

ML

Date _____
Page _____

linear Regression:

$$y = mx + c$$

$$\text{price} = m \times \text{area} + c$$

Model finds m and c by minimizing $\sum_{i=1}^n (y_i - \hat{y}_i)^2$

(best fit line)

Given a house with some area, predict the price.

* Multivariate linear regression: (multiple variables)

$$y = m_1x_1 + m_2x_2 + m_3x_3 + b$$

$$\text{price} = m_1 \times \text{area} + m_2 \times \text{bedrooms} + m_3 \times \text{age} + b$$

dependent variable
independent variables
(features)

m_1, m_2, m_3 — coefficients

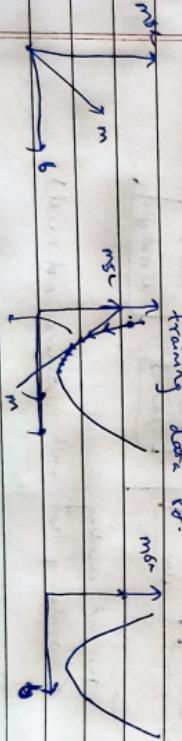
b — intercept

* Mean squared error:

$$\text{mse} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Cost function

* Gradient descent :: to find best fit line for given training data $\{d\}$.



24

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - (m\pi_i + b))^2$$

$$\frac{\partial}{\partial m} = \frac{2}{n} \sum_{i=1}^n -m_i (y_i - (m\pi_i + b))$$

$$\frac{\partial b}{\partial b} = \frac{2}{n} \sum_{i=1}^n - (y_i - (m\pi_i + b))$$

$$m' = m - \text{learning rate} \times \frac{\partial}{\partial m}$$

For example, start with $m=0$, $b=0$, learning rate = 0.001

After 1000 iterations we get the desired m and b values
(but just j/k)

→

After training, save the model using `pickel` or
`pickle`. (No need to train again!)

* One hot encoding:

→ Categorical variables:

Nominal	Ordinal
male	green
female	red
	blue

high	soft
medium	medium
low	dissatisfied

One hot encoding
(representing text with integers not labels).

obj:	town	area	price	A	B	C
A	-	-	-	1	0	0
A	-	-	-	0	1	0
B	-	-	-	0	0	1
B	-	-	-	0	1	0
C	-	-	-	0	0	1

↳ Now drop this column

→ Then drop one of the dummy variable columns (either A or C)

e.g.: if we drop column C,
as the dropped column can be
divided from others

the other columns

$y = \text{price}$

→ train-test-split : e.g.: training data: 80%, testing data: 20%

$X_{\text{train}}, X_{\text{test}}, y_{\text{train}}, y_{\text{test}} = \text{train-test-split}(X, y, \text{test_size}=0.2)$
(randomly select rows)

Logistic Regression :

→ Predicted values is categorical. e.g.: Yes or No

→ Classification problem

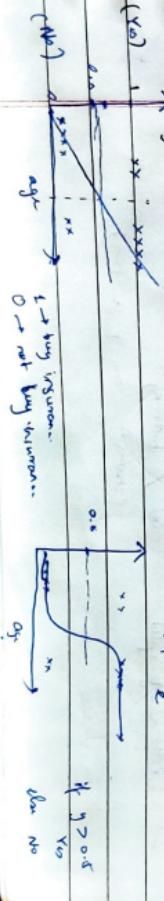
Binary classification

Multiclass classification

Q: Yes, No
e.g.: Dimensional, qualitative, independent

$$y \leftarrow \delta'(z) = \frac{1}{1 + e^{-z}} \quad y \in [0, 1]$$

$$y = mx + b \quad y = \frac{1}{1 + e^{-(mx+b)}}$$



Naive Bayes Classifier:

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(A)}$$

Naive because of assumption that features are independent of each other

Type: Bernoulli, Multinomial, Gaussian
Features are: binary - discrete, continuous

Hyper-parameter Tuning:

→ Use k-fold cross validation score for different set of parameters & select the params with the best score.

library: (i) GridSearchCV

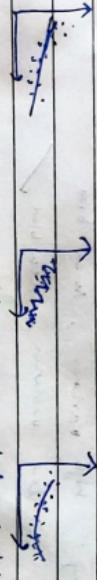
(does the same thing in one line instead of multiple for loops)

(ii) RandomizedSearchCV :

No. of permutation of parameter can be very high.
selecting a fixed no. of random (parameters)
and then select the one with best score

L1 & L2 regularization:

Lasso Ridge regularization



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

Penalty: make θ_3 and θ_4 close to 0

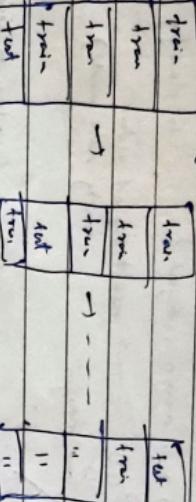
$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - h(\theta, x_i))^2 + \lambda \sum_{i=1}^n \theta_i$$

new term added that penalizes higher θ values

$$\text{MSE} = \frac{1}{n} + \lambda \sum_{i=1}^n |\theta_i| \rightarrow \text{L1 regularized}$$

K-fold cross-validation:

divides into k folds

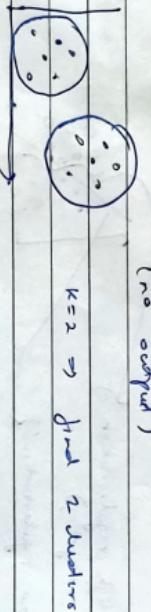


take the average of the error

→ every algo studied so far (regression, classification)

- * ML → supervised learning (we know target variable / class)
- └ unsupervised " " "
- └ reinforcement " "

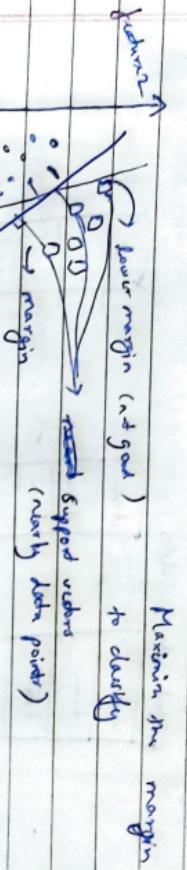
K-means clustering: (Unsupervised learning)



Algorithm:

- ① Take K random centroids. The point closest to a centroid
- ② become part of that cluster.
- ③ Recalculate the centroid by finding the center
- of all the points of that cluster.
- ④ Now repeating step 1 will give new set of point belonging to that cluster.
- ⑤ Repeat the above process until no new points added to any cluster. Finally we get K diff. clusters
- ⑥ To find \sum (sum of squared error) vs K
- ↓ and the above point is desired K

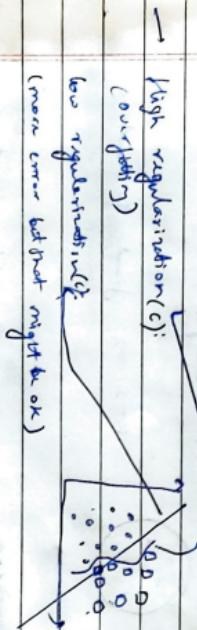
Support Vector Machine :



feature → 2D → boundary is line
3D → boundary is plane.

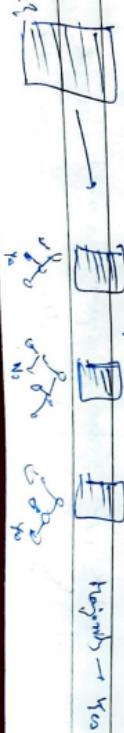
Support vector machine draws a hyperplane in ' n ' dimensional space : maximizes the margin between classification groups.

→ High gamma : consider only the nearby points for margin
low gamma : " all the points



Random forest :

- randomly sample the training data into multiple training sets & make a decision tree of each set.
- Then to predict, take the output from each of the decision trees & select the one with majority vote



Date _____
Page _____

→ Multiclass example : X : image of digit (represented in 2D)

y : 0, 1, 2, 3, ..., 9

→ Confusion matrix:

$cm = \text{confusion matrix} (y_{\text{true}}, y_{\text{predicted}})$



Out of all the test data, 37 times the model predicted the image as 0 when the image was actually 0.

Decision Predicted (y-predicted) | 0 times, the model predicted the image as 0 when the image was actually 3.
Correct predictions | 3

Decision Tree :

eg. Feature: Computer

Salary > 100k \$?

Computer program

Business program

Facebook

Yes

Bachelor Master

No

How to select ordering of features?

y: on basis of concern
Facebook: 6 Yes 0 No 6/0 (low entropy)

Worst: 3 Yes, 4 No

High Information gain

→ Criteria of split : Min / Entropy

(choose one of the two)

K nearest neighbour Classification :

↓
feature → Nearest

↓
... → ? g: k=3

↓
↓ → Nearest → find nearest k points (minimum distance)
↓ → Nearest → find the class having the maximum kth point near to the test point.

→ find K → trial & error (g (threshold) (typically = 5))

PCA: Principal Component Analysis:

→ dimensionality reduction technique

(~~eg~~ select the features having the maximum impact)

(eg: 100 features to 6 features (PC_1, PC_2, \dots, PC_6))

→ the features are completely new (not a subset of the columns)

* Bias: how accurately a model can capture a pattern in training data, high train accuracy → high bias

Variance: high test error → high variance

Oversampling: high variance, high bias

Balanced: low variance, low bias ✓

Basis → set of vectors in \mathbb{R}^n that are LI and can represent any point in \mathbb{R}^n

Orthonormal basis: rows are mutually orthogonal
columns " "

$$A^T A = A A^T = I \Rightarrow A^T = A^{-1}$$

$$g: A = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

Eigen decomposition:

$$Av = \lambda v$$

(vector v is scaled)

Spanning eigen values called eigenvalues

g: $A = \begin{bmatrix} 1.25 & 0.75 \\ 0.75 & 1.25 \end{bmatrix}$

(symmetric)

\Rightarrow

(diagonal)

$$\begin{bmatrix} 1.25 & 0.75 \\ 0.75 & 1.25 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \lambda \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

we get x_1, x_2

$$1.25x_1 + 0.75x_2 = \lambda x_1$$

$$0.75x_1 + 1.25x_2 = \lambda x_2$$

$$x_1 + x_2 = 1$$

we get x_1, x_2 for λ

$$v_1 = \begin{bmatrix} 0.707 \\ 0.707 \end{bmatrix}, \quad v_2 = \begin{bmatrix} 0.707 \\ -0.707 \end{bmatrix}, \quad \lambda_1 = 0.5 \quad (\text{scaled row})$$

$$v_2 = \begin{bmatrix} 0.707 \\ -0.707 \end{bmatrix}, \quad \lambda_2 = 0.5 \quad (\text{scaled row})$$

Exercise: $\|x\|_1 = 1$ find Ax

We know there such x : $x = v_1 + v_2$

$x = (0.707, 0.707) \rightarrow (1.414, 1.414)$



Eigenvectors

$$Ax = \begin{bmatrix} 1.414 \\ 1.414 \end{bmatrix}$$

$$Ax = \begin{bmatrix} 1.414 \\ 1.414 \end{bmatrix}$$

$$Ax = \begin{bmatrix} 1.414 \\ 1.414 \end{bmatrix}$$

→ set of vectors is linearly independent if no vector in the set is a linear combination of other vectors

Date _____
Page _____

→ If column space is $\mathbb{R}^m \Rightarrow$ exactly m linearly independent columns

→ Square matrix with $1|1$ columns is called 'singular'.

