
DATE:

$$v(S_{1,1}) = v_{\text{ar}}\left(\sum_{i=1}^n w_i x_i\right) = \sum_{i=1}^n v_{\text{ar}}(w_i x_i)$$

Now $(\sum w_i)$ now rest $\leftarrow (w_i)$ now \sum

$$\text{variance minimization} = \sum_{i=1}^n \left(E[w_i] \right)^2 v_{\text{ar}}(x_i)$$

$$\text{and } \dots \text{ now } (c, 0) \text{ or } + (E[x_i])^2 v_{\text{ar}}(w_i)$$

$$+ v_{\text{ar}}(x_i) v_{\text{ar}}(w_i)$$

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

lets say x_i & w_i

both wt are 0 mean 1

standard deviation

$$v_{\text{ar}}(x_i) = v_{\text{ar}}(x), v_{\text{ar}}(w_i) = v_{\text{ar}}(w)$$

$$\left(\frac{\sum}{n}\right) = w \text{ same}$$

$$= \sum_{i=1}^n v_{\text{ar}}(x_i) v_{\text{ar}}(w_i)$$

$$(S)_{\text{now}} = (w)_{\text{now}} = \sum_{i=1}^n v_{\text{ar}}(w) v_{\text{ar}}(x)$$

$$\therefore V(S_{1,1}) = (n v_{\text{ar}}(w))(v_{\text{ar}}(x)) - \textcircled{1}$$

Now,

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

$$= n V(S_{1,1}) v_{\text{ar}}(w_2)$$

Now we can see that if w_2 is very large then $S_{2,1}$ will be large
(and so $V(S_{2,1})$ will be large) substituting \textcircled{1}

$$V(S_{2,1}) \propto \left[\underbrace{[n v_{\text{ar}}(w)]^2 v_{\text{ar}}(x)}_{\text{if}} \right]$$

thus low values & high values

would lead to a variance of $S_{2,1}$
or very high variance of $S_{2,1}$

Ishan Modi



$$\text{Var}(s_{ki}) = [n \text{Var}(w)]^k \text{Var}(x_i) = (0.2)^{k+1}$$

if

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

→ If $(z)_{avg} = 0$ and $\text{var}(z) = 1$ and we take weights such that

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

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

thus

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

$$\text{Now } = n \text{Var}(\frac{z}{\sqrt{n}}) \quad \left. \begin{array}{l} \text{var}(az) = a^2 \text{var}(z) \\ (z, w)_{\text{cov}} = 0 \end{array} \right\}$$

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

$$= (z)_{\text{cov}}^2 = (0.2)^{\text{cov}}$$

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

Applicable for Sigmoid & Tanh (0.2)^{\text{cov}}

$$(0.2)^{\text{cov}} (0.2)^{\text{cov}}$$

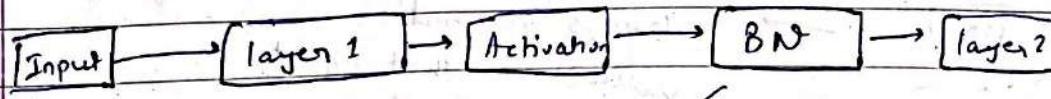
→ For ReLU (since it is 0 for negative values that is approx half of time)

$$w = \frac{[z]_{\text{cov}}^2 [(w)_{\text{cov}}]}{\sqrt{\frac{n}{2}}} \propto (0.2)^{\text{cov}}$$

Ishan Modi
Author of a book on
Machine Learning
for beginners and
researchers



* Batch Normalization



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

Note - Batch Normalization as the name suggests is applied on the ^{each} ~~entire~~ batch that passes through the network.

Also, batch norm is applied across individual features of the input & not the entire input

e.g -

$x_1 \quad x_2 \quad x_3 \quad x_4$

Let say we have a batch with n inputs

{	input 1 \rightarrow	0.3	0.9	0.7	0.2
	input 2 \rightarrow	0.4	0.7	0.2	0.1
	:	:	:	:	:
	input n \rightarrow	0.9	0.8	0.8	0.4

Each column in the above case should be normalized separately

Aman Modi

→ Formula

As seen earlier, s_{ki} is the result of dot product of ~~over~~ input & weights.

