

Types of attributes / features / dimensions / variables.
→ we want to transform features those have non-numeric values in the way so we may able to perform ML on them.

Q what is attribute?

A data field that representing a characteristic or feature of a data object

E.g. customer-ID, name etc

Types of attributes w.r.t to data they contained.

Nominal: categories, states qualities, labels, names of things.

Hair-color, marital status, occupation, ID numbers, zip codes

Binary: Nominal attribute with only 2 states (0 and 1)

(i) Symmetric (ii) Asymmetric

Ordinal: values have meaningful order (ranking). but magnitude between them is not known/quantifiable. e.g. size, grades, army ranks

Numeric: Quantitative.

Base value is not 0.

Interval: No true-zero point exist / starting point is not 0
°C / F°, calendar dates.

Ratio

True-zero point exist / starting point 0 /
Base value is 0

Kelvin

Discrete attributes → Represented by integers

Continuous attributes → Represented by real values (floating values)

→ we assign id to nominal attributes.

Types of attributes / features / dimensions / variables.

→ we want to transform features have non-numeric values in the way so we may able

to perform ML on them.

Q what is attribute?

A data field that representing a characteristic or feature of a data object.

e.g. customer-ID, name etc

Types of attributes w.r.t to data they contained.

Nominal: categories, labels, qualities, names of things.

order of
is not
hair-color, marital status, occupation, ID numbers, zip codes.

Important Binary: Nominal attributes with only 2 states (0 and 1)

(i) Symmetric (ii) Asymmetric

Ordinal: values have meaningful order (ranking), but magnitude

is important of difference between them is not known/quantifiable. e.g. size, grades, army rank.

Numeric: Quantitative

Interval: No true-zero point exist / floating point is not 0

0° / 10°, calendar dates.

Ratio: True-zero point exist / floating point 0 /

Kelvin represented by integers

Discrete attributes → Represented by integers

Continuous attributes → Represented by real values (floating values)

→ we assign id to nominal attributes.

Similarity and Dissimilarity between attributes

- Proximity - refers to a similarity or dissimilarity

- Similarity are

Numerical measure that tells how alike two objects value is higher when objects are more alike

values fall in range [0, 1]

0 means ^{objects} not matched / similar at all.

1 means ^{objects} exactly matched / similar

Dissimilarity

Numerical measure that tells how different two objects are.

value is higher when objects are more different.

Range for dissimilarity [0, varies

0 means objects exactly matched / similar or less dissimilar upper limit values.

Data matrix

n = no. of example point. in dataset.

example point.

p = no. of features / dimension / attributes / variables of every

$x_{11} \dots x_{1p} \dots x_{1p}$		1	... 5 ... 10
...	...	2	
$x_{21} \dots x_{2p} \dots x_{2p}$	e.g	50	
...	...	1	
$x_{n1} \dots x_{np} \dots x_{np}$	nxp	100	100x10.

we make dissimilarity matrix on the basis of n .

Dissimilarity Matrix

- n example points but registers only the distance

- A triangular matrix of size $n \times n$.

Distance formula is dissimilarity measure.
 Example point = Data point = Data object.

If $n=4$ then dissimilarity matrix is a triangular matrix of size $n \times n = 4 \times 4$.

$E_1, E_2, E_3, E_4, i, j = 1, 2, 3, 4$.

	E_1	E_2	E_3	E_4
E_1	0			
E_2		0		
E_3			0	
E_4				0

$n \times n$.

Proximity measure for nominal attributes.

- Attributes can take 2 or more states.

(Generalization of Binary attributes)

Method:- Simple matching.

(1) First to calculate total no. of Nominal attributes in dataset

P = Total no. of nominal attributes in dataset.

m = No. of matches.

then

$$d(i, j) = P - m \quad S(i, j) = m$$

Proximity measure for Binary attributes P .

contingency table for binary attribute objects:

Total matches = $a + t$

	1	0	Sum	Total non-matches = $r + s$
1	a	r	$a + r$	Total = $a + r + s + t$
0	s	t	$s + t$	
Sum	$a + s$	$r + t$		

for Symmetric Binary attributes for Asymmetric Binary attributes

$$d(i, j) = s + r / a + r + s + t$$

$$S(i, j) = a + t / a + r + s + t$$

$$d(i, j) = s + r / a + r + s + t$$

$$S(i, j) = a / a + r + s$$

Jaccard coefficient (Similarity measure)

Standardization/Normalization of Numeric Data

① Min-Max Normalization: to $[new_min_A, new_max_A]$

$$\hat{v} = \frac{v - min_A}{(new\ max_A - new\ min_A)} + new\ min_A$$

Here: $\frac{max_A - min_A}{}$

v represents numeric feature value

\hat{v} represents transformed/normalized numeric feature value.

min_A : minimum value of feature before normalization

max_A : maximum value of feature before normalization

$new\ min_A$: new minimum value of feature that you want to set after normalization

$new\ max_A$: new maximum value of feature that you want to set after normalization

Example let feature(income) range 12000 to 98000

normalized to $[0.0, 1.0]$. Then 73600 is mapped to

$$\frac{73600 - 12000}{98000 - 12000} (1.0 - 0.0) + 0.0 = 0.716$$

② Z-score normalization (μ : mean, σ : standard deviation)

$$\hat{v} = \frac{v - \mu_A}{\sigma_A}$$

Here μ_A

μ_A : mean of attribute / feature

σ_A : standard deviation of attribute / feature

Example let $\mu = 54000$, $\sigma = 16000$ then 73600 is mapped to

$$= \frac{73600 - 54000}{16000} = 1.225$$

③ Normalization by decimal scaling

$$\hat{v} = \frac{v}{10^j}$$
 where j is the smallest integer such that $\max(|\hat{v}|) \leq 1$

Special cases of Minkowski Distance. $d(i, j) = \left(\sum_{f=1}^p |x_{if} - x_{jf}|^n \right)^{1/n}$

① $n=1$: Manhattan (city block, L₁ norm) distance.

E.g., the Hamming distance: the number of bits that are different between two binary vectors.

$$d(i, j) = \left[\sum_{f=1}^p |x_{if} - x_{jf}|^1 \right]^{1/1} = \sum_{f=1}^p |x_{if} - x_{jf}|$$

$$= |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

② $n=2$: Euclidean (L₂ norm) distance.

$$d(i, j) = \left(\sum_{f=1}^p |x_{if} - x_{jf}|^2 \right)^{1/2} = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2}$$

③ $n \rightarrow \infty$: Supremum (L_{max} norm, L_∞ norm) distance.

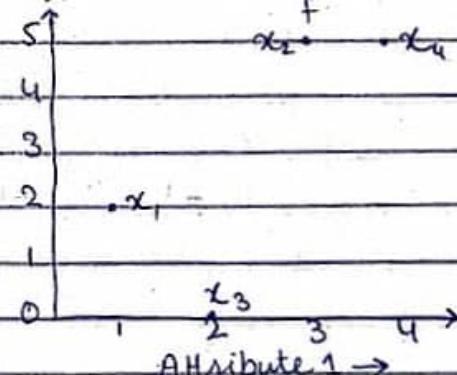
This is the maximum difference between any component attribute of the vectors.

$$d(i, j) = \lim_{n \rightarrow \infty} \left(\sum_{f=1}^p |x_{if} - x_{jf}|^n \right)^{1/n} = \max_f |x_{if} - x_{jf}|$$

Example

Points attribute1 attribute2

	x ₁	x ₂	x ₃	x ₄
x ₁	1	2	3	
x ₂	3	5	Attribute 2	
x ₃	2	0	2	1
x ₄	4	5	0	



Dissimilarity Matrices.