Norms:

→ intuitive meaning → distance of x from origin

$$l' = \|x\|_p = \left(\sum_i |x_i|^p \right)^{\frac{1}{p}} \quad p \in \mathbb{R}, p \geq 1$$

$$l^2 = (\sum x_i^2)^{\frac{1}{2}} = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

(analogous to magnitude of the vector)

$$\text{Square } l^2_{\text{norm}} = x^T x$$

Ex: $\begin{pmatrix} x_1 & x_2 & \dots & x_n \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = x_1^2 + x_2^2 + \dots + x_n^2$

$$\rightarrow l^1 = \|x\| = \sum_i |x_i| \quad (\text{mod is used})$$

$$\rightarrow \text{Frobenius norm} \quad \|A\|_F = \sqrt{\sum_{i,j} A_{ij}^2} \quad (\text{for matrix, } \|A\|_{\text{norm}})$$

$$\rightarrow \text{Diagonal matrix: } D_{ij} = 0 \text{ where } i \neq j$$

Distance matrix: $A_{ij} = D_{ij}$

$$\rightarrow \text{Unit vector: } \|v\|_2 = 1$$

→ In \mathbb{R}^n , atmost n vectors can be mutually orthogonal (90°) with non-zero norms.

g_i for \mathbb{R}^3 : $\vec{i}, \vec{j}, \vec{k}$ are mutually orthogonal

→ Vectors orthogonal & having unit norm are orthonormal

$$\frac{a_1}{a_2} = \frac{b_1}{b_2} = \frac{c_1}{c_2} \rightarrow \text{Lyrical}$$

$\neq \frac{c_1}{c_2} \rightarrow \text{no soln}$

→ Exactly one soln $\frac{a_1}{a_2} \neq \frac{b_1}{b_2}$

⇒ Span is entire 2D plane

(linear combination of the column vectors spans the entire plane)

$$A^n = b \Rightarrow n = A^{-1}b$$

→ Matrix non-invertible if

not a square matrix or columns are linearly dependent
(have rank)

$A^n \rightarrow$ linear comb' of column vectors of A

$A^n = b \rightarrow$ test whether b is in the span of A
or column space

→ For 3D: If all 3 vectors in same plane

i.e. they are linearly dependent (any vector can be obtained by linear combination of other 2),
then span is a plane.

Else span is entire 3D space

→ $A^n = b$ to have solution $A \in \mathbb{R}^{m \times n}$,
column space of A must be \mathbb{R}^m .

$$A: m \times n$$

A must have atleast m columns i.e. $n \geq m$

$$g_i: A: 3 \times 2 \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} \quad a: 2 \times 1 \begin{bmatrix} \dots \\ \dots \end{bmatrix} \quad b: 3 \times 1 \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix}$$

$A \in \mathbb{R}^{m \times n} \rightarrow$ 2D plane is soln only when b lies on that plane.

$\therefore [m \geq n]$ is necessary to have soln if b lies on that plane.

Linear Algebra

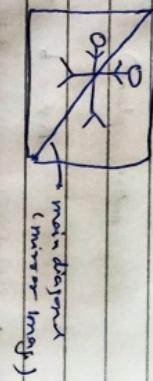


- 0D : scalar
- 1D : vector
- tensor : 0D or 1D or 2D or 3D : -

→ Matrix transpose

$$B_{ij}^T = A_{ji}$$

$$(AB)^T = B^T A^T$$



→ Matrix multiplication:

$$A \cdot B = A_{ij} \cdot B_{ji}^T$$

(dot product of vectors)

→ Covariance matrix: $A \cdot A^T$

feature

vector

$$A: \begin{bmatrix} \text{row 1} \\ \text{row 2} \\ \vdots \\ \text{row n} \end{bmatrix} \quad A^T: \begin{bmatrix} \text{feature 1} \\ \text{feature 2} \\ \vdots \\ \text{feature n} \end{bmatrix}$$

$A \cdot A^T = C_{ij}^T$: how much similar i item is to j item

$A^T \cdot A = D_{ij}$: feature i and feature j

→ System of Equations: $A \cdot x = b$

$$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} = A = \begin{bmatrix} -4 & -5 \\ -2 & 3 \\ \vdots & \vdots \\ m_1 & m_2 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad b = \begin{bmatrix} -13 \\ 9 \\ \vdots \\ g \end{bmatrix}$$

is discriminant = 0,
non span or column
option is a straight line
⇒ No soln or infinite soln.

$$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} = A = \begin{bmatrix} -4 & -5 \\ -2 & 3 \\ \vdots & \vdots \\ m_1 & m_2 \end{bmatrix} + m_2 \begin{bmatrix} -5 \\ 3 \end{bmatrix} = \begin{bmatrix} -13 \\ 9 \\ \vdots \\ g \end{bmatrix}$$

Feature Engineering

PCA → Unsupervised
Date: _____
Page: _____

linear regression:

→ $\text{f(x)} = x^{\alpha} \cdot \dots \cdot x^{\beta} \rightarrow$

$$y = \omega_1 + \omega_2 m - \omega_m$$

$$\text{Objection : to minimize } \text{MSE}_{\text{train}} \\ \nabla_{\theta} \text{ MSE}_{\text{train}} = 0 \\ \nabla_{\theta} \frac{1}{m} \sum_{i=1}^m \left[y^{(\text{train})}_i - \hat{y}^{(\text{train})}_i \right]^2 = 0$$

$$\begin{aligned} & \left[\begin{array}{cc|c} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{array} \right] \xrightarrow{\text{Row } 2 \leftrightarrow \text{Row } 3} \left[\begin{array}{cc|c} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{array} \right] \xrightarrow{\text{Row } 2 - \text{Row } 1} \left[\begin{array}{cc|c} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{array} \right] \\ & \text{The system of equations is:} \\ & x_1 = 1 \\ & x_2 = 0 \\ & \text{There is one free variable: } x_3 \end{aligned}$$

$$\omega = \begin{pmatrix} \omega_{11} & \omega_{12} \\ \omega_{21} & \omega_{22} \end{pmatrix} = \begin{pmatrix} \omega_{11} & \omega_{12} \\ \omega_{21} & \omega_{22} \end{pmatrix} X^{-1} \begin{pmatrix} \omega_{11} & \omega_{12} \\ \omega_{21} & \omega_{22} \end{pmatrix}$$

Finely
 (precise) = ~~order~~ $w_1, w_2 + \dots + w_n +$
~~bias~~ + b

one part of the system and another part of the system.

$$\tau_{\text{new}} = \tau_0 - m\alpha$$

104

Probability, Bayes theorem problems, Gaussian Distributions

$$\rho = \begin{bmatrix} \rho_1 & \rho_2 & \dots & \rho_n \end{bmatrix} \quad \rho \rightarrow \text{new basis}$$

new

The new are old rates

$$\rho_1, \rho_2, \dots, \rho_n$$

$$(\tau_1, \tau_2, \dots, \tau_n) \begin{bmatrix} \rho_1 & \rho_2 & \dots & \rho_n \end{bmatrix}$$

unnormalized

product give coordinate along ρ

$$\hat{\mathbf{A}} = \mathbf{A}\rho \quad (\hat{\mathbf{A}} \rightarrow \text{matrix of transformed rates})$$

$m \times n$

diagonal matrix

$$\frac{1}{m} \hat{\mathbf{A}}^T \hat{\mathbf{A}} = \rho^T \Sigma \rho = \mathbb{D} \quad \begin{array}{l} \text{covariance matrix} \\ \text{Let } A = m \times n \end{array}$$

covariance matrix

(cov. matrix)

$A^T A$ transformation

of original

$$C_{ij} = \frac{1}{m} \sum_{i=1}^m (x_{ij} - \mu_i)(x_{ij} - \mu_i)$$

$$\therefore \hat{\mathbf{P}} \leq \rho = \mathbb{D} \quad \text{If } \mu_i \text{ are } 0$$

& which orthogonal matrix ρ diagonalise Σ ?

$$C_{ij} = \frac{1}{m} \sum_{i=1}^m x_{ii} \cdot x_{ij}$$

Matrix P where columns

eigen vectors of $\Sigma = \hat{\mathbf{A}}^T \hat{\mathbf{A}}$

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

$$\left[\begin{array}{c} \rho_1 \\ \vdots \\ \rho_n \end{array} \right] = \mathbf{P} \left[\begin{array}{c} \mu_1 \\ \vdots \\ \mu_n \end{array} \right]$$

\mathbf{P}^{-1}

\mathbf{P}^{-1}

#

Determinant :

 $\det(A)$ — only for square mat \rightarrow

= product of all eigen value of matrix

→ Matrix idemus ?