$s_{ki} \rightarrow$ output of i^{th} neuron of k^{th} layer

$A_{ki} \rightarrow$ output of i^{th} neurons of k^{th} layer after activation

$$\bullet A_{hi} = \frac{A_{ki} - E[A_{ki}]}{\sqrt{\text{var}(A_{ki})}} \quad \text{o mean, unit variance}$$

The above expression is differentiable so we can easily backpropagate the error

- Now, if we don't want normalization because normalizing values is having not allowing us to learn effectively we can convert batch norm to original

$$y^{(k)} = \gamma^k \cdot \hat{A}_{ki} + \beta^{(k)}$$

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

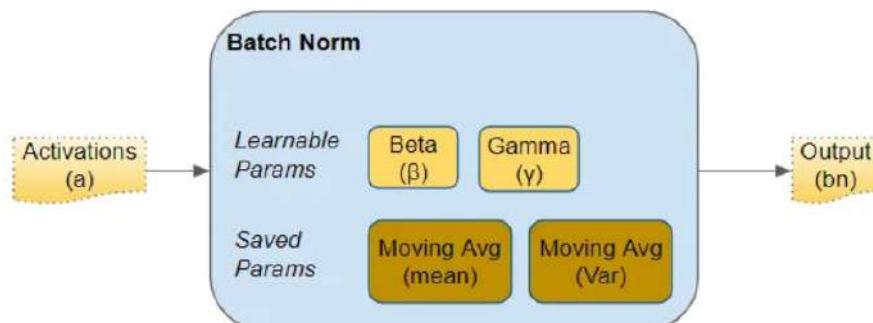
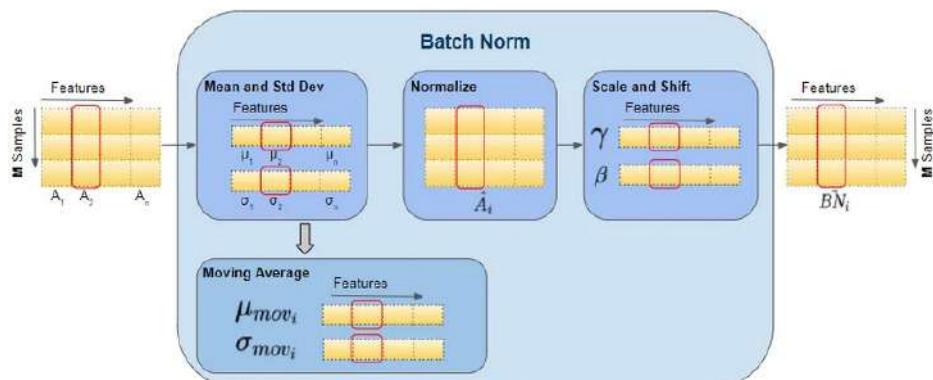
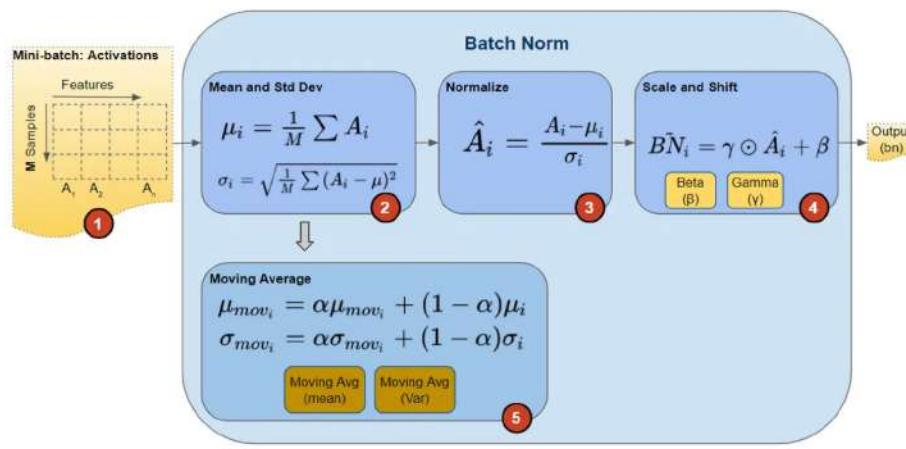
$$\gamma^k = \sqrt{\text{var}(A_{ki})} \quad \text{and} \quad \beta^k = E[A_{ki}]$$