Manhattan (L₁) Euclidean (L₂) Supremum (L_∞)

L ₁ , L ₀ , L _∞	x ₁	x ₂	x ₃	x ₄
x ₁	0,0,0	-	-	-
x ₂	5,3,6,3	0,0,0	-	-
x ₃	3,2,4,2	6,5,1,5	0,0,0	-
x ₄	6,4,2,3	1,1,1	7,5,3,9,5	0,0,0

Ordinal variable/attributes.

- An ordinal variable can be discrete or continuous
- order is important e.g. rank
- can be treated like interval-scaled
 - Assign id's to all ranks or to all possible values of ordinal
 - for example universities teachers designation feature

designation lecturer Assistant-Lecturer Associate Prof Professor

id's 1 2 3 4

- $r_{if} \in \{1, \dots, M_f\}$

- r_{if} represents corresponding id or rank value to i^{th} example point in f^{th} feature (ordinal feature)

- 1 represents lowest rank

- M_f represents Maximum/highest rank.

- Replace x_{if} by r_{if}

x_{if} represents i^{th} example point in f^{th} feature.

- After replacing map f^{th} feature onto $[0, 1]$ by min-max Normalization

$$z_{if} = (r_{if} - 1) / (M_f - 1)$$

- compute the dissimilarity using methods for interval scaled like Euclidean Distance method

Attributes of Mixed Type

- A database may contain all attribute types (Nominal, Binary, Asymmetric binary, numeric, ordinal). then compute distance/dissimilarity matrices for all types
- one may use a weighted formula to combine their effects:

$$d(i, j) = \sum_{f=1}^F \delta_f d_{ij}(f) / \sum_{f=1}^F \delta_f$$

→ f is binary or nominal $d_{ij}(f) = 0$ if $x_{if} = x_{jf}$ or $d_{ij} = 1$ otherwise

→ f is numeric use normalized distance $d_{ij}(f) = |x_{if} - x_{jf}| / \max_i x_{if} - \min_i x_{if}$

→ f is ordinal

compute ranks r_{if} and r_{jf} as interval scaled $z_{if} = \frac{r_{if} - 1}{M_f - 1}$

Simultaneous evaluation is not scaleable but GD is
Gradient Descent

- Gradient Descent is just like Agile Methodology.

- Agile Methodology (Iterative approach).

① Build something quickly (Prototype).

② Get it out there (Deploy or Test)

③ Get some feedback (Shared with customer)

④ Make changes depending upon the feedback

⑤ Repeat.

- Similarly agile methodology we use iterative approach in gradient descent first we quickly build model and then evaluate it and compare results with actuals results and repeat this process.

- Example.

- lets have some function say $j(\theta)$. θ may be one parameter or list of parameters.

- objective function

Either we maximize or minimize $\theta = [\theta_0, \theta_1, \dots, \theta_n]$

the function ($j(\theta)$) by finding specific values of θ . $\max j(\theta)$ or $\min j(\theta)$.

- lets we want to find such values of θ for which $j(\theta)$ return minimum value. so to achieve this objective we use gradient descent

- working of GDA

search for random initialization techniques.
for weights.

- Randomly initialize θ 's. { usually range is $-0.5 \text{ to } +0.5$ }

- Keep changing θ 's to reduce/minimize $j(\theta)$ until we hopefully end up at a minimum value of $j(\theta)$ by using weight updated equations.

- repeat until convergence.

$$\theta_i := \theta_i - \alpha \frac{\partial j(\theta)}{\partial \theta_i} \rightarrow \text{weight update equation}$$

convex concave

2.E 3.1

Weight update equation for minimization

$$\theta_i := \theta_i - \alpha \frac{\partial}{\partial \theta_i} j(\theta)$$

Weight update equation for maximization.

$$\theta_i := \theta_i + \alpha \frac{\partial}{\partial \theta_i} j(\theta)$$

α is called eta

- step size.

- learning rate.

operator.	In programming languages	In Mathematics
Assignment	=	$::=$
Equal	$==$	$=$

→ Derivative return rate of change or slope of function.

Dry run

minimization

tells there is convex problem and function is $j(\theta) = (\theta - 3)^2 + 5$.

θ_j $j(\theta_j)$

↑

-4 54

$j(\theta)$

-3 41

100

-2 30

90

-1 21

80

0 14

70

1 9

60

2 6

50

3 5

40

4 6

30

5 9

20

6 14

10

7 21

1

8 30

2

9 41

3

10 54

4

-5 -4 -3 -2 -1 0 1 2 3 4 5 6 7 8 9 10

$\leftarrow \theta_i \rightarrow$

unknown parameter

→ No. of weight update equation = No. of weights in function.

So by using GD for minimization

$$\theta_i := \theta_i - \alpha \frac{\partial f(\theta_i)}{\partial \theta_i}$$

$$\frac{\partial f(\theta_i)}{\partial \theta_i} = \frac{\partial}{\partial \theta_i} (\theta_i - 3)^2 + 5 = 2(\theta_i - 3) \cdot \frac{\partial}{\partial \theta_i} (\theta_i - 3) + \frac{\partial}{\partial \theta_i} (5) = 2(\theta_i - 3)$$

$$\text{let } \alpha = 0.1, \text{ so}$$

$$\theta_i := \theta_i - 0.1 \times 2(\theta_i - 3).$$

B: Randomly initialize θ_i and run two iteration of GD

lets $\theta_i = 10$ and $\alpha = 0.1$ and objective function $\min f(\theta_i)$.

$$\theta_i = 10 - 0.1 \times 2(10 - 3).$$

$$= 8.6.$$

$$= 8.6 \times 0.1 \times 2(8.6 - 3)$$

$$= 7.48.$$

lets $\theta_i = -5$ and $\alpha = 0.1$

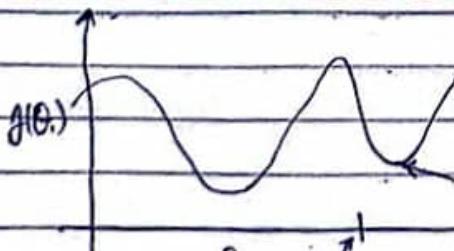
$$\theta_i = -5 - 0.1 \times 2(-5 - 3)$$

$$= -3.4.$$

$$\theta_i = -3.4 \times 0.1 \times 2(-3.4 - 3).$$

$$= -2.12$$

so you see θ_i value approaches 3 and after many iteration θ_i value converged to 3 and $\frac{\partial f(\theta_i)}{\partial \theta_i}$ this factor return 0 because slope at $\theta_i = 3$ is zero.