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \quad \det(A) = A = \begin{bmatrix} A_{11} = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \\ A_{12} = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \end{bmatrix}$$

Principles Component Analysis:

$$\frac{\partial L_{\text{sum}}}{\partial x_i} = \begin{bmatrix} \frac{\partial L_{\text{sum}}}{\partial x_{i1}} & \frac{\partial L_{\text{sum}}}{\partial x_{i2}} \\ \frac{\partial L_{\text{sum}}}{\partial x_{i3}} & \frac{\partial L_{\text{sum}}}{\partial x_{i4}} \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = A$$

$$\text{new matrix with } \# \text{ rows} = \# \text{ features} \quad \# \text{ columns} = \# \text{ variables} \quad \text{from } L_{\text{sum}}$$

$$\text{row sum operation} \quad \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\begin{aligned} \text{row } 1: & \quad a_{11}x_1 + a_{12}x_2 \\ \text{row } 2: & \quad a_{21}x_1 + a_{22}x_2 \\ = & \quad a_{11}x_1 + a_{12}x_2 + a_{21}x_1 + a_{22}x_2 \\ = & \quad a_{11}x_1 + 2a_{12}x_2 + a_{21}x_1 \\ = & \quad f(a_1, a_2) \end{aligned}$$

$$\begin{aligned} \frac{\partial (x^T A x)}{\partial x_i} &= [a_1 \ a_2] \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= [a_1 \ a_2] \begin{bmatrix} a_{11}x_1 + a_{12}x_2 \\ a_{21}x_1 + a_{22}x_2 \end{bmatrix} \\ &= a_{11}x_1^2 + a_{12}x_1x_2 + a_{21}x_1x_2 + a_{22}x_2^2 \\ &= 2x_1^2 + 2x_1x_2 + 2x_2^2 \\ &= 2x^T A x \end{aligned}$$

$$\frac{\partial (x^T A x)}{\partial x_i} = 2x^T A x$$

$\text{Tr} \rightarrow$ Eigen vectors of $A \in \mathbb{R}^{n \times n}$ having distinct eigen values are \perp

$\text{Tr} \rightarrow$ Eigen values of square symmetric matrix are

$$Q = V^T V : Q_{ij} = U_i^T U_j = 0 \text{ if } i \neq j = 1 \text{ if } i=j$$

$$\Rightarrow V^T V = I \Rightarrow V^{-1} = V^T$$

Matrix is said singular if any of eigen value = 0

$$f(\lambda) = \lambda^T A \lambda \quad \text{when} \quad \|A\|_2 = 1$$

When α is eigen vector of A , it takes corresponding eigen value

$$\alpha^T A \alpha = \lambda^T \alpha \Rightarrow \alpha \lambda^T \alpha = \lambda (\sum \alpha_i) = \lambda$$

\rightarrow Eigen decomposition

but A has n independent eigen vectors $\{\alpha_1, \alpha_2, \dots, \alpha_n\}$ with corresp. eigen val. $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$

$$\text{Let } \mathbf{V} = [\alpha_1, \alpha_2, \dots, \alpha_n]$$

$\lambda = \sqrt{\lambda_1, \lambda_2, \dots, \lambda_n}$ eigen value decom.

$$A \mathbf{V} = \mathbf{V} \text{diag}(\lambda) \Rightarrow A = \mathbf{V} \text{diag}(\lambda) \mathbf{V}^{-1}$$

(λ is eigen-symmetric, then $B = V \text{diag}(m) V^{-1}$)

$$\text{eg: } \begin{bmatrix} 1.2 & 0.45 \\ 0.45 & 1.15 \end{bmatrix} \begin{bmatrix} 1.014 & -0.359 \\ 1.414 & 0.282 \end{bmatrix} = \begin{bmatrix} 1.014 & -0.359 \\ 1.414 & 0.282 \end{bmatrix}$$

$$(A) \quad (V) \quad (U)$$

$\text{diag}(A) = \sqrt{|A|} \mathbf{V}$ diag of A

(diag(λ))

Trace Operator :

$$\rightarrow \text{Tr}(A) = \sum_{i,j} A_{ii} \quad (\text{sum of main diagonal elements})$$

$$\rightarrow \|A\|_F^2 = \sum_{i,j} (A_{ij})^2 = \sqrt{\text{Tr}(A^T A)} \quad \begin{bmatrix} \cdots & \cdots \\ \cdots & \cdots \end{bmatrix} = \boxed{\text{---}}$$

$$\rightarrow \text{Tr}(B) = \text{Tr}(B^T)$$

$$\text{Tr}(AB) = \sum_{i,j} A_{ij} B_{ji}$$

$$\rightarrow \text{Tr}(ABC) = \text{Tr}(BAC) = \text{Tr}(CAB) = \text{Tr}(BCA)$$

4

Training error, Generalization error
(1992 on training set) (see - 1)

(1979 on training set) (error on testing set) (mean std.)

Overfitting → overtraining error, high generalization error

Capacity:

\rightarrow high capacity \Rightarrow overfitting (eg: $a_n + b \rightarrow a_n^{-1} + b$)

U.S. Senate

General Naldehra : Vajreshwar

Non-parametric methods (no parameter vector)

For a given x , $\hat{y} = y_i$ where $i = \arg\min \|x_i - x\|_2$

(exclusion due to y_i)

→ مکانیزم

24

Logistic Regression :

$$\text{Probability} = \frac{1 - e^{-(\mu + \sigma^2)}}{1 - p(x)}$$

$$P(\omega_1, \omega_2) = \prod_{i=1}^n P(\omega_i)$$

$$= \prod_{i=1}^n p(x_i)^{y_i} \cdot \left(1 - p(x_i)\right)^{1-y_i} + \sum_{j=1}^{n-1} \log(p(x_i))$$

NPTEL:

Sequence \rightarrow RNN, LSTM

Image \rightarrow CNN

$$\vdots \rightarrow j \rightarrow g$$

Date _____
Page _____

McCulloch Pitts Neuron:

$$y = 1 \text{ if } \sum_{i=0}^n w_i n_i \geq \theta \quad \text{otherwise}$$

$$\Rightarrow y = \begin{cases} 1 & \sum_{i=0}^n w_i n_i \geq \theta \\ 0 & \text{otherwise} \end{cases}$$

Perceptron:

$$y = 1 \text{ if } \sum_{i=0}^n w_i n_i \geq \theta \quad \text{otherwise}$$

$$= 0 \text{ if } \sum_{i=0}^n w_i n_i < \theta$$

AND gate:

$$y = 1 \text{ if } \sum_{i=0}^n w_i n_i \geq \theta \quad \text{otherwise}$$

$$= 0 \text{ if } \sum_{i=0}^n w_i n_i < \theta$$

\rightarrow red colored inputs allowed

\rightarrow threshold θ can be learned

Single perceptron: can not classify if the function is not linearly separable.

Perceptron learning algo:

$$\text{Perceptron learning algorithm: } w^{t+1} = w^t + \eta d_n$$

where d_n is the error term

$$w_0 + w_1 n_1 - 1 = 0$$

(one function)

$$w_0 + w_1 n_2 - 1 = 0$$

(one function)

$$w_0 + w_1 n_3 - 1 = 0$$

(one function)

$$w_0 + w_1 n_4 - 1 = 0$$

(one function)

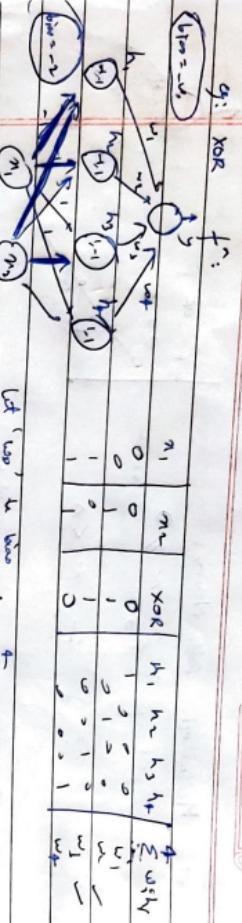
Training: Not linearly separable e.g. XOR

$$x_1 \oplus x_2 = y$$

\rightarrow Multilayered network of perceptrons (MLP) can represent

any function f (may be linearly inseparable)

* Any function f of n inputs can be represented
linearly by a network of perceptrons containing 1 hidden layer with
2 perceptions in 1 output layer with 1 perceptron



Data
Page

四

وَيُكَلِّفُهُمْ أَنْ يَرْجِعُوا إِلَى الْأَنْوَارِ (الْأَنْوَارُ الْمُبَشِّرَةُ)

smooth, continuous, differentiation

Input: x^i → real values $x \in \mathbb{R}^n$
 Output: y^i → real numbers 0 and 1
 We are aiming for
 (vector quantization for
 real valued \mathbf{x}^i)

1) Data : $\{x_i, y_i\}_{i=1}^n$ Compounds :

$$y = \frac{1}{1 + e^{-(w^T x)}}$$

Learning algo : algo to learn w (e.g. prediction learning algo)
↳ gradient descent

Learning alg. should minimize the loss function:

卷之三

→ At the output layer

$$f(a) = h_L(a) = O(a_L(m))$$

L output activation fn

(e.g. softmax, linear, etc.)

Parameters : (for all supervised learning algos),

$$\text{i)} \quad \text{Data} : \{x_i, y_i\}_{i=1}^N \quad (\text{N training examples})$$

$$\text{ii)} \quad \text{Feature Model} : \hat{y}_i = f(x_i) \quad \hat{y}_i = O(w^3 g(w^2 g(w^1 x_i + b_1) + b_2) + b_3)$$

$$\text{iii)} \quad \text{Parameter} : \theta = w_1, \dots, w_3, b_1, b_2, \dots, b_3 \quad (L=3)$$

iv) Algorithm: Gradient descent with back propagation

v) Objective / loss function:

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^k (\hat{y}_{ij} - y_{ij})^2$$

vi) minimize $L(\theta)$.

$$\text{vii)} \quad \theta = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ b_1 \\ b_2 \\ b_3 \end{bmatrix} \quad \nabla \theta = \begin{bmatrix} \frac{\partial L(\theta)}{\partial w_1} \\ \vdots \\ \frac{\partial L(\theta)}{\partial b_3} \end{bmatrix}$$

$$\theta_{t+1} = \theta_t - \eta \in \nabla \theta_t$$

viii) \rightarrow backpropagation derivatives w.r.t all weights

$$w_j : \frac{\partial L(\theta)}{\partial w_{j1}} \rightarrow \frac{\partial L(\theta)}{\partial w_{j2}} \rightarrow \dots \rightarrow \frac{\partial L(\theta)}{\partial w_{jn}},$$

$$\text{Total derivatives} = \underbrace{(L-1) \times (n \times n)}_{(L-1) \times n} + \underbrace{1 \times (n \times k)}_{1 \times (k \times 1)}$$

$$y \in \delta^{(m)}$$

$$f(x) = \frac{1}{1 + e^{-(wx+b)}}$$

$$l = \frac{1}{2} (f(x) - y)^2$$

$$\nabla_w = \frac{1}{2} \cdot 2 (f(x) - y) \frac{\partial l}{\partial w}$$

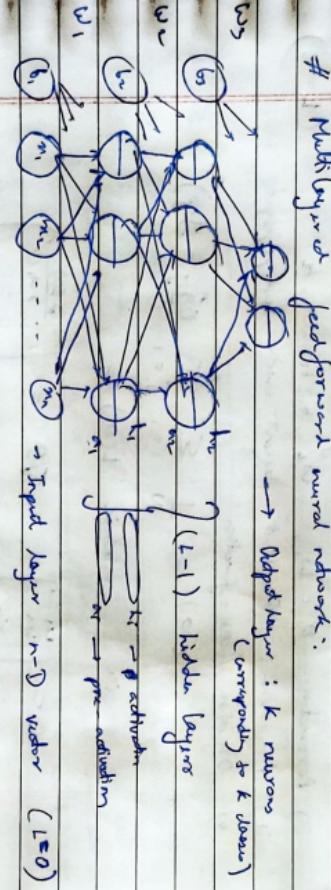
$$\nabla_b = (f(x) - y) \frac{\partial l}{\partial b} \cdot (1 - f(x)) \cdot x$$