Note - Batch Norms helps network converge faster in most cases

Ishan Modi

Batch Normalization During Training

→ Example shows one mini batch at a time

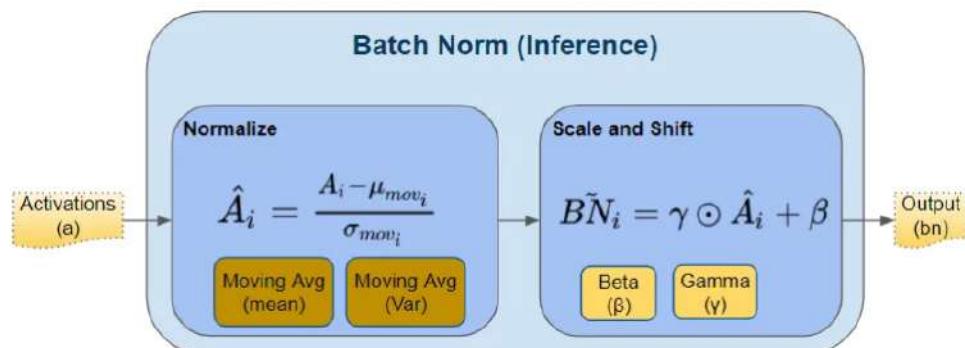


Parameters of a Batch Norm layer (Image by Author)

Ishan Modi

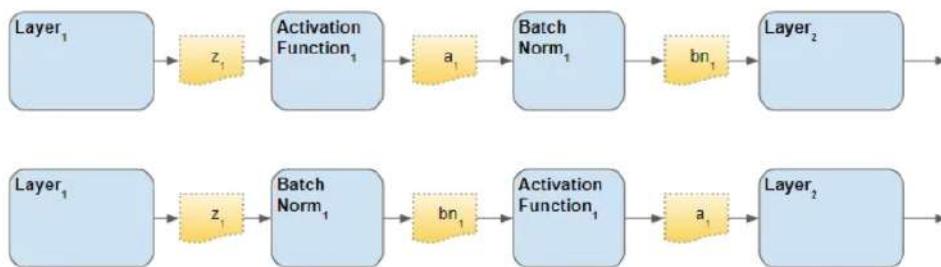
Batch Normalization During Inference

→ Example shows one mini batch at a time



Batch Norm calculation during Inference (Image by Author)

→ Finally, below is the overall functioning of batch normalization in a neural network



Batch Norm can be used before or after activation (Image by Author)

Ishan Modi

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 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

Ishan Modi

(work in (n-1) hours | cont.)
Table 7 - Co-occurrence matrix $\Rightarrow k=1$ TM9 & (cont.)

context

	Cat	eats	Dog	Bone	Mouse	TM9
Cat	0	2	0	0	0	0
eats	2	6	3	1	0	0
Dog	0	0	0	0	0	0
Bone	0	0	0	0	0	0
Mouse	0	0	0	0	0	0

f1

• Here k is not for a current 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

- 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

Chand Modi for marks less than 50%
assisted minitab analysis

- Use $\text{PMI } l = \begin{cases} \text{high if } \text{count}(w, c) \text{ is low} \\ \text{low if } \text{count}(w, c) \text{ is high} \end{cases}$

$$\text{- PMI}(w, c) = \log \frac{p(c|w)}{p(c)} = \log \left(\frac{\text{count}(w, c)}{\text{count}(c) / N} \right)$$
$$= \log \left(\frac{\text{count}(w, c) + N}{\text{count}(c) + \text{count}(w)} \right)$$

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 $\text{PMI}_0(w, c) = 0$, or $\pm \infty$ otherwise

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

* Using SVD for dimensionality reduction

$$\hat{X}_{m \times n} = \underbrace{(u_{m \times k} \sigma_{k \times k} v_{k \times n}^T)}_{\text{approximate}} \text{ where } u \text{ and } v \text{ are orthogonal matrices}$$

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

RANK k approx
we construct it using k matrix

Now (u, σ, v^T) rows = x_i

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

if we ignore denominator

Ishan Modi

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

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

$\therefore X \hat{X}^T < \hat{X} \hat{X}^T$ } similarity between various

+ terms, so elements of $\hat{X} \hat{X}^T$ increases by it is able to learn relation