θ_i at local optima

current value
of θ_i

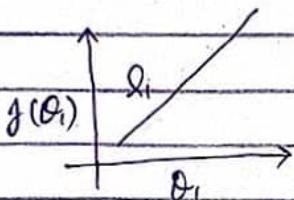
$$\theta_i := \theta_i - \alpha \frac{d}{d\theta_i} f(\theta_i).$$

If you have $f(\theta)$ as plot a: you reach at local optima, and want to update θ_i you can't do that

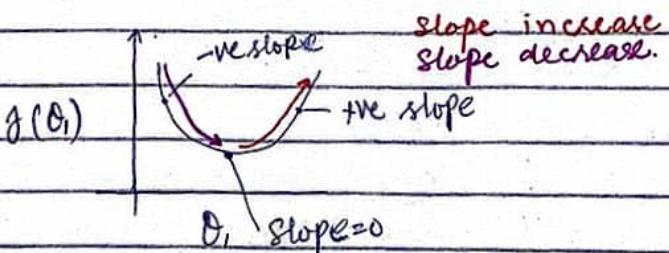
Impact of learning rate in GD (usually in range $(0-1)$)

- If α is too small, GD can be slow. More time and iterations are required to achieve goal.

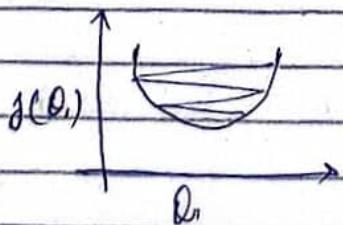
- As slope of line is same at every point and by taking ^{too} small value of α GD slowly approached its goal (minimum value of $J(\theta)$) by taking more time and iterations but approached its goal with same speed.



- As slope of curve is not same at every point and by taking too small value of α GD slowly approached its goal by taking more time and iterations but initially GD ^{try to} approach its goal fastly but becomes slow over the period of time



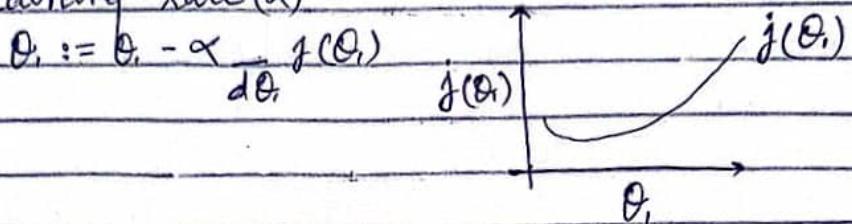
- If α is too large, GD can overshoot the minimum. It may fail to converge, or even diverge.



- if $\alpha = 0$ weights will not update
if $\alpha < 0$ flip the direction of update essentially moving away from the optimal solution.

- α value depend on function. $\rightarrow \alpha$ is hyperparameter.
- eta value should be large.
- eta value should be small.

- GD can converge to a local minimum, even with the fixed learning rate (α)



As we approach a local minimum, GD will take automatically smaller steps due to the impact of slope. so no need to decrease α over time. Therefore we take initially small the value of eta and gradually increase it over the period of time.

- There is some cost function $j(\theta_0, \theta_1)$

- Objective function is to minimize cost function $j(\theta_0, \theta_1)$

$$\min_{\theta_0, \theta_1} j(\theta_0, \theta_1)$$

- we minimize $j(\theta_0, \theta_1)$ by GD.

For using GD to minimize $j(\theta_0, \theta_1)$ following steps should be taken.

① Find weight update equation for every weight.

- Partial derivative of cost function $j(\theta_0, \theta_1)$ w.r.t weights e.g. for θ_0, θ_1

$$\theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} j(\theta_0, \theta_1).$$

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} j(\theta_0, \theta_1).$$

② Randomly initialize weights.

③ Repeat until weights not converged, by using weight update equation.

$$\theta_i := \theta_i - \alpha \frac{\partial}{\partial \theta_i} j(\theta_i)$$

- 3.1 There are two approaches to stop iteration
- ① Check desired result (MSE is 0 or minimum) is achieved.
 - ② Check unknown parameters are converged or not.

④ Simultaneously update weights.

Correct way.

$$\text{temp}_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} j(\theta_0, \theta_1)$$

$$\text{temp}_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} j(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp}_0$$

$$\theta_1 := \text{temp}_1$$

Incorrect way.

$$\text{temp} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} j(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp}$$

$$\text{temp} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} j(\theta_0, \theta_1)$$

$$\theta_1 := \text{temp}$$

Gradient Descent for Linear Regression.

Hypothesis function: tell how to predict values.

Cost function: which unknown parameters are best.

Objective function: tell either minimize or maximize the cost

Approach (GD, Simultaneous equation -) : To find unknown parameters.

Linear Regression.

Hypothesis function $y^{(i)} = h_{\theta}(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$

Hypothesis function for multiple features $y^{(i)} = h_{\theta}(x^{(i)}) = \theta_0 + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \dots$

Cost function = MSE = $j(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$

Objective function = $\min_{\theta_0, \theta_1} j(\theta_0, \theta_1)$

Approach = Gradient Descent.

① weight update equation for each weight.

$$\theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} j(\theta_0, \theta_1); \quad \theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} j(\theta_0, \theta_1).$$

$$\frac{\partial j(\theta_0, \theta_1)}{\partial \theta_0} = \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

$$\frac{\partial j(\theta_0, \theta_1)}{\partial \theta_1} = \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

These already solved.
In linear regression
lecture.

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^i - y^i) = \frac{1}{m} \sum_{i=1}^m (y'^i - y^i)$$

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^i - y^i) x^i = \frac{1}{m} \sum_{i=1}^m (y'^i - y^i) x^i$$

so.

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (y'^i - y^i) \quad \text{and} \quad \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (y'^i - y^i) x^i$$

$\sum_{i=1}^m$ represent all example points that why this is example of Batch Gradient Descent.

① Randomly initialize weights

③ Repeat until weights not converged by using derived weight update equations

④ Simultaneously update weights.

Dry run on training data.

lets $\alpha = 0.1$ $\theta_0 = 3, \theta_1 = 7$ $\theta_0 = 2.6, \theta_1 = 5.4$, $\theta_0 = 2.72, \theta_1 = 5.68$.

x = inputs	y = output	y'	y'	y'
1	10	10	8	8.38
2	15	17	13.4	14.06
3	20	24	18.8	19.74
4	25	31	24.2	25.42
5	30	38	29.6	31.1

Assignment:

Implement Linear Regression by using GD. for
 $\alpha = 0.1, 0.001, 0.01$. and range for random initialization
of weight is [0 to 1].