#

Gradient descent:

$$\begin{aligned} n' &= n - \epsilon \underbrace{\nabla_x f(x)}_{\text{numerical gradient}} \\ &\quad \leftarrow \text{gradient update} \quad \leftarrow (n') < \delta^{(n)} \end{aligned}$$

$$u: \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix} = \begin{bmatrix} n_1 \\ n_2 \\ \vdots \\ n_n \end{bmatrix} - \epsilon \begin{bmatrix} \frac{\partial l}{\partial n_1} \\ \frac{\partial l}{\partial n_2} \\ \vdots \\ \frac{\partial l}{\partial n_n} \end{bmatrix}$$



→ Hidden layers & output layers:
pre-activation \Rightarrow aggregation (Σ)

Activation: non-linearity

\rightarrow $w_i \in R^{n \times n}$ ($i < 0 < L$) & $b_i \in R^n$

$w_L \in R^{n \times k}$ and $b_L \in R^k$ (for output layer)

$\rightarrow h_L = \hat{y} = f(x)$

$$a_i(n) = \sum_j w_{ij} h_{j-1}(n) \quad \left| \quad h_i(n) = g(a_i(n)) \right.$$

\downarrow $\sum_j w_{ij}$ \downarrow $g(a_i(n))$
 $\left| \quad \text{addition of } j^{\text{th}} \text{ neuron}$
every element of vector

Simple FFN with hidden layer:

$$\begin{aligned}
 w &= \begin{bmatrix} w_1 & w_{12} \\ w_{21} & w_{22} \end{bmatrix} & m &= \begin{bmatrix} m_1 \\ m_2 \end{bmatrix} \\
 &\quad \text{(2x2)} & h_1 &= \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} \\
 && \text{# neurons in hidden layer} & \\
 && \text{optimal weight} & \\
 & w^T m + c & = \begin{bmatrix} w_{11} & w_{12} & \begin{bmatrix} m_1 \\ m_2 \end{bmatrix} \\ w_{21} & w_{22} & \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} \end{bmatrix} \\
 &\quad \text{# neurons in hidden layer} & & \text{bias} \\
 h &= g(w^T m + c) & &
 \end{aligned}$$

activation function (sigmoid, tanh, ReLU)

$$\therefore \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = g \left(\begin{bmatrix} w_{11}m_1 + w_{12}m_2 + c_1 \\ w_{21}m_1 + w_{22}m_2 + c_2 \end{bmatrix} \right)$$

$$f \rightarrow \text{let } g(z) = \max\{0, z\} \quad (\text{ReLU activation fn})$$

$$\begin{aligned}
 \text{Now, } y &= \begin{bmatrix} w_1 & w_2 \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} + b \\
 &\Rightarrow y = w^T \max\{0, w^T m + c\} + b
 \end{aligned}$$

Example: Solution for xor problem:

$$w = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}, \quad c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \quad w = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \rightarrow b = 0$$

$$\text{let } \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$w^T m = \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \quad w^T m + c = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \quad g(w^T m + c) = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$y = [1 - 2] \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 0 = 0 + 0 = 0 \quad (\because 1 \wedge 1 = 0)$$

Full forward neural network example:

$$\text{xor } f^* \quad x = \{[1, 0], [1, 1], [0, 1], [0, 0]\}$$

Target function : $y = f^*(x)$
(actual)

Model : $y = f(x; \theta)$

Let cost function = MSE.

$$J(\theta) = \frac{1}{4} \sum_{\text{nex}} (f^*(m) - f(m; \theta))^2$$

Let us consider a linear model

$$y^2 \quad f(m; \theta) = f(m; w, b) = w^T u + b$$

$$J(\theta) = \frac{1}{4} \left[\left(0 - (w_1 u_1 + w_2 u_2 + b) \right)^2 + \left(1 - (w_1 u_1 + w_2 u_2 + b) \right)^2 + \left(1 - (w_1 u_1 + w_2 u_2 + b) \right)^2 + \left(0 - (w_1 u_1 + w_2 u_2 + b) \right)^2 \right]$$

$$\begin{aligned} \frac{\partial J}{\partial w_1} &= \frac{1}{4} \left[-2(u_1 - (w_1 u_1 + w_2 u_2 + b)) + (-u_1 - (w_1 u_1 + w_2 u_2 + b)) \right. \\ &\quad \left. + (-u_1 - (w_1 u_1 + w_2 u_2 + b)) + (0 - (w_1 u_1 + w_2 u_2 + b)) \right] \\ &= -2(u_1 - (w_1 u_1 + w_2 u_2 + b)) - 2(-u_1 - (w_1 u_1 + w_2 u_2 + b)) + 2(0 - (w_1 u_1 + w_2 u_2 + b)) \\ &= -2(1 - u_1 - b) - 2(1 - u_1 - u_2 - b) - 2b = 0 \end{aligned}$$

$$\cdot \quad 2 - 2u_1 - u_2 - 2b = 0$$

$$\cdot \quad 2 - 2u_2 - u_1 - 2b = 0$$

$$\cdot \quad 2 - 2u_1 - u_1 - 2b = 0$$

$$\cdot \quad 2 - 2u_2 - u_2 - 2b = 0$$

	Output	
	Real Values	Probabilities
Output Activation	Linear	Softmax
Loss function	Squared error	Cross-entropy

→ Output are real values (y regression), then activation f. is linear (bounded by required)

$$\text{sqd. error} \rightarrow \text{loss}(y_{\text{pred}}) = (y_{\text{pred}} - y_{\text{actual}})^2$$

e.g:

$$p_i = \frac{1}{N} \sum_{j=1}^N (\hat{y}_{ij} - y_{ij})^2$$

More generally cost function used is squared error = $\sum_{i=1}^N (y_{\text{pred}} - y_{\text{actual}})^2$

→ When output are probabilities:

$$y_j: \text{classifications} \quad \geq p_i = 1 \\ \text{softmax} \quad \hat{y}_j = D(a_L)_j = e^{a_L j} \quad \text{sigmoid not used as we don't want value to be zero and 1}$$

$$\hat{y}_{ij} = 1 \text{ if class value is true 0 and 0 otherwise}$$

→ Information content: $I(A) = -\log_2(P(A))$
(lower the probability, more the info content)

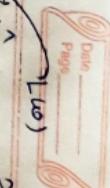
$$\text{Entropy} = -\sum p(i) I(i)$$

$$= -\sum p_i \log(p_i)$$

$$\text{Cross-entropy} : -\sum p_i \log(\hat{p}_i)$$

$p_i \rightarrow \text{true distribution}$
 $\hat{p}_i \rightarrow \text{predicted distribution}$
(minimum when true distribution is as close as predicted distribution.)

Backpropagation in FFNN:



$$\frac{\partial L(\theta)}{\partial w_{11}} = \frac{\partial L(\theta)}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial a_3} \cdot \frac{\partial a_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial a_2} \cdot \frac{\partial a_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial a_1} \cdot \frac{\partial a_1}{\partial w_{11}}$$

$$g: L(0) = -\log \hat{y}_1 \quad (L = \text{true loss})$$

$$(i) \rightarrow \frac{\partial L(\theta)}{\partial \hat{y}_1} = -\frac{\partial}{\partial \hat{y}_1} (-\log \hat{y}_1) = \frac{-1}{\hat{y}_1} \quad \text{if } i=1$$

= 0 otherwise

$$\rightarrow \frac{\partial L(\theta)}{\partial \hat{y}_i} = \frac{-1_{i=1}}{\hat{y}_i} \quad *$$

Gradient

$$\frac{\partial L}{\partial \hat{y}_i} L(\theta) = \begin{bmatrix} \frac{\partial L}{\partial \hat{y}_1} \\ \vdots \\ \frac{\partial L}{\partial \hat{y}_k} \end{bmatrix} = \frac{-1}{\hat{y}_i} \begin{bmatrix} 1_{i=1} \\ 0_{i \neq 1} \end{bmatrix}$$

$$(ii) \frac{\partial L}{\partial a_L} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial a_L} = -\frac{\partial}{\partial (\log \hat{y}_i)} \cdot \frac{\partial \hat{y}_i}{\partial a_L} \quad \text{depends on } a_L \text{ and softmax used}$$

$$= -\frac{1}{\hat{y}_i} \cdot \frac{\partial}{\partial a_L} \text{softmax}(a_L) \leq \frac{e^{a_L}}{e^{a_L}}$$

$$= -\frac{1}{\hat{y}_i} \left(\sum_{k=1}^L \hat{y}_k - \hat{y}_i \hat{y}_i \right)$$

$$= -(1_{i=1} - \hat{y}_i)$$

$$\nabla_{a_L} = \left[\frac{\partial L}{\partial a_L} \right] = \begin{bmatrix} -(1_{i=1} - \hat{y}_i) \\ -(1_{i \neq 1} - \hat{y}_i) \\ \vdots \\ -(1_{i=L} - \hat{y}_i) \end{bmatrix}$$

* "represent" power of MLP:
 → single hidden layer can be used to represent any function ~~precisely~~ (no errors)
 → hidden layer neurons = 2^n (as inputs)

* Represent power of multilayered network of sigmoid neurons:

→ a single hidden layer can be used to approximate any continuous function to any desired precision.
 i.e. given for any f : $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$, we can find neural network with enough neurons in 1 hidden layer where output $g(x)$ satisfies $|g(x) - f(x)| < \epsilon$

Contour

→ more the no. of layers f , better approximation.

