

1. CONSTRAINT BASED SYSTEMS

↳ user specifies his initial preference

↳ user is displayed a set of matching items with reason why an item was recommended

↳ 2 sets of constraints

↳ User

↳ Customer

↳ uses attributes

✓ constraints
↓ objectives

X

Components of Constraint-based Systems

- Customer properties
- Product properties
- Constraints
- Filter conditions and product
- Products

V_c
 V_{prod}
 C_p (on customer properties)
 C_r - relationship between customer
 C_{prod} - possible instantiations

$V_c =$ { $v_{c1} : [vegan, no_dairy, no_requirements] \dots /* dietary requirements */$ $v_{c2} : [fast_food, traditional, continental, none] \dots /* food preference */$ $v_{c3} : [high, medium, low] \dots /* spending style */$ } ... }	$V_{prod} =$ { $v_{p1} : [pizza, BBQ, burger, biryani, italian] \dots /* type of food */$ $v_{p2} : [300, \dots, 4000] \dots /* food price */$ $v_{p3} : [plaintext] \dots /* name of food */$... }	$C_{prod} =$ { $c_{prod1} : v_{p3} = afghani\ tikka\ pizza \wedge v_{p2} = 1300$ $c_{prod2} : v_{p3} = chicken\ burger \wedge v_{p2} = 750$ $c_{prod3} : v_{p3} = chicken\ biryani \wedge v_{p2} = 400$ $c_{prod4} : v_{p3} = alfredo\ pasta \wedge v_{p2} = 1200$... } less than 1000
$C_R =$ { $c_{R1} : v_{c1} = vegan \rightarrow v_{c2} \neq fast\ food$ $c_{R2} : v_{c2} = continental \rightarrow v_{c3} \neq low$... }	$C_f =$ { $c_{f1} : v_{c1} = vegan \rightarrow v_{p1} \neq BBQ \wedge v_{p1} \neq burger$ $c_{f2} : v_{c3} = low \rightarrow v_{p2} < 1000$... }	

if user preferences = No dairy, none, low

recommended items: chicken burger, chicken biryani

CONS

↳ can get 0 recommendations as customer given req. might be such that no product shown

Solution

1. make a default product ↗

2. relax constraints ↗ 3 violations recommended

↳ customer might be unsure about the values of their attributes

↳ business rules change frequently ↗ will have to hire for every change → EXPENSIVE

↳ no serendipity

↳ optimization algos are slow

↳ domain expertise needed

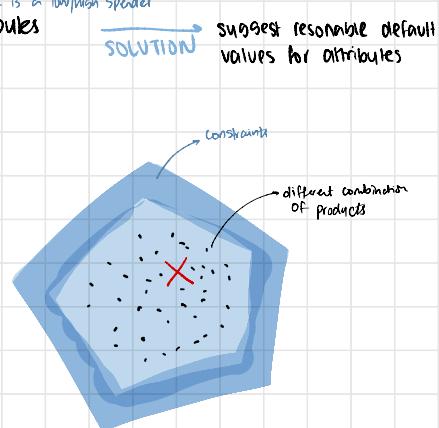
It's a Subset Selection Problem

Different techniques to find subsets ↗ to optimise

1. QuickPlan

2. Linear Programming

3. Gradient Descent ↗ terrible idea



2. CASE BASED

↳ uses attributes

↳ items are retrieved using similarity measures

can't be sparse
↓
should have all info
↳ unlike collaborative filtering



Q) Allowed attributes:

$$P = \{12, 25, 55\}$$

$$M = \{8, 10, 48, 54, 60\}$$

$$L = \{3, 4, 5, 6\}$$

$$W = \{\text{Yes, No}\}$$

weights attributes:

$$W_{\text{Price}} = 1.5$$

$$W_{\text{megapixel}} = 1$$

$$W_{\text{LCD}} = 1$$

$$\text{sim}(p, \text{REQ}) = \frac{\sum_{r \in \text{REQ}} w_r \times \text{sim}(p, r)}{\sum_{r \in \text{REQ}} w_r}$$

weight for requirement

* if this data not given
get max. value from table

↳ CASES

camera	Price	Megapixel	LCD	Waterproof
C1	12	8	3	No
C2	25	48	4	No
C3	55	60	6	Yes

User comes with
this query

NORMALISE

camera	Price	Megapixel	LCD	Waterproof
C1	0.218	0.133	0.428	0
C2	0.454	0.8	0.511	0
C3	1	1	0.857	1