(Now, check) for $\hat{X} \hat{X}^T$ we have

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

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

(1) $\text{base: } U\Sigma V^T \cdot V\Sigma^T U^T$

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

$(U\Sigma)^T = \text{shape } U\Sigma \text{ on } (U\Sigma)^T$ belongs to $\mathbb{R}^{k \times k}$

$m \times n \xrightarrow{n \times m} m \times k \text{ & } k \times m$ dimension formula $k \times m$ and $n \times k$

if let $U\Sigma V^T$ is W_{word} } $\xrightarrow{\text{dimension}} \text{shape of } W_{word} = m \times k$

$X \hat{X}^T = W_{word} W_{word}^T$ both have same cosine similarity (dot product)

both have same cosine similarity (dot product) ignoring denominator

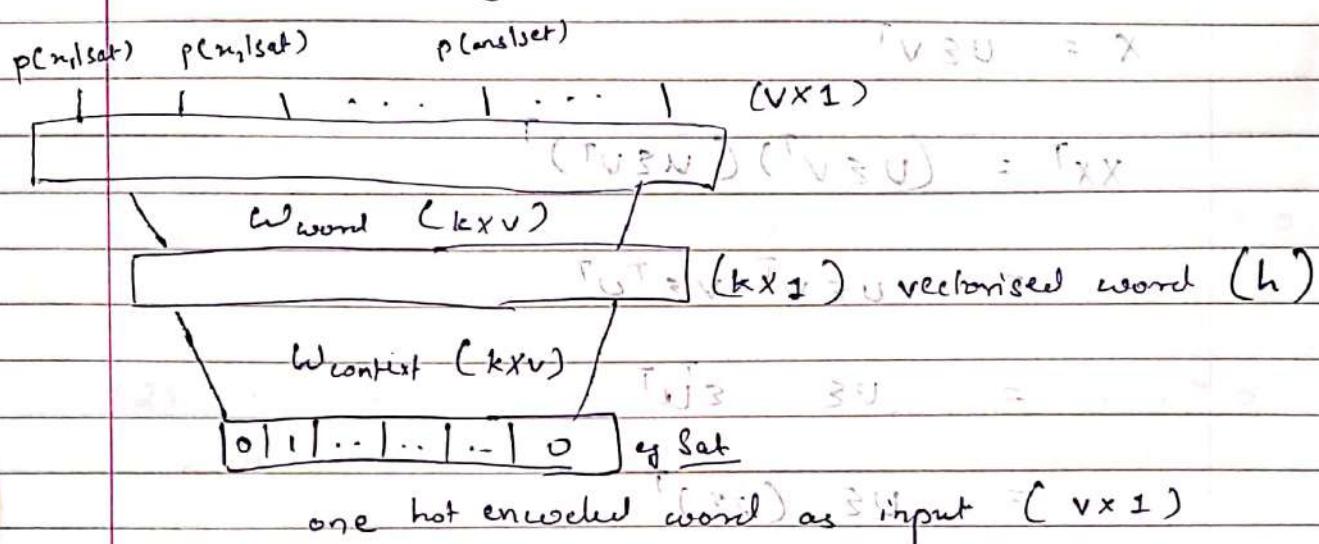
though dimension of W_{word} is very less than X

Ishan Modi

What happens is you add all the softmax outputs for words.

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

X Working

Ishan Modi

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

(Word context w is $x = [x_i \dots x_{i+h}]$)

g-

$$\begin{matrix} & w_{\text{context}} \\ x & = h \end{matrix}$$

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

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

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

→ Now,

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

lets see any one entry from output

$$p(\text{Con} | \text{sat}) = \frac{e^{(W_{\text{word}} \times h)(i)}}{\sum_j e^{(W_{\text{word}} \times h)(j)}}$$

Word

$$W_{\text{word}}^T \times h = \text{o/p}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$

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

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

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

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



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

$$h = \mathbf{W}_{\text{context}} \cdot \mathbf{u}_c = u_c$$

now i^{th} row is for word w_i

$$\hat{y}_w = \frac{\exp(u_c \cdot v_w)}{\sum_{w' \in V} \exp(u_c \cdot v_{w'})}$$

$$\exp(u_c \cdot v_w) = (\mathbf{u}_c \cdot \mathbf{v}_w)$$

Now we want to learn v_w

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

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

$$= -(\mathbf{u}_c \cdot \mathbf{v}_w - \log \sum_{w' \in V} \exp(u_c \cdot v_{w'}))$$

$$\nabla_{v_w} = \frac{\partial}{\partial} \left(\mathbf{u}_c \cdot \left[\frac{\exp(u_c \cdot v_w)}{\sum_{w' \in V} \exp(u_c \cdot v_{w'})} \cdot \mathbf{u}_c \right] \right)$$

loss funct

wrt v_w

$$\text{Loss} = -\mathbf{u}_c (1 - \hat{y}_w)$$

Update rule:

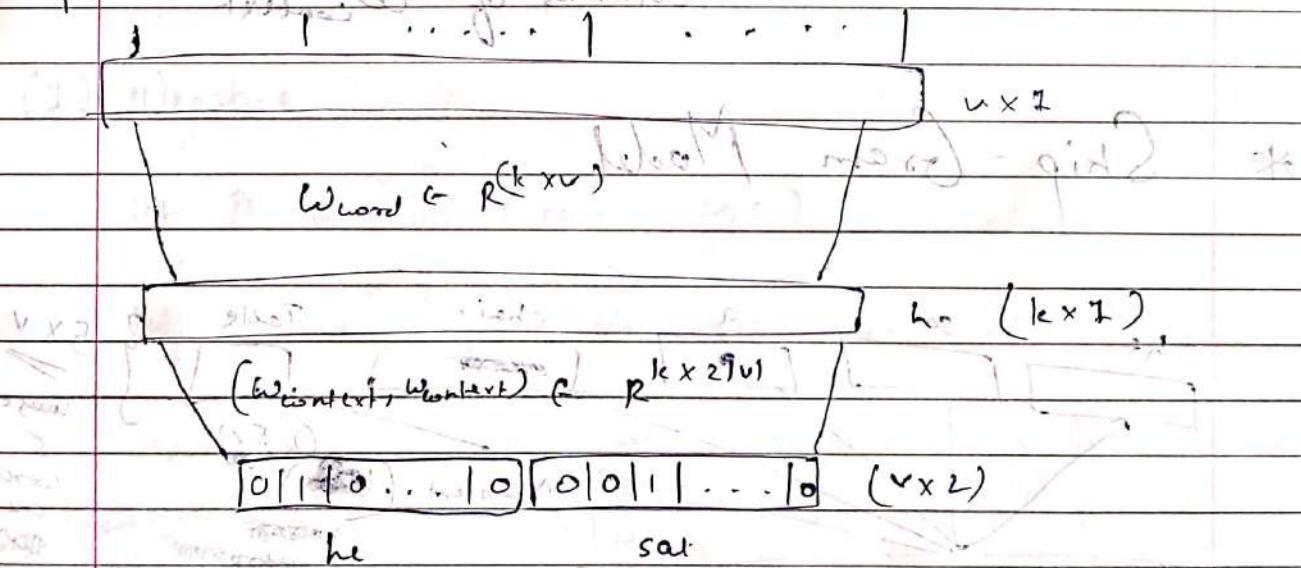
$$\Delta V_w = V_w - \eta \nabla V_w$$

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

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

→ More than one word are given

multiple storage & querying present
 written in two sets of words
 (plurals, past) & (singular, present)
 based on context



Now,

number of words ($k \times 2$)

$$d-1 \text{ columns}$$

$$h = \sum_{j=1}^{d-1} u_{c_j} \rightarrow \text{sum of columns in } w_{context}$$

corresponding to one lot of words given

Ishan Modi

$$\begin{array}{c}
 A+B \\
 \begin{array}{ccccccc}
 -1 & \{0.5\} & 2 & -1 & 0.5 & \{2\} & 0 \\
 3 & -1 & -2 & 3 & -1 & -2 & 1 \\
 -2 & \{1.7\} & 3 & -2 & 1.7 & \{3\} & 0
 \end{array}
 \end{array}
 \quad \text{using backprop}$$