→ $f(x) = \text{sum of tower function}$

having $\sigma(x) = \frac{1}{1 + e^{-x}}$ (sigmoid) as output

Neural network:

1D input \rightarrow 2 neurons construct from \rightarrow 4 neurons \rightarrow $O(2^n)$

1	\rightarrow	1 hr
2	\rightarrow	1 hr 10 min
3	\rightarrow	1 hr 90 min
4	\rightarrow	3 hrs
5	\rightarrow	6 hrs
6	\rightarrow	2 hrs
7	\rightarrow	2 hrs
8	\rightarrow	2 hrs

$$2 \cdot 1 + 4 \cdot 2 + 6 \cdot 3 \\ 2 + 8 + 18 \\ + 32$$



SGD example:

Point : $(x_i, y_i) : (1, 2), (2, 4), (3, 6), (4, 8)$

$$y = w_0 + w_1 x$$

(we need to find w that minimizes the sum function MSE)

$$w_0 = 4, \quad E = 0.1$$

$$\nabla_{\theta} J(w) = \frac{1}{4} \left[\sum_{i=1}^4 \nabla_{\theta} L(x^{(i)}, y^{(i)}, \theta) \right]$$

$$L = (y_i - w_0 - w_1 x_i)^2$$

$$\frac{\partial L}{\partial w_0} = 2(y_i - w_0 - w_1 x_i)$$

Gradient descent:

Gradient descent rule :

→ move θ in the direction opposite the gradient.

$$\theta_t = \theta - \eta w_t$$

$$\theta_{t+1} = \begin{bmatrix} w_{t+1} \\ b_{t+1} \end{bmatrix} = \begin{bmatrix} w_t \\ b_t \end{bmatrix} - \eta \begin{bmatrix} \nabla_{w_t} \\ \nabla_{b_t} \end{bmatrix}$$

$$\text{where } \nabla_{w_t} = \frac{\partial L}{\partial w} \Big|_{w=w_t, b=b_t}$$

$$\nabla_{b_t} = \frac{\partial L}{\partial b} \Big|_{w=w_t, b=b_t}$$

$$- \eta \nabla C(w^{(0)}, w^{(1)}, \dots, w^{(t)})$$

→ Randomly split training data into mini-batches.

- larger gradient descent for all minibatch

w away from .

$$w_1' = w_1 + 0.12$$

$$L(\theta) = - \sum_{c=1}^C y_c \log \hat{y}_c$$

$y_c = 1$ if $c = \hat{c}$

$$\Rightarrow L(\theta) = - \log \hat{y}_{\hat{c}}$$

minimize $L(\theta)$ ~~maximise~~ predicted probability of true class

$$\text{or maximum } -L(\theta) = \underbrace{\log \hat{y}_{\hat{c}}}_{\text{log-likelihood of data}} \quad (\hat{y}_{\hat{c}} \text{ is a } \hat{y}^n \text{ in } \hat{y} \theta = [w_1, w_2, \dots, b_1, \dots, b_n])$$

(Maximum log-likelihood is equivalent to minimize M.S.E.)

$$\frac{\partial L(\theta)}{\partial h_{k,j}} = \frac{\partial L(\theta)}{\partial a_{k+1,m}} \cdot \frac{\partial a_{k+1,m}}{\partial h_{k,j}}$$



a gradient w.r.t hidden layer:

$$\begin{aligned} \text{(iii)} \quad \frac{\partial L(\theta)}{\partial h_{k,j}} &= \sum_{m=1}^r \frac{\partial L}{\partial a_{k+1,m}} \cdot \frac{\partial a_{k+1,m}}{\partial h_{k,j}} \quad \text{as} \\ &= \sum_{m=1}^r \frac{\partial L}{\partial a_{k+1,m}} \left(w_{j+1,m,j} \right) \end{aligned}$$

$$\text{(iv)} \quad \nabla_{a_{i,j}} L(\theta) = \begin{bmatrix} \frac{\partial L}{\partial a_{1,j}} \\ \vdots \\ \frac{\partial L}{\partial a_{n,j}} \end{bmatrix} = (w_{i+1})^\top (\nabla_{a_{i+1}} L(\theta))$$

$$\text{(v)} \quad \frac{\partial L(\theta)}{\partial a_{i,j}} = \frac{\partial L(\theta)}{\partial h_{i,j}} \cdot \frac{\partial h_{i,j}}{\partial a_{i,j}} = \frac{\partial L(\theta)}{\partial h_{i,j}} \cdot g'(a_{i,j})$$

+ gradient w.r.t weight:

$$a_k = b_k + w_k h_{k-1}$$

$$\text{(vi)} \quad \frac{\partial a_{k,i}}{\partial w_{k,j}} = h_{k-1,j}$$

$$\nabla_{w_{k,j}} L(\theta) = \left[\frac{\partial}{\partial w_{k,0,0}} \cdots \frac{\partial}{\partial w_{k,n-1}} \right] \cdot \frac{\partial L(\theta)}{\partial a_{k,j}}$$

$$\frac{\partial L(\theta)}{\partial w_{k,j}} = \frac{\partial L(\theta)}{\partial a_{k,j}} \cdot \frac{\partial a_{k,j}}{\partial w_{k,0,0}}$$

$$a_{k,i} = b_{k,i} + \sum_{j=1}^{n_{k-1}} w_{k,i,j} h_{k-1,j} \quad \rightarrow \quad \frac{\partial L(\theta)}{\partial w_{k,i,j}} = \frac{\partial L(\theta)}{\partial a_{k,i}} \cdot \frac{\partial a_{k,i}}{\partial h_{k-1,j}}$$

Stochastic gradient descent.

→ Normal gradient descent:

Use all our entire training data to find Δw and ~~Δb~~

$$\text{Let then update } w = w - \alpha \Delta w, \quad b = b - \alpha \Delta b$$

→ Stochastic gradient descent:

Update in batches (mini-batch) (gradient descent)

e.g.: B (batch size = 1) after every training data encountered.

update w and b .

$$\Delta w = \sum_{i=1}^B (f(w) - y) f'(w) + ((1-f(w))) \cdot \alpha \quad \text{for } \alpha = \sigma(a)$$

($y \rightarrow$

Mean through entire training data (N items)

& then updating w in case of normal gradient descent

→ In 1 epoch (one pass over entire data) : N datapoints

Gradient descent \rightarrow 1 \int $B \rightarrow$ mini-batch size

Stochastic " " $\rightarrow N$ \int N steps in 1 epoch

Mini-batch " " $\rightarrow N_B$ ($1 \leq N_B \leq N$ steps = 1 update of params)



Normal distribution:

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2}$$

mean
standard deviation

Bias & Variance:

- * Bias : let $\hat{f}(n)$ → true model
 $\hat{f}(n) \rightarrow$ estimate of the model

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

avg. value of model

→ simple model has high bias (underfitting)
complex model → overfitting (overfitting)

Variance:

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

→ simple model → low variance
complex " → high variance

overfitting

error margin

MSE:

$$L(\omega) = E[(y - \hat{f}(n))^2] = \text{Bias}^2 + \text{Variance} + \sigma^2 \text{ (brischke error)}$$

model components

→ while training instead of minimizing training error $L(\theta)$, minimize:

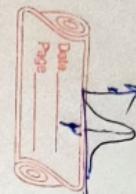
$$\min_{\omega \in \Theta} L_{\text{train}}(\theta) + \underline{\Omega}(\theta) = L(\theta)$$

big for complex model

$$L_2 \text{ regularization: } \tilde{L}(\omega) = L(\omega) + \frac{\alpha}{2} \|\omega\|^2$$

$$\text{for SGD: } \nabla \tilde{L}(\omega) = \nabla L(\omega) + \alpha \omega$$

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



→ Solving we get that to minimize.

$$\int \frac{1}{n} \mathbf{x}^T \mathbf{x} = c$$

$$\min_{\rho_1, \dots, \rho_n} \sum_{j=k+1}^n p_j^T \mathbf{f}_j \quad ; \quad \mathbf{p}_j^T \mathbf{p}_j = 1 \\ \text{(minimizing } \mathbf{p}^T \Sigma \mathbf{p}) \rightarrow \text{ make minimum eigen value})$$

Sol: pick (rk) smallest eigen values of Σ and direct them or pick the k largest eigen vectors.

e.g. let $\mathbf{u}_1 = [1, 1]$, $\mathbf{u}_2 = [-1, 1]$ are new basis vectors.

$$\text{Norm}(\mathbf{p}_2) = \sqrt{1^2 + 1^2} = \sqrt{2}$$

For the point $\mathbf{x} = [3, 3, 3]$

$$\begin{aligned} \hat{\mathbf{x}} &= \mathbf{x}\mathbf{p} = [3, 3] \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} \\ &= \begin{bmatrix} 3\frac{1}{\sqrt{2}} & -3\frac{1}{\sqrt{2}} \end{bmatrix} = [3\sqrt{2}, 3\sqrt{2}] \end{aligned}$$

(i.e. \mathbf{x} can be reconstructed by $a_1 \mathbf{u}_1 + a_2 \mathbf{u}_2 = [3, 3, 3]$)

Note with fewer dimensions ($k = 1 < n$)

(using only \mathbf{u}_1 dimension)

$$\hat{\mathbf{x}} = \mathbf{x}\mathbf{p} = [3, 3, 3] \begin{bmatrix} \frac{1}{\sqrt{n}} \end{bmatrix} = \left[\frac{3}{\sqrt{n}} \right]$$

$$\mathbf{x} = \mathbf{x}_1 \mathbf{u}_1 = \frac{6}{\sqrt{n}} \begin{bmatrix} \frac{1}{\sqrt{n}} \end{bmatrix} = \left[\frac{6}{\sqrt{n}}, \dots, \frac{6}{\sqrt{n}} \right]$$

(reconstruction with minimum error)

→ The n^{th} dimension of transformed point $\hat{\mathbf{x}}$:

$$\mathbf{x}_i = \mathbf{X}\mathbf{p}_i \\ \text{variance along this dimension} = \frac{\sum_{i=1}^m \mathbf{x}_i^T \mathbf{x}_i}{m} = \frac{1}{m} \mathbf{p}_i^T \mathbf{X}^T \mathbf{X} \mathbf{p}_i = \frac{1}{m} \mathbf{p}_i^T \Sigma \mathbf{p}_i \\ = \frac{2}{m} \quad (\text{large eigen value} \Rightarrow \text{high variance})$$

$$\rightarrow \rho = \frac{1}{m} X^T X$$

$\rightarrow X^T X$ is symmetric $\therefore (X^T X)^T = X^T X$

\rightarrow P is orthogonal matrix

$$\text{we want } P^T \Sigma P = D$$

but Σ is a square symmetric matrix

\Rightarrow eigen values of Σ is orthogonal

which orthogonal matrix P diagonalizes Σ ?