↳ USER req

Price	Megapixel	LCD	Waterproof
20	40	5	No

NORMALISE

Price	Megapixel	LCD	Waterproof
0.363	0.666	0.714	0

↳ Find Similarities

user-price case-price

$$\text{Sim}(p_e, p_{c1}) = 1 - |0.363 - 0.218| = 0.855$$

$$\text{Sim}(M_e, M_{c2}) = 1 - |0.666 - 0.133| = 0.667$$

$$\text{Sim}(L_e, L_{c3}) = 1 - |0.714 - 0.428| = 0.714$$

User
req
cases

↳ Find Similarity of user requirements with all

$$\text{Sim}(\text{User}_{\text{REQ}}, C1) = \frac{1(0.855) + 1(0.667) + 1(0.714)}{1+1+1} = 0.703$$

Ans is highest score

$$\text{Sim}(\text{User}_{\text{REQ}}, C2) =$$

$$\text{Sim}(\text{User}_{\text{REQ}}, C3) =$$



PRO

- ↳ no O recommendations
- ↳ has no constraints
so no domain expert needed → ?
- ↳ low cost technique

CON

- ↳ attributes aren't always independent in reality
but we assume they are

Knowledge-based Systems: **Limitations**

- ❑ Cost of knowledge acquisition and domain expertise
- ❑ Attributes (preferences) are not always independent

- a) Consider the given table depicting the features of laptops available for purchase.

	HD Display	RAM (GB)	Storage (GB)	Price (k)
L1	Yes	4	512	52
L2	No	8	1020	80
L3	Yes	16	1020	96

Use case-based knowledge driven technique to recommend the most suitable laptop for the following user requirements:

HD Display	RAM (GB)	Storage (GB)	Price (k)
No	16	512	60

Note: Assume all attribute weights to be 1.

- b) List down some limitations of constraint-based systems.

Normalise cases

	HD	RAM	STO	PRI
L1	1	0.25	0.502	0.542
L2	0	0.5	1	0.833
L3	1	1	1	1

Normalise User

	HD	RAM	STO	PRI
User	0	1	0.502	0.625

Find Similarities

With L1

$$HD = 1 - |0 - 1| = 0$$

$$RAM = 1 - |1 - 0.25| = 0.25$$

$$STO = 1 - |0.502 - 0.502| = 1$$

$$PRI = 1 - |0.625 - 0.542| = 0.917$$

$$\text{Similarity} = \frac{0 + 0.25 + 1 + 0.917}{1 + 1 + 1 + 1} = 0.542$$

With L2

$$HD = 1 - |0 - 0| = 1$$

$$RAM = 1 - |1 - 0.5| = 0.5$$

$$STO = 1 - |0.502 - 1| = 0.502$$

$$PRI = 1 - |0.625 - 0.833| = 0.792$$

$$\text{Similarity} = \frac{1 + 0.5 + 0.502 + 0.792}{1 + 1 + 1 + 1} = 0.68$$

will be recommended

With L3

$$HD = 1 - |0 - 1| = 1$$

$$RAM = 1 - |1 - 1| = 1$$

$$STO = 1 - |0.502 - 1| = 0.502$$

$$PRI = 1 - |0.625 - 1| = 0.625$$

$$\text{Similarity} = \frac{1 + 1 + 0.502 + 0.625}{1 + 1 + 1 + 1} = 0.531$$

b) Limitations of Constraint-based Recommender Systems:

- It may generate zero recommendations in some cases.
- Users are not always sure of what their full requirements are.
- If business rules often change, updates to the knowledge-base and constraints need to be made.
- Domain expertise is required.

Consider the given table depicting the features of processors available for purchase.

	Cores	Price (in thousands)	L3 Cache (in MB)
CPU 1	6	69	32
CPU 2	12	78	32
CPU 3	8	74	64

Use case-based knowledge driven technique to recommend the most suitable processor for the following user requirements.