(1) $E^{(w)} = w_1^2 + w_2^2$ run two iteration of GDA.

assume eta = $\gamma = \alpha = 0.1$, $w_1 = 8$ $w_2 = 10$.

solve for minimization.

$$w_1 := w_1 - \alpha \frac{\partial E(w)}{\partial w_1} \quad \left. \begin{array}{l} \\ \end{array} \right\} \quad w_2 := w_2 - \alpha \frac{\partial E(w)}{\partial w_2}$$

$$w_1 := w_1 - \alpha 2w_1 \quad \left. \begin{array}{l} \\ \end{array} \right\} \quad w_2 := w_2 - \alpha 2w_2$$

$$= 8 - 0.1 \times 2 \times 8 = 10 - 0.1 \times 2 \times 10$$

$$= 6.4 \text{ (after 1st iteration)} = 8 \text{ (after 1st iteration)}$$

$$= 6.4 - 0.1 \times 2 \times 6.4 = 8 - 0.1 \times 2 \times 8$$

$$= 5.12 \text{ (after 2nd iteration)} = 6.4 \text{ (after 2nd iteration)}$$

Derivation of weight update equation by chain rule in Linear Regression

Hypothesis function $h_{\theta}(x) = y' = \theta_0 + \theta_1 x^2$

Cost function = MSE = $j(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (y'^i - y^i)^2$

Objective function $\min_{\theta_0, \theta_1} j(\theta_0, \theta_1)$

chain Rule. $\frac{dy}{dx} = \frac{dy}{dz} \cdot \frac{dz}{dx}$

Required weight update equations

$$\theta_0 := \theta_0 - \eta \nabla \Delta \theta_0 \quad \text{and} \quad \theta_1 := \theta_1 - \eta \nabla \Delta \theta_1$$

$$\Delta \theta_0 = \frac{\partial j(\theta_0, \theta_1)}{\partial \theta_0} \quad \Delta \theta_1 = \frac{\partial j(\theta_0, \theta_1)}{\partial \theta_1}$$

$$= \frac{\partial j(\theta_0, \theta_1)}{\partial y'} \cdot \frac{\partial y'}{\partial \theta_0}$$

$$= \frac{\partial j(\theta_0, \theta_1)}{\partial y'} \cdot \frac{\partial y'}{\partial \theta_1}$$

$$= \frac{\partial j(\theta_0, \theta_1)}{\partial y'} \cdot \frac{1}{m} \sum_{i=1}^m (y'^i - y^i)$$

$$= \frac{\partial j(\theta_0, \theta_1)}{\partial y'} \cdot \frac{1}{m} \sum_{i=1}^m (y'^i - y^i) \cdot \alpha^i$$

so

$$\theta_0 := \theta_0 - \eta \frac{1}{m} \sum_{i=1}^m (y'^i - y^i)$$

$$\theta_1 := \theta_1 - \eta \frac{1}{m} \sum_{i=1}^m (y'^i - y^i) \cdot \alpha^i$$

Simultaneously updates weights.

Reuse these factors values $\frac{1}{m} \sum_{i=1}^m (y'^i - y^i)$, $\frac{1}{m} \sum_{i=1}^m (y'^i - y^i) \cdot x^i$

$$\text{temp}_0 = \frac{1}{m} \sum_{i=1}^m (y'^i - y^i) \quad \text{temp}_1 = \frac{1}{m} \sum_{i=1}^m (y'^i - y^i) \cdot x^i$$

NOW updates weights.

$$\theta_0 := \theta_0 - \gamma \cdot \text{temp}_0 \quad \text{and} \quad \theta_1 := \theta_1 - \gamma \cdot \text{temp}_1$$

Final Exam

- Extensive discussion on hyperparameters, model overfitting, model underfitting, cross validation, information leaking.
- Question on GD
- Pre-mid not included
- Types of attributes topic included in final.
- Softmax related question

Team project

- assignment
- 4 times effort required than
- always split data into training validation and testing
- validation data set for tuning hyperparameter.
- select your topic
- Extensive use of library
- data set and libraries required for project
- NOT used already learned algorithm
- Kaggle, Hugging face.

Classification problems: ① In classification problems output is categorical, countable or discrete. ② we classify data into pre-defined classes.

e.g. binary classification problems, multiclassification problems, multi-level classification problems.

Linear classification problem.

- classes are linearly separable
- classes can be separated by linear decision boundary

Non Linear classification problems.

- classes are not linearly separable

- classes can be separated by non-linear decision boundary

→ we can relate AND gate logic and OR gate logic to a linear classification problem with two input features and binary outputs.

$$x_1 \quad x_2 \quad \text{AND} = y$$

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

$$0 \quad 0 \quad 0$$

→ Range of $\theta_0, \theta_1, \theta_2, x_1$, and x_2 is $(-\infty, +\infty)$

$$0 \quad 1 \quad 0$$

but if we fix x_1 and x_2 in range $[0, 1]$ it will not

$$1 \quad 0 \quad 0$$

affect θ_0, θ_1 , and θ_2 range. $(-\infty, +\infty)$

$$1 \quad 1 \quad 1$$

Threshold :- the level or point at which you start to experiment something.

Defining Threshold: Threshold required for classification problem but not for regression problem because in case of regression output is continuous value while in case of classification output is discrete and we require threshold to classify the data.

Iff for above example threshold is 3.

If $y' \geq 3$

class 1

else

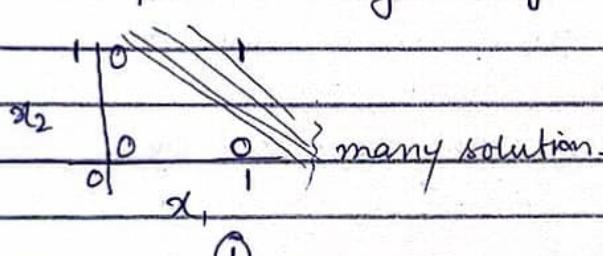
class 0.

So for what values of θ_0 , θ_1 , and θ_2 following equation is valid. $y' = \theta_0 + \theta_1 x_1 + \theta_2 x_2$ and threshold is 3.

$$x_1, x_2 \text{ AND} = y \quad y' = \theta_0 + \theta_1 x_1 + \theta_2 x_2 \quad \text{if } y' \geq 3$$

$$\begin{array}{ccc|c} 0 & 0 & 0 & \theta_0 < 3 \\ 0 & 1 & 0 & \theta_0 + \theta_1 < 3 \\ 1 & 0 & 0 & \theta_0 + \theta_2 < 3 \\ 1 & 1 & 1 & \theta_0 + \theta_1 + \theta_2 \geq 3 \end{array} \quad \begin{array}{l} \text{class 1} \\ \text{else} \\ \text{class 0} \end{array}$$

so there are many solutions as we see when we plot AND gate logic



θ_0	1	0
θ_1	1	1
θ_2	1	2

(1)

(2)

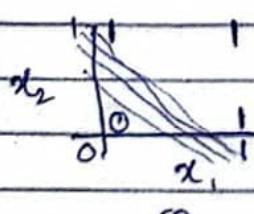
OR Gate

lets Threshold is 7.

$$x_1, x_2 \text{ OR} = y \quad y' = \theta_0 + \theta_1 x_1 + \theta_2 x_2 \quad \text{if } y' \geq 7$$

$$\begin{array}{ccc|c} 0 & 0 & 0 & \theta_0 < 7 \\ 0 & 1 & 1 & \theta_0 + \theta_1 > 7 \\ 1 & 0 & 1 & \theta_0 + \theta_2 > 7 \\ 1 & 1 & 1 & \theta_0 + \theta_1 + \theta_2 > 7 \end{array} \quad \begin{array}{l} \text{class 1} \\ \text{else} \\ \text{class 0} \end{array}$$

$$\begin{array}{ccc|c} 1 & 1 & 1 & \theta_0 & 3 & 5 \\ 0 & 1 & 1 & \theta_1 & 4 & 2 \\ 1 & 0 & 1 & \theta_2 & 5 & 2 \end{array} \quad \begin{array}{l} \text{logic} \\ \text{So there are many solutions. as we plot OR gate} \end{array}$$