Sol^r \rightarrow Matrix P whose columns are eigen vectors of $\Sigma = X^T X$

$$D = V^{-1} \Sigma V = P^T \Sigma P$$

$V \rightarrow$ orthogonal matrix whose columns are eigen vectors of $(\Sigma = X^T X)$

\rightarrow ~~Recall~~ eigen values of $X^T X$ are orthogonal since $(X^T X)$ is symmetric

procedure:

i) Find $\Sigma = X^T X$ (before it make only X orthonormal)

ii) First eigen vectors of Σ : p_1, p_2, \dots, p_n At time P (normalize them vector)

$$\boxed{X = X P}$$

iv) To represent in K dimensions, select eigen vectors corresponding to largest eigen values.

v) To reconstruct X : $\underbrace{X = X P}_{X \rightarrow m \times n} \approx X = m \times K$

$$x_i = \sum_{j=1}^K \alpha_{ij} p_j$$

Ranking the top K dimensions (top K of most non-zero entries)

$$\hat{x}_1 = \sum_{j=1}^K \alpha_{1j} p_j$$

$$\text{Solve } \rho_{11} \text{ to minimize } \rho = \sum_{i=1}^n (x_i - \hat{x}_1)^T (x_i - \hat{x}_1)$$

Dr. Praveen Kumar

A⁻¹ = A^T

Date _____

Page _____

1

It is often necessary to make a
long distance trip.

n_1, n_2, \dots, n_m are data

卷之三

$$\pi_i = \alpha_{i1} p_1 + \alpha_{i2} p_2 + \dots + \alpha_{in} p_n$$

law is abnormal, we can find the dis' we

$$e_1 = \left[\pi_1 - \pi_2 \right] f$$

Spirillum

$$x_i = m_i^T \begin{bmatrix} p \\ \vdots \\ 1 \end{bmatrix} \quad X \rightarrow \text{min}$$

(in direction of transform point)

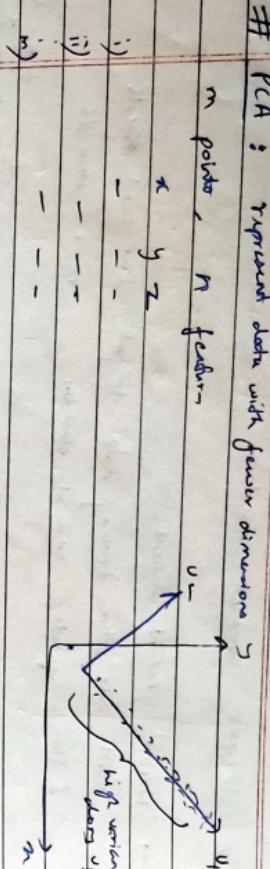
$$\begin{aligned} \hat{x}_i &= x_i^\top p \quad \Rightarrow \quad \hat{X} = Xp = \begin{bmatrix} -m_1 \\ -m_2 \\ \vdots \\ -m_n \end{bmatrix} = \begin{bmatrix} 1 & p_1 \\ 1 & p_2 \\ \vdots & \vdots \\ 1 & p_n \end{bmatrix} \\ \hat{x}_m &= x_m^\top p \\ \text{Now, } X &= Xp \\ \text{covariance of } X &= \frac{1}{m} \sum_{i=1}^m X_i^\top X_i = \frac{1}{m} (Xp)^\top Xp = p^\top \left(\frac{1}{m} X^\top X \right) p \\ &= p^\top \Sigma p \end{aligned}$$

$\Sigma \rightarrow$ convergence mode of X

$$(3) \quad \hat{X}^T \hat{X}_j = 0 \quad i \neq j \quad \text{Condition: } C_{ii} = 0$$

$$\frac{1}{n} \sum_{i=1}^n x_i = \rho \rightarrow \rho = D - \text{second moment}$$

Standardize dataset mean making 0 mean & unit variance



correlation b/w y & z

$$\rho_{yz} = \text{covariance}(y, z) / (\sigma_y \sigma_z)$$

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

$$\begin{aligned} X_{mn} &= \left[\begin{array}{cccc} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{array} \right] & P &= \left[\begin{array}{ccc} \uparrow & \uparrow & \uparrow \\ p_1 & p_2 & \dots & p_n \\ \downarrow & \downarrow & \ddots & \downarrow \\ \vdots & \vdots & \ddots & \vdots \\ m & m & \dots & m \end{array} \right] \\ &&&\text{Error } (y_i - \hat{y}) \\ &&&= \frac{1}{m} \sum_{i=1}^m (y_i - \bar{y})^2 \\ &&&= \frac{1}{m} \sum_{i=1}^m (z_i - \bar{z})^2 \end{aligned}$$

+ Covariance matrix : (c_{ij}) stores sum covariance b/w column i & j

→ If the columns of X are 0 mean

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

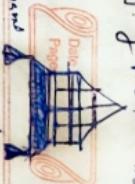
$$= \frac{1}{m} \sum_{k=1}^m x_{ki} x_{kj} \quad (\because \mu_i = \mu_j = 0)$$

$$= \frac{1}{m} \sum_{k=1}^m x_i^T x_j = \frac{1}{m} (X^T X)_{ij}$$

$$\therefore \text{Cov. } C = \sum_{i,j} = \boxed{\frac{1}{m} X^T X} \rightsquigarrow \text{covariance matrix}$$

- We update by taking the weighted avg. of previous moment
- be gradient at the previous position
 $(\theta_t + \alpha v_{t-1})$
- Momentum term adds the influence of previous gradient.
- more robust against oscillations & slow convergence compared to simple momentum term

Batch Normalization:



$$L = \frac{1}{N} \sum_i J(f(x_i, \theta), y_i)$$

$$L = \frac{1}{m} \sum_j f(\mu_j, \sigma_j x_j)$$

$$\bar{\theta}_t = \theta_t + \alpha v_t$$

$$\theta^{t+1} \leftarrow \theta^t - \gamma \nabla L$$

(Problem: High variance.
K.F. are unstable
but we need)

$$\text{Consider } \beta = (x_1, x_2, \dots, x_m)$$

$$\rightarrow L = \frac{1}{m} \sum_i J(c_i)$$

* Trained covariate shift:

when the distribution of the input features vary across the mini-batch.

g: feature extraction
cat:

In our minibatch with center

- layer is normalized
- incurred by standardizing the pre-activation
- at each unit.

expanding values over entire

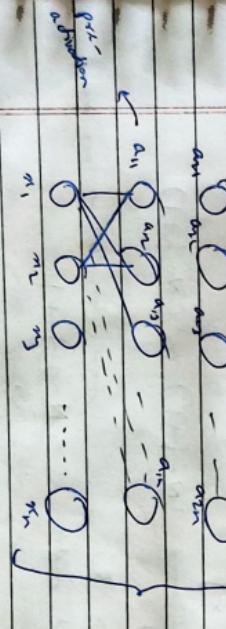
in another "step back" into

Convergence will take time

SE [air]

Scaling and transformation

(θ scaling and transforma)



$$\text{act}(x) = \text{act}(Bw + b) = y^{(k)} = \gamma^{(k)} \hat{w} + \beta^{(k)}$$

(rescaling and

retransformation)

In final output layer

↓ not motion

Homework: With batch normalization, how gradient descent to converge

propagation with worth

Optimization:

i)

Gradient descent:

$$\theta^{t+1} = \theta^t - \epsilon \nabla_{\theta} J$$

$$J = \frac{1}{N} \sum_{i=1}^N J(F(x_i, \theta), y_i)$$

Problem → high computation time as all N samples training data contain

ii)

SGD: batch size = 1 → choose a random point → let me

$$J(x, \theta) = \frac{1}{N} \sum_{i=1}^N J(F(x_i, \theta), y_i)$$

Problem → slow convergence

iii)

$$\text{Minibatch} \rightarrow L = \frac{1}{k} \sum_{i=1}^k J(F(x_i, \theta), y_i)$$

+

Batchsize: optimizing E (overheads) or slow convergence

Solution: decay learning rate:

$$E_K = (1-\alpha) E_0 + \alpha E_2 \quad \text{where } \alpha = \frac{k}{K}$$

$$E_2 \ll E_0 \quad (\text{typically } E_2 < 0.01 E_0) \quad \frac{1}{k}$$

k → no. of iterations → K → total no. of iterations

$$\rightarrow d\phi = \left(\frac{\partial \phi}{\partial x} \uparrow + \frac{\partial \phi}{\partial y} \uparrow + \frac{\partial \phi}{\partial z} \uparrow \right) \cdot \left(dx \uparrow + dy \uparrow + dz \uparrow \right)$$

$$= \nabla \phi \cdot \vec{n}$$

→ $\nabla \phi$ points to a direction where do the steepest surface $\|\nabla \phi\| = |\nabla \phi| / N$

$$d\phi = \nabla \phi \cdot \vec{n}$$

$$= 1.0 \uparrow |1.0| \cos \theta$$

$|\vec{D}\phi| \rightarrow$ direction of steepest ascent

+ Handling oscillations: momentum-based technique:

~~Momentum~~

$$V(t) = \omega \sqrt{t+1} - E g_t$$

$\theta^{t+1} = \theta^t + V^t$ θ^0 → randomly initialized

V^t → momentum (remember $\vec{v} = m \vec{a}$)

Terminal velocity: $E(g)$ → $(-g)$ direction of learning rate.

$$(1-\alpha)$$

$$V = -E g_t \quad E(g) = E(g_1) - E(g_2)$$

$$V_1 = -E g_1 \quad E(g_1) = E(g_1) - E(g_2)$$

$$V_2 = -E g_2 \quad E(g_2) = E(g_2) - E(g_3)$$

$$V_3 = -E g_3 \quad E(g_3) = E(g_3) - E(g_4)$$

$$V_4 = -E g_4 \quad E(g_4) = E(g_4) - E(g_5)$$

$$V_5 = -E g_5 \quad E(g_5) = E(g_5) - E(g_6)$$

$$V_6 = -E g_6 \quad E(g_6) = E(g_6) - E(g_7)$$

$$V_7 = -E g_7 \quad E(g_7) = E(g_7) - E(g_8)$$

$$V_8 = -E g_8 \quad E(g_8) = E(g_8) - E(g_9)$$

$$V_9 = -E g_9 \quad E(g_9) = E(g_9) - E(g_{10})$$

$$V_{10} = -E g_{10} \quad E(g_{10}) = E(g_{10}) - E(g_1)$$

$$V_1 = -E g_1 \quad E(g_1) = E(g_1) - E(g_2)$$

$$V_2 = -E g_2 \quad E(g_2) = E(g_2) - E(g_3)$$

$$V_3 = -E g_3 \quad E(g_3) = E(g_3) - E(g_4)$$

$$V_4 = -E g_4 \quad E(g_4) = E(g_4) - E(g_5)$$

$$V_5 = -E g_5 \quad E(g_5) = E(g_5) - E(g_6)$$

$$V_6 = -E g_6 \quad E(g_6) = E(g_6) - E(g_7)$$

$$V_7 = -E g_7 \quad E(g_7) = E(g_7) - E(g_8)$$

$$V_8 = -E g_8 \quad E(g_8) = E(g_8) - E(g_9)$$

$$V_9 = -E g_9 \quad E(g_9) = E(g_9) - E(g_{10})$$

$$V_{10} = -E g_{10} \quad E(g_{10}) = E(g_{10}) - E(g_1)$$

$$E(\xi_i) = V$$



→ If perfectly correlated $C = V$
⇒ $\text{var} = \frac{V}{K} - \frac{(K-1)V}{K} = V$ (biasing downside)

Dropout:

→ dropping neurons of neural network (random node of the network having connection)
→ for each round, retain nodes with a fixed prob. (say 0.5 for hidden nodes) and ($P = 0.8$) for visible (input) nodes.