Feature	Cores	Price (in thousands)	L3 Cache (in MB)
Requirement	12	68	64
Weight	1.25	1.40	1.12

Normalised cases

	Cores	Price	L3 cache
CPU 1	0.5	0.885	0.5
CPU 2	1	1	0.5
CPU 3	0.667	0.949	1

Normalised user

	Cores	Price	L3 cache
Requirement	1	0.872	1
Weight	0.104	0.018	0.0175

sum of weight: 0.1395

Find Similarity

CPU 1 Requirements

$$\text{Sim}(R_C, C) = 1 - |1 - 0.5| = 0.5$$

$$\text{Sim}(R_P, P) = 1 - |0.872 - 0.885| = 0.987$$

$$\text{Sim}(R_{L3}, L3) = 1 - |1 - 0.5| = 0.5$$

$$\text{Sim}(\text{CPU1}) = \frac{(0.104)(0.5) + (0.018)(0.987) + (0.0175)(0.5)}{0.1395}$$

$$= 0.563$$

CPU 2

$$\text{Sim}(R_C, C) = 1 - |1 - 1| = 0.872$$

$$\text{Sim}(\text{CPU2}) = \frac{(0.104)(0.872) + (0.018)(0.872) + (0.0175)(0.5)}{0.1395}$$

$$\text{Sim}(R_P, P) = 1 - |0.872 - 1| = 0.872$$

$$0.1395$$

$$\text{Sim}(R_{L3}, L3) = 1 - |1 - 0.5| =$$

$$= 0.825 \rightarrow \text{will be recommended}$$

CPU 3

$$\text{Sim}(R_C, C) = 1 - |1 - 0.667| = 0.667$$

$$\text{Sim}(\text{CPU3}) = \frac{(0.104)(0.667) + (0.018)(0.949) + (0.0175)(1)}{0.1395}$$

$$\text{Sim}(R_P, P) = 1 - |0.872 - 0.949| =$$

$$0.1395$$

$$\text{Sim}(R_{L3}, L3) = 1 - |1 - 1| =$$

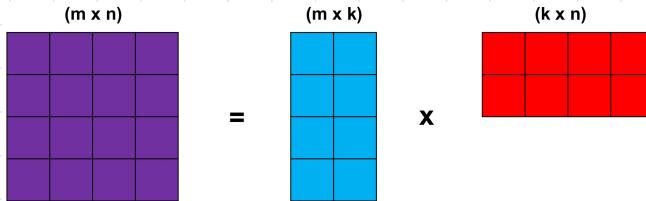
$$= 0.742$$

Matrix Factorization

Focuses
on
STORAGE

↳ decomposes interaction matrix into a product of

2 or 3 matrices → you can have any number of factors



$$4 \times 4 = 16$$

$$5 \times 5 = 25$$

$$100 \times 100 = 10,000$$

$$4 \times 2 = 8$$

$$5 \times 2 = 10$$

$$100 \times 5 = 500$$

$$4 \times 2 = 8$$

$$2 \times 5 = 10$$

$$\times 100 \times 5 = 500$$

$$= 16$$

$$= 20$$

$$= 1000$$

SAVING

$$16 - 16 = 0$$

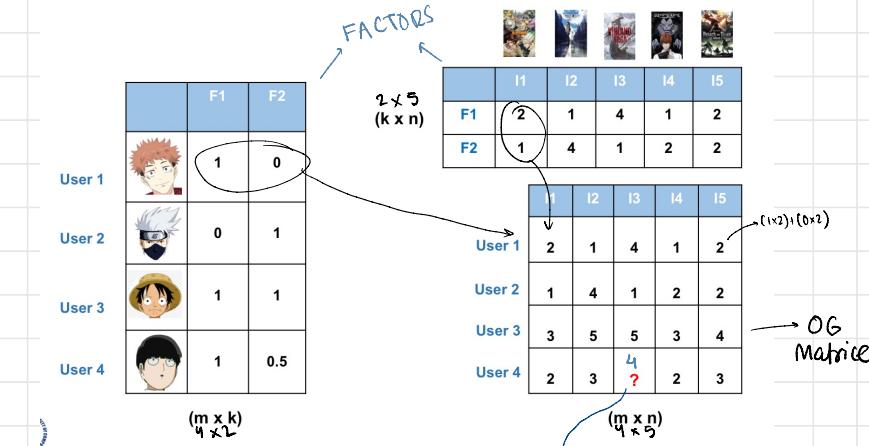
$$25 - 20 = 5$$

$$10,000 - 1000 = 9,000$$

m : rows

n : columns

Matrix Factorization



Con

↳ struggles with sparsity

↳ suffers from cold start problem

Finding Factors

1. Gradient Descent
2. Genetic Algorithm
3. Linear Programming
4. PSO

	I	II	III	IV	V	
(k x n)	F1	2	1	4	1	2
	F2	1	4	1	2	2
User 1		1				
User 2		0	1			
User 3		1	1			
User 4		1	0.5			

(m x k)

Issue: Where do these values come from? In other words, how do we find these two matrices (factors)?