$w_{context}$ w_{target}

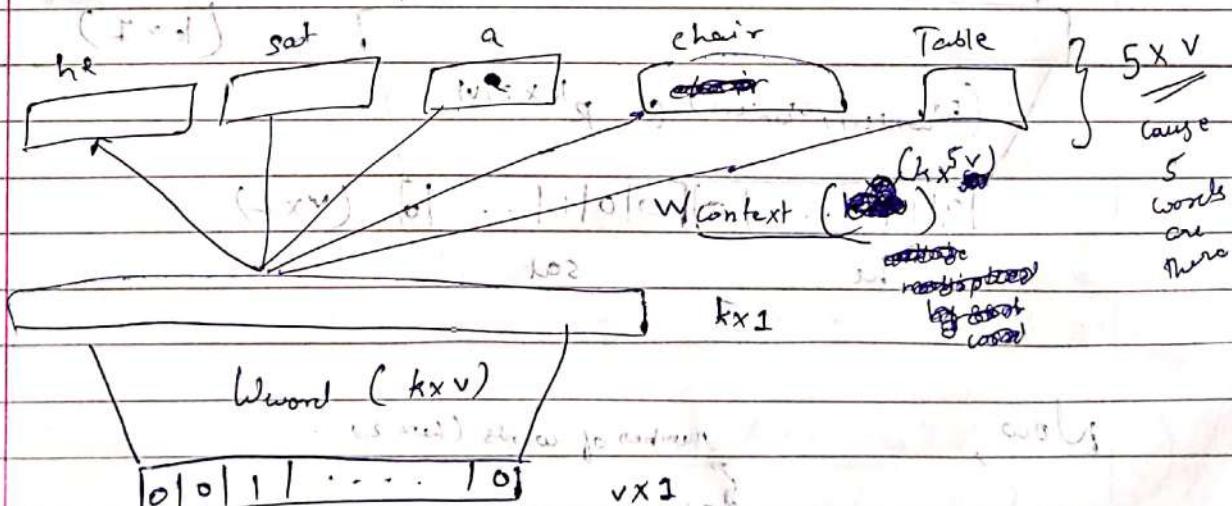
$$\left[\begin{array}{c} 0 \\ 1 \\ 0 \\ 0 \end{array} \right] \left\{ \begin{array}{c} sat \\ he \end{array} \right\}$$

$$= \begin{bmatrix} 2.5 \\ -3 \\ 4.7 \end{bmatrix} = \underline{\underline{A+B}}$$

Now,

During Backprop, update all values of w_{target} & only the participating columns of $w_{context}$

* Skip-Gram Model.



Ishan Modi

→ 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 randomly, 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)

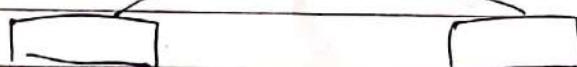
let D' be random incorrect pair (w, r)

$$u_c = c, v_w = w, u_r = r$$

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

\uparrow sigmoid \uparrow

\uparrow dot product \uparrow



Ishan Modi

Now a red arrow between training of two words before soft \rightarrow
 $p(z=1|w, c) = \sigma(u_c^T v_w)$ we have been
 doing this in back propagation

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

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

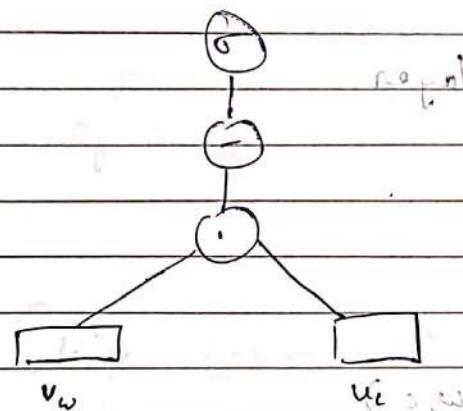
between thoughts to know \rightarrow instances & also new word if
 maximize $\prod p(z=1|w, c)$ for all words

is maximum Θ over $(w, c) \in D$ exist and choose Θ and w

and No. of mistakes

~~incorrect context words~~

Now,



$$\text{normal form } p(z=0|w, c) = 1 - \sigma(u_c^T v_w)$$

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

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

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

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

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

$$\therefore \text{maximize } \prod_{(w, c) \in D} p(z=0|w, c)$$

Ishan Modi

Total

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

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

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

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

bringing u_c & v_ω moving u_r & v_ω away

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

→ Selection of (ω, r)