- * In each round, take one training instance (mini-batch)
- "Dropout proposal", compute loss & back propagate
- (update weights which were active in that round)

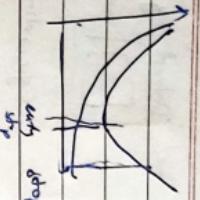
- * At test time,
use full neural network & scale the output of each node
by P (fraction of times it was "on" during training)

Regularization

* Early stopping,

$$w_{t+1} = w_t - \eta \sum_{i=1}^t \nabla w_i$$

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



- let $t \in \mathbb{N}$ be the max value of ∇w_i
- for early stopping, allow only 't' no. of updates
- Early stopping is also related to L2 regularized -

Ensemble methods:

- combine output of different models to reduce test error
- e.g: $\logistic\ regression$ { some weights }

never been

- Training can be done with diff. hyperparams, features,

→ or diff. samples of training data

* Bagging:

- ensemble method using diff. instances of same classifier
- construct multiple training sets by sampling with replacement (T_1, T_2, \dots, T_k)
- Train i'm classifier with training set T_i
- better if the ~~classifiers~~ errors across classifiers are independent

~~Ex:~~ k models of linear regression

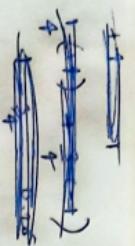
Error for each model $\rightarrow E_i$

$$E \rightarrow \text{average mean multivariate normal distribution}$$

$$E(C_i c_j) = V \quad E(E_i c_j) = C$$

$$mse = E \left(\left(\frac{1}{K} \sum_{i=1}^K (y_i - \hat{y}_i)^2 \right) \right) = \frac{V}{K} + \frac{(K-1)}{K} C$$

If no correlation, $mse = \frac{V}{K}$ (error is now reduced by factor of K)



→ Dataset augmentation:



Generating more training samples from existing one
not training model is more generalized.

* Adding noise to the inputs:

$$\tilde{x}_i = x_i + \epsilon_i$$

$$\hat{y} = \sum_i w_i x_i \rightarrow \text{predicted without noise } w_i + \epsilon_i \text{ due to noise, then } \hat{y} + \epsilon$$

$\tilde{y} = \sum_i w_i \tilde{x}_i \rightarrow \text{with noise}$

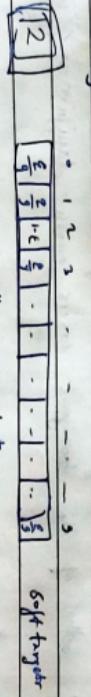
$$E(\epsilon_i) = 0$$

$$E[(\hat{y} - y)^2] = E[(\sum_i w_i x_i + \epsilon - y)^2] = \sum_i w_i^2 E(x_i^2)$$

Non-square error	Same as L2
without noise	minimum penalty
(training error)	$\frac{1}{n} \sum \epsilon_i^2$

∴ Adding noise to input is equivalent to L2 regularization
for gaussian normal distribution

* Adding noise to outputs:



$\epsilon \rightarrow$ small & un constant

$$\text{minimize } \sum_{i=1}^n l(y_i, \hat{y}_i) \quad (\text{loss function})$$

$$\text{true distribution + noise: } p = \left[\frac{E}{3}, \frac{E}{3}, 1 - \frac{E}{3}, \frac{E}{3} \dots \right]$$

estimated distribution

Noise minimization

Formation backpropagation: overcomes obstacles, activates convergence, provides stability by modifying a constant direction & accumulating influence of past gradient.

→ Backpropagation idea: decrease learning rate

→ Idea: decrease learning rate to prevent backpropagation from having large gradient & increase it for learning with small gradients.

Problems with DFNN:

- too many pixels ($10^3 \times 10^3$) → each pixel → $R \times G \times B$ (too many inputs)
- localizing features (low bias) → spatial / temporal dependence
- transformation doesn't affect overall result (equivariance)

Image properties:

- i) locality
- ii) stationarity
- iii) hierarchy (image is made up of multiple smaller sub-songs)

Tasks in neural:

i) convolution layer:

→ Impose kernel on the image & do element wise multiplication

Image — $N \times N$

$$(N-f+1) \times (N-f+1)$$

Kernel — $f \times f$

$$(N-f+1) \times (N-f+1)$$

Resultant matrix — (when stride = 1)

If option = 5 (move 5 steps in one row)
then resultant matrix = $\frac{(N-f)}{5} + 1$

Image: $32 \times 32 \times 3$ (3 layers)

Image size \rightarrow 3 channels (RGB)

Kernel \rightarrow e.g. one particular feature
Dot product of image & kernel → similarity value

→ Use multiple kernels for diff. features.

→ predefined kernels

→ learned kernels

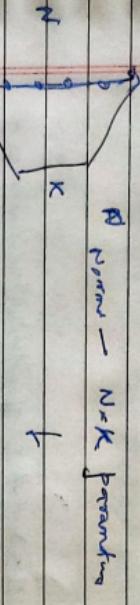


Advantage:

(i) → significantly reduces no. of ~~feature~~ parameters

CNN → compressing the receptive field of a unit (neuron)

receptive field → No. of units in lower layer that affect the output of higher layer



• Normal - $N \times K$ receptive field

$f \times f \times k$

(very CNN)

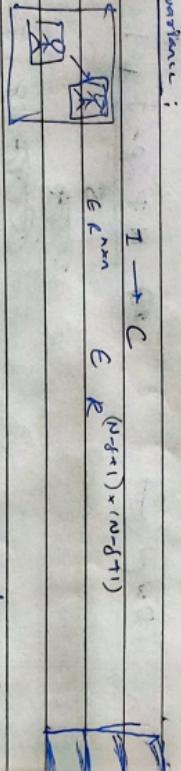
→ To prevent loss, increase size of kernel stride

(iii) regions similar in values will output similar values
(locality)

* Invariance:

$I \rightarrow C$

$C_{\text{CNN}} \in \mathbb{R}^{(N \times d \times 1) \times (N \times d \times 1)}$



→ translational invariance → movement of objects within the image shouldn't change the convolutional output (C) much

$$f(g(x)) = f(x)$$

$f \rightarrow$ convolution

→ Equivariance: Applying f after transform \equiv transform the output of f

$$f(g(cx)) \approx g(f(cx))$$

Assumption: Order of (layering) will not change if kernels are swapped.

Even with varying height, the no. of partition to
seam remains same.

2 (softmax)

forward prop:

start with random values of all partition

and no. complete all the values

h₀ → h₁ → h₂ → ... → h_n

backpropagation

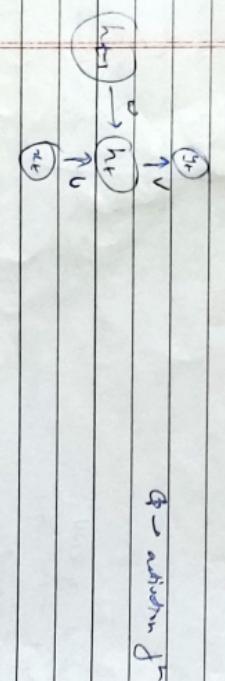
→ N. of stage needs to be decided

$$\text{Loss function: } \frac{1}{2} \sum_{k=1}^n (\hat{y}_k - y_k)^2$$

$$(h_0) \xrightarrow{\theta_0} (h_1) \xrightarrow{\theta_1} (h_2) \dots \xrightarrow{\theta_n} (h_n)$$

$$\text{Overall loss} \rightarrow \frac{1}{2} \sum_{k=1}^n (\hat{y}_k - y_k)^2$$

$$(\theta_0) \xrightarrow{\theta_1} (\theta_2) \dots \xrightarrow{\theta_n}$$



(θ → activation fn)

$$\frac{\partial L}{\partial h_k} = \prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}} \quad L = ab$$

$$\frac{\partial L}{\partial h_k} = \prod_{i=k+1}^t \frac{\partial h_i}{\partial w_i} \frac{\partial w_i}{\partial b_i} \frac{\partial b_i}{\partial z_i} \quad \frac{\partial L}{\partial z} = a^\top b + \frac{1}{2} \frac{\|z\|^2}{\lambda}$$

h_i → vector

$$h_i = w^\top a + b (h_{i-1})$$

$$\frac{\partial L}{\partial h_i} = w^\top \frac{\partial L}{\partial h_{i-1}}$$

WELLS



Residual Network (resnet)

- ↳ can learn identity layer
- ↳ F(x) + x

→ But as we keep on ↑ depth in layers,
after certain point → error increases due to overfitting

Recurrent Neural Network (RNN):

problem with CNN
words → vector
(fixed length)
{similar words have similar vectors}

→ want
and then apply CNN
to
text

→ problem:
2012 → variable length short (image)
(height)

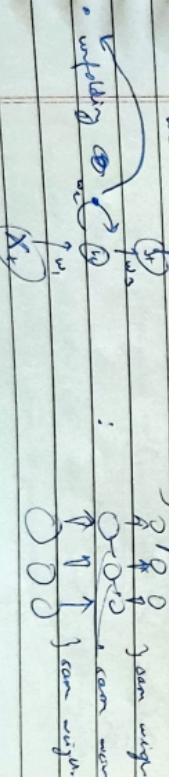
→ a word may be repeated
on a road at a long distance
This may not be captured.

→ In 2012, I went to Eiffel Tower
with generate different words.

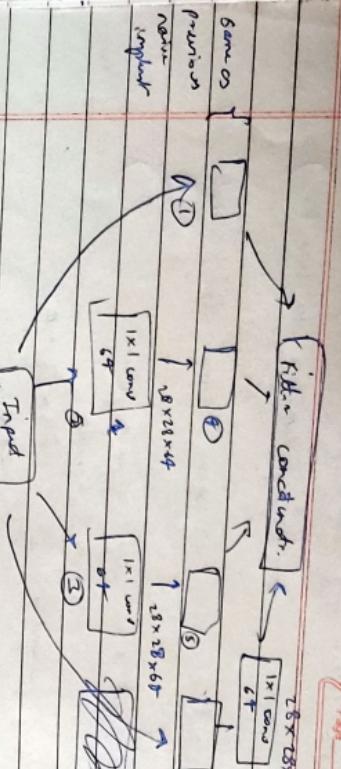
- #### Applications of RNN :
- i) Input writing, output
 - ii) prose generation
 - iii) sentiment analysis
 - iv) text - to - speech
 - v) machine translation (Machine learning 1 on language 2)

RNN →

we share parameters across long sequences



perceive the image at multiple scales simultaneously.