(CS-4950 Course Material 7)

1. Gradient Descent

- ↳ Quadratic

Matrix Factorization: Gradient Descent

	F1	F2						
User 1	0.2	0.5						
User 2	0.4	0.6						
User 3	0.1	0.05						
User 4	0.9	0.3						

Start with a rough guess

$$0.2 \times 0.1 + 0.5 \times 0.8 = 0.42$$

(Do the same to fill all cells)

We are "off" by 1.58

Error = $(2 - 0.42)^2 + (1 - 1.85)^2 + \dots$

↓ Repeat

(CS-4950 Course Material 7)

Steps

- ↳ Initially put random values

- ↳ Error function

1. Sum of errors → Used in diagram

2. Mean square error

- ↳ Try to make overall error close to 0 → Local Optima

- ↳ Ignore anomalies

- ↳ Keep repeating till error minimised



	F1	F2						
User 1	0.14	0.45						
User 2	0.32	0.5						
User 3	0.43	0.07						
User 4	0.8	0.3						

Let's try this again with different values

$$0.14 \times 0.2 + 0.45 \times 1.1 = 0.52$$

Keep repeating until error can be minimized

	I	II	III	IV	V
F1	0.2	1.7	0.8	3.7	1.9
F2	1.1	3.3	3.2	4.0	0.5
User 1	0.52	1.72	1	3	2
User 2	1	4	1	2	2
User 3	3	5	3	4	4
User 4	2	3	2	2	3

PROS

- ↳ SAVES STORAGE
- ↳ Quick and easy Predictions

CONS

- ↳ very slow
 - ↳ non deterministic
 - ↳ expensive
 - ↳ time consuming
- ↓ Solution
- use population algos
- ↳ genetic algo

cost guess
how many iterations

Consider the given $m \times n$ interaction matrix:

5	3	5	3	3	1
	2		3	3	4
5	3	1		5	
3		2	1		2

→ user 1

- a) Which one of the following values of k for the given factors $m \times k$ and $k \times n$, will be providing the most efficient storage of interactions:

$$K = 1, 2, 4$$

- b) Populate the factors 4×2 and 2×6 with some arbitrary values and calculate the mean squared error for User 1 (represented by the first row of the original matrix).

a) $m=4, n=6 \rightarrow 4 \times 6 = 24$

$$k=1 \rightarrow 4 \times 1 + 1 \times 6 = 10 \rightarrow 24 - 10 = 14 \rightarrow \text{ANS}$$

$$k=2 \rightarrow 4 \times 2 + 2 \times 6 = 20 \rightarrow 24 - 20 = 4$$

$$k=4 \rightarrow 4 \times 4 + 4 \times 6 = 40 \rightarrow 24 - 40 = 16$$

↓ asked most efficient
not which saves
most storage

→ smaller K
less complexity

b)

1	0
3	2
1	3
2	0

0	2	4	1	0	0
1	1	3	0	2	0

0	2	4	1	0	0
2	8	18	3	4	0
3	5	13	1	6	0
0	4	8	2	0	0

→ user 1

$$\begin{aligned} \text{Sum of squared error} &= (5-0)^2 + (3-2)^2 + (5-4)^2 + (3-1)^2 + (3-0)^2 + (1-0)^2 + \\ &= \text{Khud plus karlo} \\ &= \text{i am mak gaya } \Sigma \end{aligned}$$

For the following 3×5 interaction matrix, answer the questions given below:

	Item 1	Item 2	Item 3	Item 4	Item 5
User 1	3	1	1	3	1
User 2	1	2	?	1	3
User 3	4	3	5	4	4

- a) Do these two matrix factors correctly represent the original interaction matrix? Provide a reason.

	F1	F2
U1	1	0
U2	0	1
U3	0	1

	I1	I2	I3	I4	I5
F1	3	1	1	3	1
F2	1	2	4	1	3

- b) Use the matrix factors provided in Part (a) to predict missing rating $R(U2, I3)$.
 c) How much storage space does the original interaction matrix require? Compare it with the total storage space required by the two given factors and write your comments about it.
 d) Re-initialize the matrix factors in Part (a) with different arbitrary values in range [0.1—0.9] and calculate the sum of all squared errors.