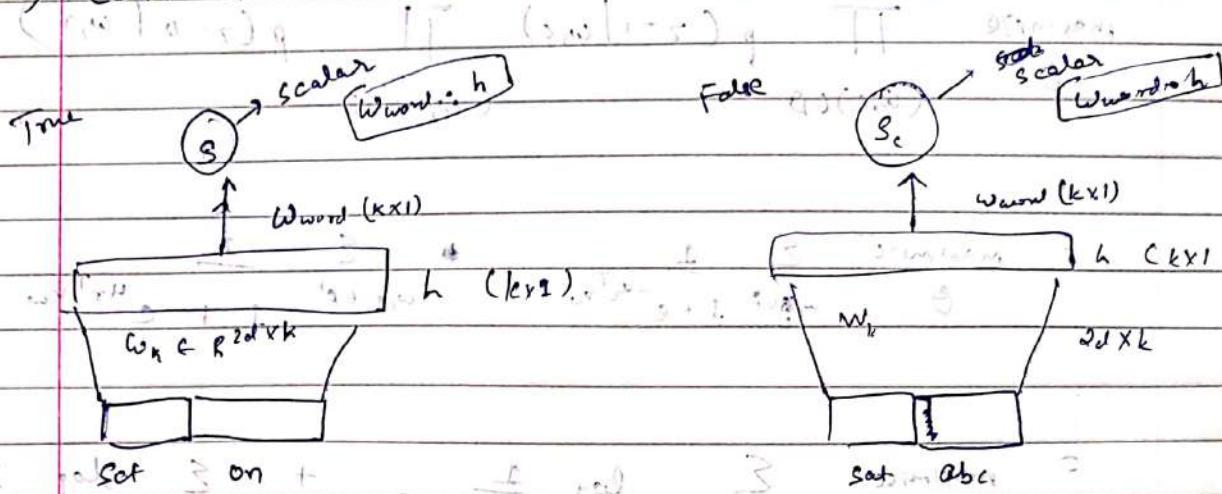
The random context word is drawn from a modified unigram distribution

$$\text{probability of } r = \frac{p(r)^{\frac{3}{4}} \cdot \text{count}(r)}{\sum_{r'} p(r')^{\frac{3}{4}} \cdot \text{count}(r')}$$

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

Ishan Modi

(2) Contrastive Estimation



We want to maximize

(large difference)

$s - s_c + \text{margin}$ (want to prevent margin between them)

Thus we maximize

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

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

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

Here

differentiate margin again

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

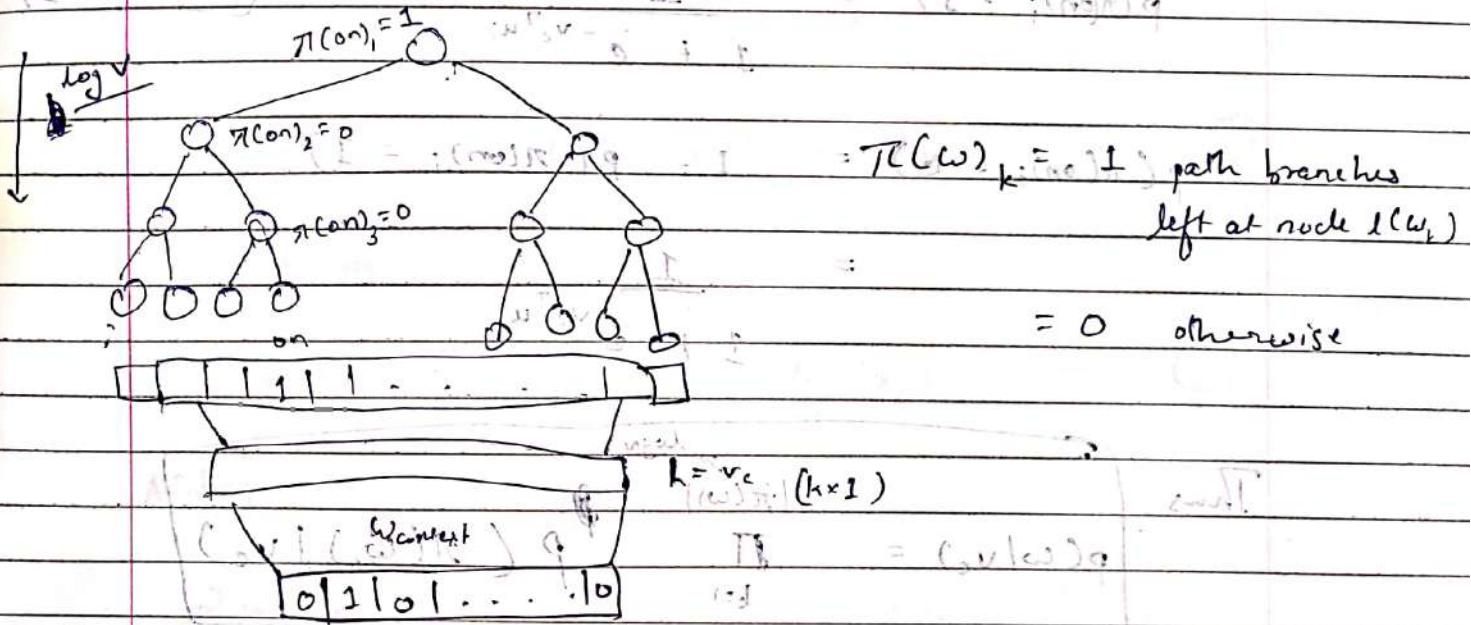
if $s - (s_c + m) < 0$, loss = $(s - (s_c + m))$