θ_0	3	5
θ_1	4	2
θ_2	5	2

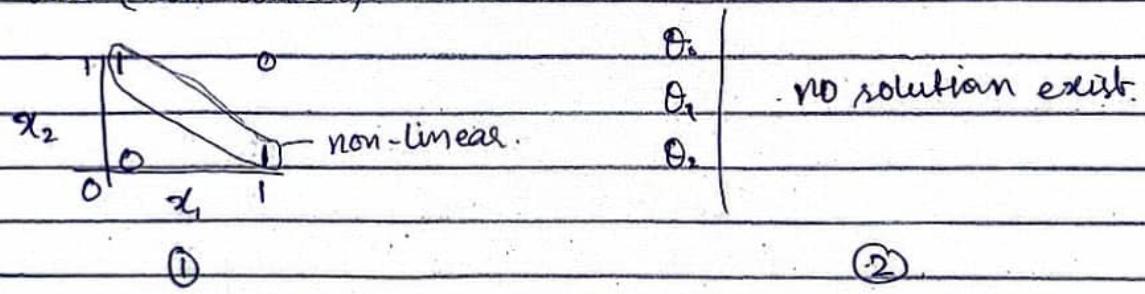
(1)

(2)

We can't relate XOR gate logic to linear classification problem because classes can't be separated by linear decision boundary rather we can relate it with non linear classification problem.

x_1	x_2	$y = \text{XOR}$	$y' = \theta_0 + \theta_1 x_1 + \theta_2 x_2$	lets threshold is 5.
0	0	0	$\theta_0 < 5$	If $y' \geq 5$
0	1	1	$\theta_0 + \theta_2 \geq 5$	class 1
1	0	1	$\theta_0 + \theta_1 \geq 5$	else class 0.
1	1	0	$\theta_0 + \theta_1 + \theta_2 < 5$.	

So there is no linear solution exist as we see when plot XOR gate logic we can't separate classes by single line (linear) but can be separated by oval (non-linear).



→ If there are more than 2 classes (multiclass problem).

- ① 1 class vs all classes approach is used for classification
 - ② Define different thresholds.

Exponential function = $f(x) = a^x$
 $a > 0, a \neq 1$ and a is any real no.
 Natural exponential function = $f(x) = e^x$

Sigmoid Derivation

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

$$\sigma'(x) = \frac{d}{dx} \left(\frac{1}{1+e^{-x}} \right)$$

$$= \frac{d}{dx} (1 + e^{-x})^{-1} = -1 (1 + e^{-x})^{-2} \frac{d}{dx} (1 + e^{-x}) = \frac{-1}{(1 + e^{-x})^2} \left(\frac{d}{dx} (1) + d(e^{-x}) \right)$$

$$= \frac{-1}{(1 + e^{-x})^2} \left(0 + e^{-x} \cdot \frac{d}{dx} (-x) \right) = \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{1 \cdot e^{-x}}{(1 + e^{-x})(1 + e^{-x})}$$

$$= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} = \sigma(x) \cdot \frac{1 + e^{-x} - 1}{1 + e^{-x}} = \sigma(x) \cdot \left(\frac{1 + e^{-x}}{1 + e^{-x}} - \frac{1}{1 + e^{-x}} \right)$$

$$= \boxed{\sigma(x) (1 - \sigma(x))}$$

Binary cross Entropy / Derivation

- For two class problem - cost / loss function for binary classification
 - cost / loss function for logistic regression

$$\text{Cost / loss function: } J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^i \log(y'^i) + (1-y^i) \log(1-y'^i)]$$

Objective function: $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$ or $\min_{\theta} J(\theta)$.

$$\text{Hypothesis function: } h_{\theta}(x) = y'^i = \sigma(\theta_0 + \theta_1 x^i)$$

Let's derive weight update equation

$$\theta_j = \theta_0 \rightarrow \theta_0 + \theta_1 - \frac{\partial}{\partial \theta_j} J(\theta)$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \left(-\frac{1}{m} \sum_{i=1}^m [y^i \log(y'^i) + (1-y^i) \log(1-y'^i)] \right)$$

$$= -\frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \theta_j} [y^i \log(y'^i) + (1-y^i) \log(1-y'^i)]$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[y^i \frac{\partial}{\partial \theta_j} \log(y'^i) + (1-y^i) \frac{\partial}{\partial \theta_j} (1-y'^i) \right]$$

$$\therefore \frac{d}{dx} \log x = \frac{1}{x} \cdot \frac{d}{dx} (x) = \frac{1}{x}$$

$$\therefore \frac{d}{dx} (\log x) = \frac{1}{x} \cdot \frac{1}{\ln 10}$$

$$\text{e.g. } \frac{d}{dx} [\log_{10} (ax^2 + bx + c)] = \frac{1}{ax^2 + bx + c} \cdot \frac{1}{\ln 10} \cdot \frac{d}{dx} (ax^2 + bx + c)$$

$$= \frac{2ax + b}{(ax^2 + bx + c) \ln 10}$$

$$\text{so} \quad = -\frac{1}{m} \sum_{i=1}^m \left[\frac{y^i \frac{\partial}{\partial \theta_j} (y'^i)}{y'^i} + (1-y^i) \frac{\partial}{\partial \theta_j} \frac{(1-y'^i)}{(1-y'^i)} \right]$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[\frac{y^i \cdot y^i \cdot (1-y'^i) \frac{\partial}{\partial \theta_j} \theta^T x}{y'^i} + (1-y^i) (0-(y'^i)) \frac{\partial}{\partial \theta_j} \theta^T x \right]$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[y^i (1-y'^i) x_j^i - (1-y^i) (y'^i) x_j^i \right]$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[y^i - y^i y'^i - y^i + y^i y'^i \right] x_j^i$$

$$= \boxed{-\frac{1}{m} \sum_{i=1}^m (y'^i - y^i) x_j^i}$$

$$\frac{\partial J(\theta)}{\partial \theta_j} \quad \text{Given } J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^i \log(y'^i) + (1-y^i) \log(1-y'^i)]$$

so by chain rule.

$$\frac{\partial J(\theta)}{\partial y^i} \cdot \frac{\partial y^i}{\partial z^i} \cdot \frac{\partial z^i}{\partial \theta_j}$$

$$\min J(\theta)$$

$$y'^i = \sigma(z^i)$$

$$z^i = \theta^T x^i = \theta_0 + \theta_1 x_1^i + \theta_2 x_2^i + \dots$$

- linear regression is line/curve fitting problem.
 - linear regression try to fit curve in the sense as data distributed
 - Decision boundary is drawn by hypothesis function.
- Classification using Gradient Descent.

Logistic Regression (SML)

- Binary classification algorithm
- Implement through supervised machine learning algorithm
- used for binary classification problems
- When we solve classification problems by using some classification algorithm ... ^{model} actually learns ^{its} decision boundary for classification by training data in training phase.
- Logistic Regression based on Sigmoid function

Q: Why we not used linear regression for classification?