- a) **Answer:** No. The factor for User 3 does not truly represent the original interaction matrix e.g., for Item 1, the dot product $U3 \cdot I1$ gives us $(0 \times 3 + 1 \times 1) = 1$ instead of the original rating 4.
- b) **Answer:** $(0 \times 1 + 1 \times 4) = 4$.
- c) **Answer:** The original interaction matrix would take $3 \times 5 = 15$ units of storage, while the two factors would take $(3 \times 2) + (2 \times 5) = 16$ units of storage. Although, for this particular example the original interaction matrix is taking up less storage than its factors, when the size of the matrix is large, there will be a significant decrease in the storage requirements for the factored matrices compared to the original matrix.
- d) Initialized matrix:

	Item 1	Item 2	Item 3	Item 4	Item 5
User 1	0.9	0.7	0.4	0.8	0.3
User 2	0.5	0.6	0.2	0.1	0.7
User 3	0.5	0.6	0.2	0.1	0.7

Sum of Squared Error:

$$(3-0.9)^2 + (1-0.7)^2 + (1-0.4)^2 + (3-0.8)^2 + (1-0.3)^2 + (1-0.5)^2 + (2-0.6)^2 + (1-0.1)^2 + (3-0.7)^2 + (4-0.5)^2 + (3-0.6)^2 + (5-0.2)^2 + (4-0.1)^2 + (4-0.7)^2 = 85.56$$

Question 3: Consider the given 7x6 interaction matrix R and its factors U and V .

5	5	5	0	0	0
5	4	5	0	0	0
1	5	2	0	0	0
5	5	3	5	5	5
-1	-4	-5	1	1	1
-2	-1	-1	1	1	1
-3	-5	-1	1	1	1

5	1
-2	1
-4	0
1	0
0	1
0	0
0	0

1	-3	-2	-5	-1	-5
3	-2	0	1	0	-6

Answer the following questions:

- If r and c represent rows and columns, do the factors contain valid values for original interactions $R(r1, c2)$ and $R(r3, c3)$?
- The given system allows negative ratings [-5 to -1], no interaction [0] and positive ratings [1 to 5]. What is the major redundancy in this approach?
- Will using gradient descent algorithm with a learning rate of 0.0005 for optimizing the given factors be a suitable approach? Justify your answer with a reason.
- Name one major disadvantage of Matrix Factorization technique.

a) $R(r1, c2) = 5 \rightarrow (5, 1), (-3, -2)$
 $\therefore (5)(-3) + (1)(-2)$
 $= 47$
 False

$R(r3, c3) = 2 \rightarrow (-4, 0), (-2, 0)$
 $\therefore (-4)(-2) + (0)(0)$
 $= -8$
 False

- Major redundancy is attempting to capture all interaction types with only 2 factors, potentially overfitting
- a learning rate of 0.0005 is very small, so slow convergence but it also helps avoid large updates, that leads to instabilities
 so suitable approach
- struggles with sparsity when there are few observed interactions

Consider the following matrix for anime recommendation where 1 and -1 represent *like* and *dislike*, respectively. Answer the given questions:

	Bleach	Berserk	Erased	Akira
User 1	1	?	1	-1
User 2	-1	1	1	1
User 3	?	?	-1	-1
User 4	-1	1	1	?

- a) Is it possible to make more than three factors of this matrix? Defend your answer with a proof.
- b) For what values of m , k , l and n , will the three factors $(m \times k)$, $(k \times l)$ and $(l \times n)$ consume less memory in total than the original interaction matrix?
- c) We used interactions instead of ratings in the given matrix. Can this approach help gradient descent converge quicker?
- d) Does matrix factorization suffer from cold-start problem as severely as Collaborative Filtering? Justify your answer.

a) Yes, as we can make any amount of factors

b) $(m \times k)$ $(k \times 1)$ $(1 \times n)$
 (4×2) (2×1) (1×4)