maximize

Ishan Modi

Prashant Jaiswal
Prashant Jaiswal is a software engineer

(3) Hierarchical Softmax



number of nodes to be evaluated = V
 number of leaf nodes = $V - k$

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

$$u_i = V_i - 1$$

↓
 no. of leaf nodes

Thus the parameters if this network will be

ω_{context} & u_1, u_2, \dots, u_{V-1}

vector dimension of ~~all the nodes in tree except leaf nodes~~
 every non-leaf ~~node~~ node has a vector associated with it

→ Now,

$$\underline{i} \underline{x} = \underline{i} \underline{x} = (i \underline{f} \underline{g})$$

$$p(\text{con}|v_c) = \prod_k \pi(\omega_k | v_c)$$

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

Ishan Modi

Here

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

 Vectors u_i 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,

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

 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

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

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

$$(x_{ij} \cdot x_{ij}) = (x_{ij})^2 = (x_{ij})_q \cdot (x_{ij})_q$$

Ishan Modi

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 $p(c_{ilj})$

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

Adding ① & ②

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

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

\downarrow \downarrow
 a b \downarrow \downarrow
learnable parameters b value

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

learnable parameters known

→ loss function

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

Ashish Modi

* Evaluating Word Representations

→ Semantic Relatedness

Comparing computer predictions with human perception

→ Synonym Detection

Term - derived

Candidates - S imposed, believed, requested, correlated

Synonym = $\arg \max_{v \in C} \cosine(v, v_{\text{term}})$

→ Analogy

- Semantic Analogy

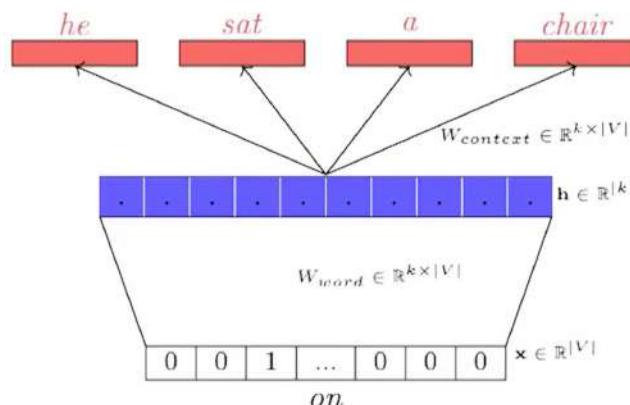
$v_{\text{sister}} - v_{\text{brother}} + v_{\text{grandson}} = v_{\text{granddaughter}}$

- Syntactic Analogy

$v_{\text{works}} = v_{\text{light}} + v_{\text{spade}} + \dots + v_{\text{species}}$

* Relation Between SVD & Word2Vec

→



- Recall that SVD does a matrix factorization of the co-occurrence matrix
- Levy et.al [2015] show that word2vec also implicitly does a matrix factorization
- What does this mean ?
- Recall that word2vec gives us W_{context} & W_{word} .
- Turns out that we can also show that $M = W_{\text{context}} * W_{\text{word}}$

where

$$M_{ij} = PMI(w_i, c_j) - \log(k)$$

k = number of negative samples

- So essentially, word2vec gives matrix M which is related to the co-occurrence matrix (ve SVD does)

Ishan Modi

X.T