Ans: One method of classification is to use linear regression (either implemented by simultaneously equation or by Gradient descent) and map all prediction using some threshold lets say greater than 0.5 as a 1 and less than 0.5 as a 0. But this method doesn't work well because classification is not actually a linear function. So we used sigmoid function to introduce non-linearity. Sigmoid function takes any value between $-\infty$ to $+\infty$ and map between 0 and 1. 0, and 1 both are inclusive.

- Hypothesis function should satisfy: $0 \leq h_{\theta}(x) \leq 1$
So we use sigmoid function also called logistic function to achieve this

$$h_{\theta}(x) = g(\theta^T x) \quad \theta^T x = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots$$

$$z = \theta^T x. \quad \{ \text{range } -\infty \text{ to } +\infty.$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad \text{linear regression.}$$

- we start with our old hypothesis (linear regression), except we want to restrict the range to 0 and 1

- This is accomplished by plugging $(\theta^T x)$ into logistic function.

- $\sigma(h_{\theta}(x))$ will give us the probability that our output is 1.

- for example, $\sigma(h_{\theta}(x)) = 0.7$ gives us the probability of 70% that our output is 1.

$$h_{\theta}(x) = P(Y=1|x;\theta) = 1 - P(Y=0|x;\theta).$$

$$P(Y=1|x;\theta) + P(Y=0|x;\theta) = 1.$$

By default $\sigma(h_{\theta}(x))$ return the probability to belong to class 1

Logistic Regression Decision Boundary.

- Decision boundary may be of any shape (straight line, curve)

- In order to get our discrete 0 and 1 classification, we can translate the output of hypothesis function as follow.

$$h_{\theta}(x) \geq 0.5 \rightarrow Y=1$$

$$h_{\theta}(x) < 0.5 \rightarrow Y=0$$

- The way our logistic function g behaves is that when its input is greater than or equal to zero, its output is greater than or equal to 0.5.

$$g(z) \geq 0.5 \text{ when } z \geq 0$$

$$g(z) < 0.5 \text{ when } z < 0.$$

- So if our input to g is $\theta^T x$ then that means:

$$h_{\theta}(x) = g(\theta^T x) \geq 0.5 \text{ when } \theta^T x \geq 0$$

- from these statement we can now say,

$$\theta^T x \geq 0 \rightarrow Y=1$$

$$\theta^T x < 0 \rightarrow Y=0$$

- The decision boundary is the line that separates the area where $y=0$ and where $y=1$. It is created by our hypothesis function.

lets we learn our unknown parameters by training data. which has 2 input features x_1, x_2 .

$$\theta = \begin{cases} 1 \\ -1 \\ 0 \end{cases} \text{ so } g(z) \text{ if } \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$$

$$S + (-1)x_1 + 0x_2 \geq 0$$

$$S - x_1 \geq 0$$

$$-x_1 \geq -S$$

$$x_1 \leq S \text{ - decision boundary}$$

- In this case our decision boundary is a straight line place on the graph where $x_1=S$ and everything on the left of it that denotes $y=1$ while everything to the right denotes $y=0$.
- Again the input to the sigmoid function $g(z)$ (e.g. $\theta^T x$) doesn't need to be linear and could be a function that describes a circle (e.g. $z = \theta_0 + \theta_1 x_1^2 + \theta_2 x_2^2$) or any shape to fit our data.

Logistic Regression Cost function.

- The more our hypothesis is off from y , the larger the cost function output. If our hypothesis is equal to y , then our cost is 0...
- Rate of change is greater when there is greater difference between actual and predicted mean more.

- we define cost function for both classes.

so our cost function is

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_\theta(x^i), y^i)$$

→ cross entropy.

$$\begin{aligned}\text{cost}(h_\theta(x^i), y^i) &= -\log(h_\theta(x^i)) && \text{if } y^i = 1 \\ \text{cost}(h_\theta(x^i), y^i) &= -\log(1 - h_\theta(x^i)) && \text{if } y^i = 0.\end{aligned}$$

$$\text{cost}(h_\theta(x), y) = 0 \text{ if } h_\theta(x) = y.$$

$$\text{cost}(h_\theta(x), y) \rightarrow \infty \text{ if } y = 1 \text{ and } h_\theta(x) \rightarrow 1$$

$$\text{cost}(h_\theta(x), y) \rightarrow \infty \text{ if } y = 0 \text{ and } h_\theta(x) \rightarrow 0$$

we can compress our cost functions for class 1 and class 0 into one case:

$$\text{cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

Notice that when $y=1$ then second term $-(1-y) \log(1-h_\theta(x))$ will be zero and will not affect the result. If $y=0$ then the first term $-y \log(h_\theta(x))$ will be zero and will not affect the result.

So our full version of cost function is

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^i \log(h_\theta(x^i)) + (1-y^i) \log(1-h_\theta(x^i))]$$

Binary cross entropy.

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^i \log(y^i) + (1-y^i) \log(1-y^i)]$$

Binary cross entropy based built on 0 and 1

descent:- Moving down.

→ GD have such behaviour if automatically converged weight.
→ Sigmoid function basically compute the probabilities

of classes.

→ Cross entropy penalize model more as compared to MSE when prediction deviate from true label.

→ Derivative of cost function penalize model direction with some reference point
-Derivative / Gradient basically tells slope of function (cost function).
-Derivative of cost function w.r.t weight tells either cost function increase or decrease when we increase in weight.

Cost function

$\log(w)$

MSE

- log loss.

w_1

w_2

- MSE

w_1

w_2

weight.

- Derivative / Gradient only tells direction to move.
- GD actually move / push gradient

Logistic Regression (Gradient Descent).

- we use 0.5 as threshold because both classes (0,1) have equal chance to occur mean both classes have same probability (0.5) and sum of probability is 1

steps

- ① Randomly initialize weights between (0 and 1)
- ② calculate prediction.

$$\text{calculate } z = \theta^T x$$

$$\text{calculate } \sigma(z)$$

1 if 0.5 class to input after thresholding .

calculate accuracy

- ③ Repeat until desired accuracy achieved using weight update equations $\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$

$$\text{Hypothesis function } y'^2 = \log(x^i) = \sigma(z) \quad \therefore z = \theta^T x \\ \therefore \sigma(z) = \frac{1}{1+e^{-z}}$$

Cost function: Binary cross entropy

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^i \log(y^i) + (1-y^i) \log(1-y^i)]$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (y'^i - y^i) x_j^i$$

so weight update equation is

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (y'^i - y^i) x_j^i$$

By default vector defines cols wide.

Logistic Regression, vectorized implementation.

- For calculating prediction:

$$h_{\theta}(x) = g(x \cdot \theta)$$

and after thresholding we can find y' :

where X is data set with additional input features.

that whose all values equal to 1

and θ is vector of weights.

- For calculating cost

$$J(\theta) = \frac{1}{m} * -y^T \log(h) - (1-y)^T \log(1-h)$$

- For update weights:

$$\theta := \theta - \alpha \nabla J(\theta)$$

$$\nabla J(\theta) = \frac{1}{m} \cdot X^T \cdot (g(X \cdot \theta) - y)$$

not use threshold
in weight update

Assignment $\theta := \theta - \alpha \cdot X^T \cdot (g(X \cdot \theta) - y)$ | went to calculate accuracy

Assignment

learning rate for assignment 0.01, 0.1, 0.3, 0.5

stopping criterion:

→ 100 iteration print accuracy in each iteration

* during training

→ accuracy on training as well as on testing data.