No matrix: $4 \times 4 = 16$

factors: $4 \times 2 + 2 \times 1 + 1 \times 4 = 14$

c) Yes, as interactions are less complex
than a range of ratings

d) Yes, as if there's no training data
no factors can be generated

NEURAL NETWORKS

Derivatives

↳ rate of change wrt to something

Q1

$$f(u) = 2u$$

$$f'(u) = 2$$

Gradient

↳ vector of partial derivatives

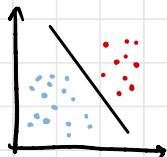
↳ overall change wrt every variable

$$y = m n + b \rightarrow \text{linear eq}$$

Linearly Separable data

↳ lines can separate data

↳ not possible in real life



Partial Derivatives

↳ everything else is constant

↳ one thing changes

Q2

$$f(u, y, z) = 2u + 3y + 5z$$

$$f'(u) = 2$$

$$f'(y) = 3$$

$$f'(z) = 5$$



constraint

$$z < 5, \quad x \leq 2, \quad y \leq 6$$

$$\begin{aligned} f(u, y, z) &= 2(2) + 3(6) + 5(4) \\ &= 42 \end{aligned}$$

easy eq so done manually

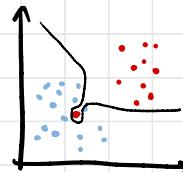
for harder eq
use gradient descent

problem

Non Linearity

↳ converts line into any shape

Activation



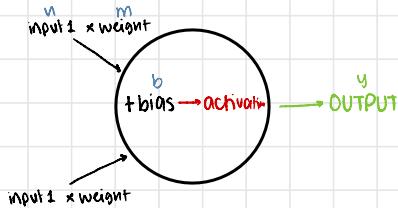
Large Regressor Models

↳ neural networks that have

↳ they are not used for classification

Neuron

- ↳ takes no. of inputs and gives an output



Neural Network

- ↳ has layers of neurons
- ↳ output of 1 neuron is input to another

↳ Virtual connections

Activation

- ↳ adds non linearity to the output of a neuron
- ↳ helps decide whether neuron should be activated or not

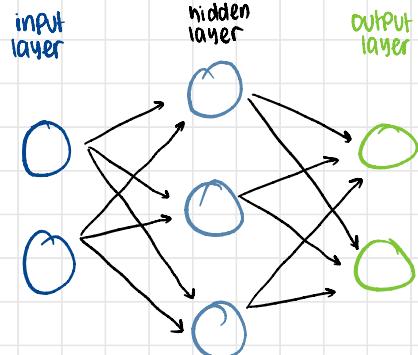
Neural Network: Why is it called a Network?

- A single neuron doesn't really help us find complex patterns (*and is not cool enough!*)
- We need a network of inter-connected neurons to make complex decisions
- The output of one neuron becomes the input of another neuron

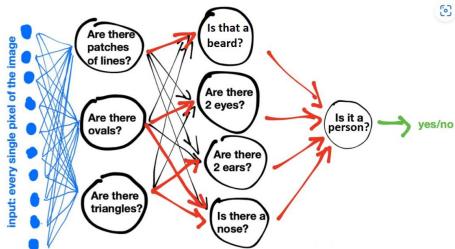


Neural Network Architecture

- ↳ It's composed of layers of neurons
- ↳ the more hidden layers the better
 - ✗ more expensive
 - ↳ can sometimes bring down accuracy
- ↳ It is a fully connected graph
- ↳ They are unidirectional only \Rightarrow one way



Neural Network: How does it work?



28

- features are abstract
- the layer closer to hidden layer has more detailed questions

Training

1. Feed raw input (features) to the input layer
2. Initialize all the weights and biases for hidden layers with random values
3. Test if the network can accurately produce the output
 - I. If it does not produce accurate results then adjust the weights and biases. It simply means we want to use optimization (e.g. gradient descent) to minimize the value of our loss function. Repeat Step 3
 - II. If it does produce accurate results then terminate training

To "adjust" the parameters (weights and biases) we use **backpropagation**

- When the neural network outputs the wrong answer, you find the slopes (derivative) of the output layer first because it was the direct cause of the incorrect answer
- Since the output layer depends on the hidden layer, you'll have to fix that too by finding the slopes and using gradient descent
- Eventually you'll work your way back (backpropagate) to the first hidden layer

32

$$y = \text{mv} + b$$

Annotations on the equation:
- 'weights' points to the 'm' term.
- 'biases' points to the 'b' term.
- 'changed after every epoch' points to the 'v' term.

↳ train till you

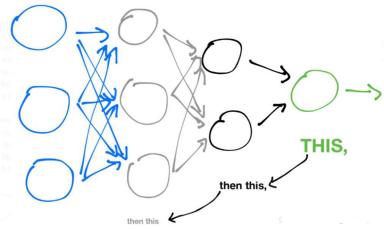
1. reach a certain accuracy
2. reach a certain epoch

BACKPROPAGATION

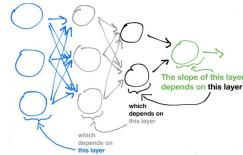
- ↳ To correct the network, you must fix the error towards the input layer → till the 1st hidden layer
- ↳ hidden layers never have any error
- ↳ start weights and biases with random values

$$\text{Error: } \frac{1}{2} (\text{Prediction} - \text{Actual})^2$$

To correct the network, you must first fix...