Date:

1

1

1

1

1

1

1

1

1

$$29 \times 28 \times 25 \times 1 - 1 \times 129$$

$$78 \times 78 \times 7 + 3 \times 3 \times 16 = 14$$

$$v_1) \quad 29 \times 23 \times 256 \times 1 \times 1 \times 64$$

H 111 - 1

(such training & teaching errors) ↳ Not due to own

100

100

100

India's Report

100

1111

hard produce sum results
but higher depth performs well
difficult to learn

→ Why increasing depths may give more terror.
(such training & today error)
↳ Not used to outfitting
as → lighter training error and also terrorist
problem may be due to frustration. ↳ (of weight)

Deeper model \rightarrow more non-linearity

\rightarrow increase receptive field

smaller filters \rightarrow fewer parameters

($n \times n$) $\rightarrow n^2$ parameters to learn

Vgg Net : (19 layers) (more depth, smaller filters)

3×3 filters \rightarrow 3 channels

$$\begin{aligned} \text{No. of parameters} &= 7 \times 3 \times 3 \times c \\ &= 21c \end{aligned}$$

\rightarrow complexity of CNN's :

measured by the # of floating point operations (FLOPs)

Image Net :

\rightarrow Native Inception :

(filtering at diff. scale
and then combining)

Native concatenation

(padding added)

Output:

$$\begin{aligned} 28 \times 28 \times 128 &\xrightarrow{\text{1x1 conv}} 28 \times 28 \times 128 \\ 128 &\xrightarrow{\text{3x3 conv}} 192 \\ \text{filter} &\xrightarrow{(i)} 192 \xrightarrow{(ii)} 96 \xrightarrow{(iii)} 3 \times 3 \text{ kernel} \\ &= 28 \times 28 \times 672 \end{aligned}$$

Input

$28 \times 28 \times 3$

Consider the no. of multiplications:

$$\begin{aligned} i) & 28 \times 28 \times 256 \times 1 \times 1 \times 128 = \text{Total FLOPs} = \\ ii) & 28 \times 28 \times 256 \times 3 \times 3 \times 192 = 854 \text{ million} \\ iii) & 28 \times 28 \times 256 \times 5 \times 5 \times 96 \\ iv) & \dots \end{aligned}$$

100
Bücher

1 hidden page can prevent any arbitrary file
but under may be executed.

— said — Korean Dept.

Image — 227 x 227 x 3

$$\text{filter} = (11 \times 11 \times 3)$$

$$\text{CON} \rightarrow \text{SIVAN} -$$

Homework

$$= \frac{55-3}{2} + 1$$

C. 15 min

~~Conv.~~ 5×5 filters \rightarrow 256 filters

~~Input = 1~~

~~pad = 2~~

$$\text{Output} \rightarrow \text{N} + \underline{2p - f} + 1$$

$$\begin{array}{r} \cancel{2} \\ - \cancel{2} \\ \hline 27 + 4 - 5 + 1 \end{array}$$

Output → 27 x 27 x 256

$$\text{Mox 100112 : } 3 \times 3 \text{ Jukts, } \sin = n \\ \text{Dabat } \rightarrow \left(\frac{27 - 3}{2} + 1 \right) = 13 \times 13 \times 255$$

$$\text{Durchschnitt} \rightarrow \left(\frac{27 - 3}{2} + 1 \right) =$$

$$\phi_{k,i-1} = \begin{cases} \phi_{k,i-1}^{(1)} & h_i = \int h_i^{(1)} \\ \phi_{k,i-1}^{(2)} & h_i^{(2)} \\ \vdots & \vdots \\ \phi_{k,i-1}^{(n)} & h_i^{(n)} \end{cases}$$

$$\frac{\partial \phi_{k,i-1}}{\partial h_i} = \left[\frac{\partial \phi_{k,i-1}^{(1)}}{\partial h_i^{(1)}} \quad \frac{\partial \phi_{k,i-1}^{(1)}}{\partial h_i^{(2)}} \quad \dots \right]$$

$$= \left[\frac{\partial \phi_{k,i-1}^{(1)}}{\partial h_{i-1}^{(1)}} \quad \frac{\partial \phi_{k,i-1}^{(1)}}{\partial h_{i-1}^{(2)}} \quad \dots \right]$$

$$= \left[\frac{\partial \phi_{k,i-1}^{(1)}}{\partial h_{i-1}^{(1)}} \quad \frac{\partial \phi_{k,i-1}^{(1)}}{\partial h_{i-1}^{(2)}} \quad \dots \right]$$

$$= \left[\frac{\partial \phi_{k,i-1}^{(1)}}{\partial h_{i-1}^{(1)}} \quad \frac{\partial \phi_{k,i-1}^{(1)}}{\partial h_{i-1}^{(2)}} \quad \dots \right]$$

$$= \begin{bmatrix} \phi_{k,i-1}^{(1)} & 0 & 0 & \dots \\ 0 & \phi_{k,i-1}^{(2)} & & \\ & \ddots & \ddots & \\ 0 & & \ddots & \ddots \end{bmatrix}$$

+ Variating gradient : due to layer RNN due

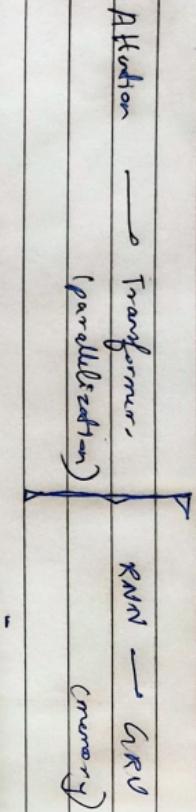
$$\frac{\partial \phi_{k,i-1}}{\partial h_i} = \frac{\partial \phi_{k,i-1}}{\partial h_{i-1}}$$

(RNN does learning)

(gradient = 0)

- i) divide them into batches
- do a back propagation over a smaller no. of steps
- gradient is calculated over a small batch only but program starts from beginning to end

Paper: Attention is all you need



BERT: Bidirectional Encoder representation using Transformers
→ pre-trained models

Paper: Pre-training of Deep bidirectional Transformers
for language Understanding
→ Jacob Devlin, et. al

Autoencoders :



ReLU

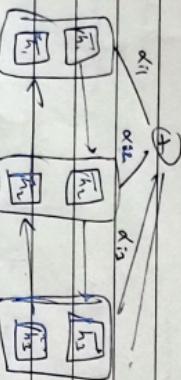
CNN, RNN, GRU, LSTM,
(+ bidirectional)

Attention & Transformer, Autoencoder

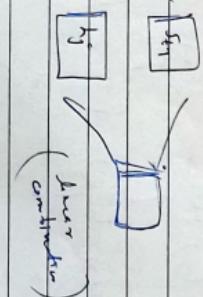
~~Attention~~

Attention layer :

$$\begin{aligned} s_{t-1} &\rightarrow s_t \\ c &= \sum_{j=1}^n \alpha_{tj} h_j \end{aligned}$$



$$\alpha_{tj} = \exp(\tilde{\alpha}_{tj})$$
$$\tilde{\alpha}_{tj} = \sum_{j=1}^n \exp(\alpha_{tj})$$
$$\alpha_{tj} = f(s_{t-1}, h_j)$$



Attention and Transformer

Date: 10/04/23
Page:

→ Paper:

Neural machine translation by jointly learning to align and translate

- Bahdanau et al ICLR 2015

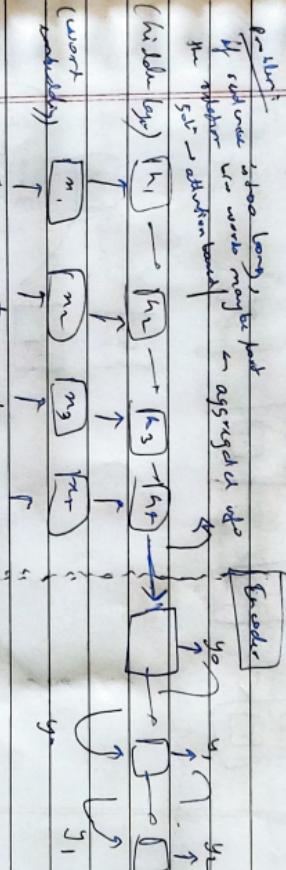
The cat eat mouse

Belle chien - kawaii mai

translation

alignment

X.T



The cat eat mouse

RNN based model:

$$\begin{aligned} X &= (x_1, x_2, \dots, x_T) \\ \text{Encoder} \quad h_t &= f(h_{t-1}, x_t) \end{aligned}$$

$$\text{Context} \quad \leftarrow c = g(h_1, h_2, \dots, h_T)$$

$$\text{Decoder} \quad p(y_t) = \prod_{t'=1}^T p(y_t | y_{t'}, s_t, c)$$

$$p(y_t | y_{t'}, s_t, c) = g(y_{t'}, s_t, c)$$

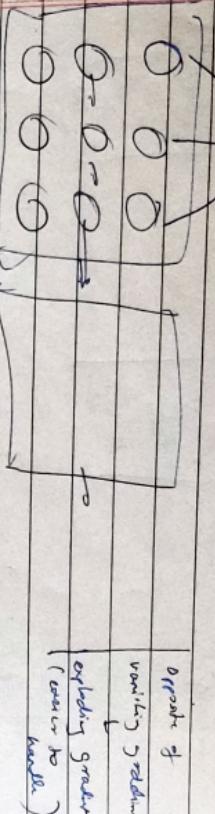


MAVELLS

Q: Wikipedia article - Summary of document
 $\begin{cases} \text{O} \rightarrow \text{prior} \\ \text{O}_t \rightarrow \text{prior} \\ \text{O} \rightarrow \text{prior} \end{cases}$

10/05

Date _____
Page _____



$$h_t = Vx_t + W\phi(h_{t-1})$$

$m \times 1$ $m \times 1$

$$\text{New input} \rightarrow T = \begin{bmatrix} h_{t-1} \\ h_t \\ h_{t+1} \end{bmatrix} \quad W_a = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix}^T$$

$2m \times 1$

$$h_t = \phi(W_a [h_{t-1}, r_t] + b_a)$$

Applications of RNN:

i) Language generation & Sequence Model.

Suppose y_1, y_2, \dots, y_n be the words in a language

Goal of a language model!

- Given the probability of a given sequence of words
- helpful in tasks like speech recognition