→ map 2 and 4 class to 0 and 1.

stopping criteria

- is desired accuracy achieved or not -

- No. of iterations

- losses / weights converged or not not suitable

- minimum loss

not suitable

K-Medoids clustering Algorithm. (USML)

- also called as partitioning around medoid (PAM)
 - A medoid can be defined as the point in the cluster, whose dissimilarities with all the other points in the cluster is minimum.
 - The dissimilarity of the medoid C_i and object or point P_i is calculated by using $E = |P_i - C_i|$
- Note that other distance measures can be used to calculate dissimilarity.
- The cost in K-Medoids algorithm is given as

$$C = \sum_{C_i} \sum_{P_i \in C_i} |P_i - C_i|$$


$C = \text{cost (a numerical number)} P_i$

Algorithm working

- Initialization of K-random points.
- Select K random points out of the N data points as points the medoid. Medoids must be selected among the data.
- Compute the closeness of each data point w.r.t each medoid. Associate each data point to the closest medoid using any common distance metric methods.
- While the cost decreases

 - for each medoid m, for each data point o which is not a medoid
 - (1) Swap m and o, associate each data points to the closest medoid, and recompute the cost
 - (2) If the total cost is more than that in the previous step, undo the swap and stop

cost tells us how better clusters are made by these medoids

K-mean is best for numeric data
K-medoid is best for categorical data.

Day Run.

Let we have input data that have two input features say X and Y and total data points are equal to 10 and $\text{I}^2 = 2$ (C_1 and C_2) with $C_1 = (3, 4)$; $C_2 = (7, 4)$

Dissimilarity from C_1

Dissimilarity from C_2

S.No.	X	Y	Dissimilarity from C_1	Dissimilarity from C_2
0	7	6	2	6
1	2	6	3	7
2	3	8	4	8
3	8	5	6	2
4	7	4	4	0
5	4	7	4	6
6	6	2	5	3
7	7	3	5	1
8	6	4	3	1
9	3	4	0	4
			11	7
				1
				9
				1

- So each point is assigned to the cluster of user 8

medoid whose dissimilarity is less

- As we see points 1, 2 and 5 go to cluster C_1 and 0, 3, 6, 7

go to cluster C_2 .

- The cost $C = (3+4+4) + (2+2+3+1+1) = 11+9 = 20$.

- Now randomly select one non-medoid point and recalculate the cost. Let randomly selected point be (7, 3). As this point has minimum dissimilarity to C_2 by C_2 .

S.No.	X	Y	Dissimilarity from C_1	Dissimilarity from C_2
0	1	6	6	3
1	6	3	3	8
2	3	8	4	6
3	5	6	9	3
4	4	4	7	1
5	2	5	2	4
6	3	2	1	2
7	4	3	5	4
8	3	4	2	5
9	0	3	1	1
			11	6
				2/11

- Each data point is assigned to next cluster whose dissimilarity is less.

- As we see data points 1, 2 and 5 go to cluster C_1 and 0, 3, 6, 7 and 8 go to cluster C_2 .
- The cost $C = (3+4+4) + (3+3+1+2+2) = 11 + 11 = 22$.

$$\text{Sweep cost} = \text{present cost} - \text{previous cost}$$
$$= 22 - 20$$
$$= 2 > 0$$

As sweep cost is not less than 0, we undo the sweep. Hence (3,4) and (7,4) are the final medoids. The clustering would be in the following way.

Advantages:

- It is simple to understand and easy to implement.
- K-medoid Algorithm is fast and converges in a fixed no. of steps
- PAM is less sensitive to outliers than other partitioning algorithms.

- The main disadvantage of K-medoid algorithm is that it is not suitable for clustering non-spherical (arbitrary shaped) groups of data points / objects. This is because it relies on minimizing the distance between the non-medoid objects and the medoid - briefly, it uses compactness as clustering criterion instead of connectivity.

Q:- what is meant by algorithm is less sensitive to outliers? It means algorithm is more robust to outliers and it is less likely to produce extreme or misleading results due to the presence of outliers in the data.

K-Medoid Assignment (Information Retrieved System)

Steps

- ① - Assign id's to each document.
- ② - Extract unique words from all docs.
By typecasting in set
- ③ - Make vocabulary of unique words and assign id's
By dictionary.

④ - Build Term Doc Matrix (TDM)

TDM dimension's = no. of docs x vocabulary size.

- ⑤ - Generate query and make query vector
- ⑥ - Compute matching score of query and docs
- ⑦ - Retrieve docs on the basis of scores

My Run

① docs = [{"I live in Lahore", "Lahore city", "Pakistan Lahore live", "Lahore is city of Pakistan"}]

② unique words = {"I", "live", "in", "Lahore", "city", "Pakistan", "is", "af", "T", "O", "live": 1, "in": 2, "Lahore": 3, "city": 4, "Pakistan": 5, "is": 6, "af": 7}

③ vocabulary = {T: 0,

	0	1	2	3	4	5	6	7
0	1	1	1	1	0	0	0	0
1	0	0	0	1	1	0	0	0
2	0	1	0	1	0	1	0	0
3	0	0	0	1	1	1	1	1

④ query = Pakistan city

0	0	0	0	1	1	0	0
---	---	---	---	---	---	---	---

⑤ matching score = [0, 1, 1, 2]

⑥ results

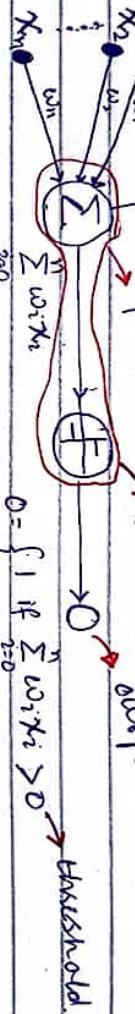
doc-4, doc-2, doc-3

Artificial Neural Network

It is a powerful learning algorithm inspired by how the brain works.

- Perception is the basic unit of artificial neural network.
- we use "bias" during learning of artificial neural network so it gave better results on unseen data.

$x_0 = 1$ perception. Activation function (step function).



$$0 = \begin{cases} 1 & \text{if } \sum_{i=0}^n w_i x_i > 0 \\ -1 & \text{otherwise} \end{cases}$$

$$0(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n > 0. \\ -1 & \text{otherwise} \end{cases}$$

$$\mathcal{T} = b + \sum_{i=1}^n w_i x_i$$

perceptron

$Z = w \cdot x + b$. (vectorize implementation).

Here $-x_1, x_2, \dots, x_n$ are input features

- b = bias = $x_0 w_0 = w_0 \because x_0 = 1$.
- Assign class 1 or -1 after thresholding
- w_0, w_1, \dots, w_n are respective weights of input features
- step function compare output with threshold.

So in this way we build a perceptron for two class problem.