□ We calculate slopes by starting from the back and moving backwards through the network until we get all the slopes for gradient descent

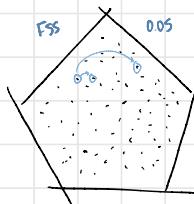


34

Exploration aka Learning Curve

- ↳ jumps to other region

CON
High exploration
might jump and skip solution



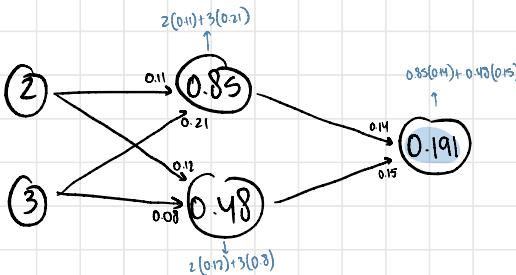
Exploitation

- ↳ trying near regions / neighbourhood

Q) $w_1: 0.11, w_2: 0.21, w_3: 0.12, w_4: 0.08, w_5: 0.14, w_6: 0.15$

Our target output is 1

FORWARD PASS



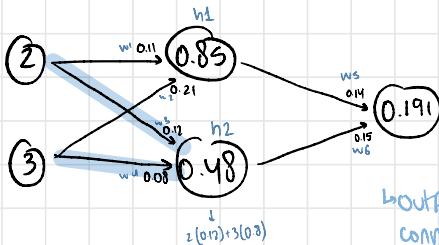
Error = $\frac{1}{2} (0.191 - 1)^2 = 0.327$

error is always five because of this square

↳ BACKWARDS PASS → to reduce the error

Gradient Descent

↳ Update the Value of Weights



↳ Output only knows
connected weights

$$\text{Prediction: } (\underbrace{i_1 w_1 + i_2 w_2}_{h1}) w_3 + (\underbrace{i_1 w_3 + i_2 w_4}_{h2}) w_5$$

$$\text{Error: } \frac{1}{2} (\text{prediction} - \text{actual})^2$$

Δ : Prediction - Actual → delta

□ The derivative of our loss can be evaluated by using Chain rule:

$$\frac{\partial \text{Error}}{\partial W_6} = \frac{\partial \text{Error}}{\partial \text{prediction}} * \frac{\partial \text{prediction}}{\partial W_6}$$

chain rule

$$\frac{\partial \text{Error}}{\partial W_6} = \frac{1}{2} (\text{predictoin} - \text{actula})^2 * \frac{\partial (i_1 w_1 + i_2 w_2) w_5 + (i_1 w_3 + i_2 w_4) w_6}{\partial W_6}$$

$$\frac{\partial \text{Error}}{\partial W_6} = 2 * \frac{1}{2} (\text{predictoin} - \text{actula}) \frac{\partial (\text{predictoin} - \text{actula})}{\partial \text{predictoin}} * (i_1 w_3 + i_2 w_4)$$

$$\frac{\partial \text{Error}}{\partial W_6} = (\text{predictoin} - \text{actula}) * (h_2)$$

$$\frac{\partial \text{Error}}{\partial W_6} = \Delta h_2$$

$\text{Error} = \frac{1}{2} (\text{prediction} - \text{actual})^2$

$\text{prediction} = (i_1 w_1 + i_2 w_2) w_5 + (i_1 w_3 + i_2 w_4) w_6$

$i_1 w_1 + i_2 w_2 \quad h1$

$i_1 w_3 + i_2 w_4 \quad h2$

$\Delta = \text{prediction} - \text{actual}$

39

$$\hookrightarrow \Delta = 0.191 - 1 = -0.809$$

$$\alpha = 0.05 \rightarrow \begin{matrix} \text{Learning} \\ \text{rate} \end{matrix} \rightarrow \begin{matrix} \text{guess} \\ \text{it} \end{matrix}$$