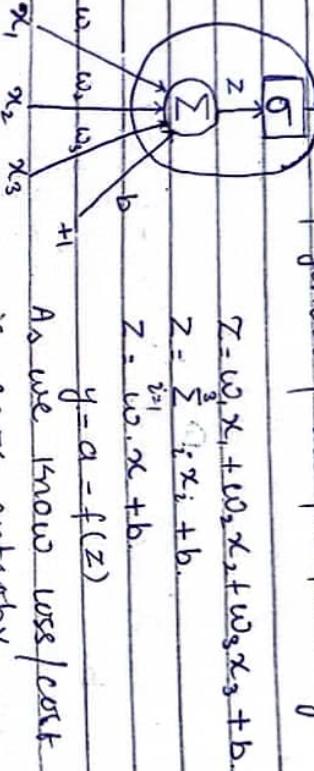
Logistic Regression as perceptron.

- Simply use sigmoid unit function as activation function to introduce non-linearity.

a generally use for hidden layer.

$$y = a = f(z)$$

y generally use for final layer



$$y = a = f(z)$$

As we know loss/cost function for LR
is cross entropy.

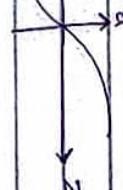
$$L_{ce}(\hat{y}, y) = -\log P(y|x) = -[y \log \hat{y} + (1-y) \log(1-\hat{y})]$$

Different Activation function a

$$\textcircled{1} \text{ Sigmoid: } a = \frac{1}{1+e^{-z}}$$



$$\textcircled{2} \text{ Tanh: } a = e^z - e^{-z} / e^z + e^{-z}$$



$$\textcircled{3} \text{ Relu: } a = \max(0, z)$$



$$\textcircled{4} \text{ Leaky Relu: } a = \max(0.1z, z)$$



perception \approx Neuron.

- Basic unit.

- No hidden layer

XOR using Artificial Neural Network

AND: First we implement AND by Artificial Neural

$x_1 \quad x_2 \quad y$

$$\begin{array}{ccc} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \end{array} \quad \begin{array}{c} x_1 \cdot w_1 \\ x_2 \cdot w_2 \\ b \end{array} \quad \sum_{i=1}^2 w_i x_i \rightarrow y = b + \sum_{i=1}^2 w_i x_i$$

$$y = b + w_1 x_1 + w_2 x_2$$

$$y = \begin{cases} 0 & \text{if } b + w \cdot x \leq 0 \\ 1 & \text{if } b + w \cdot x > 0 \end{cases}$$

perceptron.

- So required evaluations one.

AND

$x_1 \quad x_2 \quad y$	Required Evaluations	Simplified	Solutions
0 0 0	$b + w_1 x_1 + w_2 x_2 \leq 0$	$b \leq 0$	$b: -0.5$
0 1 0	$b + w_1 x_1 + w_2 x_2 \leq 0$	$b + w_2 \leq 0$	$w_2: 0.4$
1 0 0	$b + w_1 x_1 + w_2 x_2 \leq 0$	$b + w_1 \leq 0$	$w_1: 0.4$
1 1 1	$b + w_1 x_1 + w_2 x_2 > 0$	$b + w_1 + w_2 \geq 0$	0.3

- Plotting AND.

~~perceptron output~~
As you see there infinitely solutions exist.

- Decision boundary.

We can find decision boundary by $y = b + w_1 x_1 + w_2 x_2$.

$$b + w_1 x_1 + w_2 x_2 = 0$$

$$w_1 x_1 + w_2 x_2 = -b$$

$$x_2 = -\frac{b}{w_2} - \frac{w_1}{w_2} x_1 \text{ decision boundary.}$$

This line act as a decision boundary in 2D space.

HTM Theory

in which the output 0 is assigned to all inputs on one side of the line, and output 1 to all input points lying on the other side of the line.

Similarly we can implement OR by one perceptron but we can't implement XOR by one perceptron but with more perceptrons.

Alternative approach for XOR

XOR

x_1	x_2	y_i	h_i	$x_1 \wedge \neg x_2$	$h_1 \vee h_2$
0	0	0	0	0	0
0	1	1	0	1	1
1	0	1	1	0	1
1	1	0	0	0	0

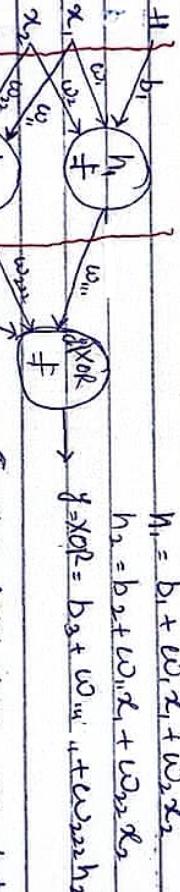
$x_1 \oplus x_2$

$x_1 \oplus x_2$

$$h_1 = b_1 + w_{11}x_1 + w_{21}x_2$$

$$h_2 = b_2 + w_{12}x_1 + w_{22}x_2$$

$$y = \text{XOR} = b_3 + w_{31}h_1 + w_{32}h_2$$



- In this ANN there are total 3

perceptrons $\binom{3}{2}$ in hidden layers

Input Layer Hidden Layer Output Layer
 - They are also called ReLU units

- True are total 9 unknown weights.
- Input, hidden and output layers may have multiple perceptrons
- An Artificial Neural Network must have input and output layers but may have 0 or more hidden layers.

$x_1 \quad x_2 \quad y = \text{XOR}$

0	0	0
0	1	1
1	0	1
1	1	0

The original x space.

$h_1 = x_1 \Delta \tau x_2, h_2 = \tau x_1 \Delta x_2, y = \text{XOR} = h_1 \vee h_2$

0	0	0
0	1	1
1	0	1
0	0	0

0	0	0
0	1	1
1	0	1
0	0	0

0	1	1
0	0	0

The new h space.

Feed forward neural network

- previous layer is fully connected with next layer

x_1

h_1

h_2

y

W_{ij} of the weight matrix w

represents the weight of
the connection from ith
input x_i to the jth hidden
unit h_j .

- Efficient implementation
by using weight matrix.

= Multiply w by x

+ b

- Add bias b.

Input layer Hidden layer Output layer

w and v are weights matrix

Similarly for next layer.

$$\begin{aligned} z &= v^T h \\ z &= \text{softmax}(z) \end{aligned}$$

softmax

- To compute $P(Y=c|x)$ where classes are more than two
- softmax takes a logits vector (z) of c (no of classes) arbitrary values and maps them to probability distribution with each value in the range (0 to 1) and all values summing to 1
- $z = [z_1, z_2, z_3, \dots, z_c]$

$$\rightarrow \text{Softmax}(z) = \left[\frac{e^{z_1}}{\sum_{k=1}^c e^{z_k}}, \frac{e^{z_2}}{\sum_{k=1}^c e^{z_k}}, \frac{e^{z_3}}{\sum_{k=1}^c e^{z_k}}, \dots, \frac{e^{z_c}}{\sum_{k=1}^c e^{z_k}} \right]$$

For Example. --

$$z = [0.6, 1.1, -1.5, 1.2, 3.2, -1.1]$$

$$\exp(z) = e^z = [1.822, 3.004, 0.223, 3.320, 24.532, 0.332]$$

$$\text{Sum of } (e^z) = 33.234$$

$$\text{softmax}(z) = [0.054, 0.090, 0.006, 0.099, 0.738, 0.010]$$

→ Softmax function generates probability distribution.