↳ Updating w_6

$$*w_6 = w_6 - \alpha \Delta h_2$$

$$= 0.15 - 0.05(-0.809)(0.08)$$

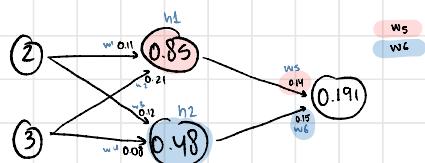
$$= 0.169 \approx 0.17$$

↳ Updating w_5

$$*w_5 = w_5 - \alpha \Delta h_1$$

$$= 0.14 - 0.05(-0.809)(0.08)$$

$$= 0.114 \approx 0.11$$



□ In order to update w_1 (existing between input layer and hidden layer):

$$\frac{\partial \text{Error}}{\partial w_1} = \frac{\partial \text{Error}}{\partial \text{prediction}} * \frac{\partial \text{prediction}}{\partial h_1} * \frac{\partial h_1}{\partial w_1} \quad \text{chain rule}$$

$$\text{Error} = \frac{1}{2}(\text{prediction} - \text{actual})^2$$

$$\text{prediction} = (h_1) w_5 + (h_2) w_6$$

$$h_1 = i_1 w_1 + i_2 w_2$$

$$\frac{\partial \text{Error}}{\partial w_1} = \frac{1}{2}(\text{predictoin} - \text{actual})^2 * \frac{\partial (\text{predictoin} - \text{actual})}{\partial \text{predictiton}} * \frac{\partial (\text{predictiton} - \text{actual})}{\partial w_1} * (w_5) * (i_1)$$

$$\frac{\partial \text{Error}}{\partial w_1} = (\text{predictoin} - \text{actual}) * (w_5) \quad \Delta = \text{prediction} - \text{actual}$$

$$\frac{\partial \text{Error}}{\partial w_1} = \Delta w_5$$

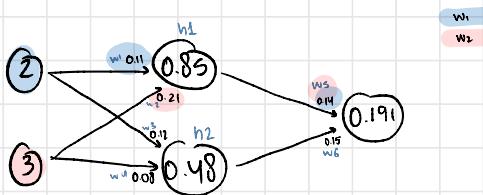
4

↳ Updating w_1

$$\begin{aligned} * w_1 &= w_1 - \alpha \Delta (w_5 \cdot i_1) \\ &= 0.11 - 0.05(-0.809)(0.14)(2) \\ &= 0.121 \end{aligned}$$

↳ Updating w_2

$$\begin{aligned} * w_2 &= w_2 - \alpha \Delta (w_5 \cdot i_1) \\ &= 0.21 - 0.05(-0.809)(0.14)(3) \\ &= 0.227 \approx 0.23 \end{aligned}$$

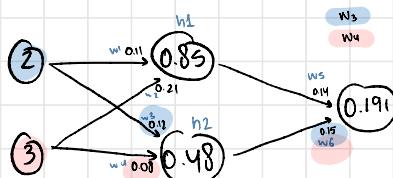


↳ Updating w_3

$$\begin{aligned} * w_3 &= w_3 - \alpha \Delta (w_6 \cdot i_1) \\ &= 0.12 - 0.05(-0.809)(0.15)(2) \\ &= 0.132 \approx 0.13 \end{aligned}$$

↳ Updating w_4

$$\begin{aligned} * w_4 &= w_4 - \alpha \Delta (w_6 \cdot i_2) \\ &= 0.08 - 0.05(-0.809)(0.15)(3) \\ &= 0.098 \approx 0.10 \end{aligned}$$

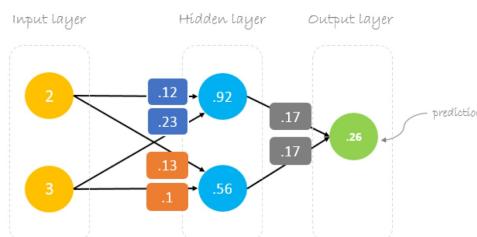


- All the updated weights (*rounded off to 2 digits*) are as follows:

- $w_1 = 0.12$
- $w_2 = 0.23$
- $w_3 = 0.13$
- $w_4 = 0.10$
- $w_5 = 0.17$
- $w_6 = 0.17$

- Task:** Re-calculate and verify the updated values of these weights

- Now, using the new weights we will repeat the forward pass



- Notice that our prediction is slightly better but still not perfect. Therefore, keep updating the weights and performing backward-forward passes until the error is minimized

45

Are we still missing something?

- Bias
- Activation
 - Sigmoid (*for binary classification*)
 - ReLU (*suffers from vanishing gradient problem*)
 - LeakyReLU (*improved ReLU*)
 - Softmax (*for multi-class classification*)
- Regularization
- Dropout
- Optimization
 - Stochastic Gradient Descent
 - Batch and Mini-batch Gradient Descent
 - Momentum
 - Adam

- a) Explain why the following ratings data cannot be properly split into training and testing samples for neural recommendations:

	I1	I2	I3	I4	I5	I6	I7	I8	I9	I10
U1	2	1		3	5	4		2	1	2
U2										
U3	4	5		3	2	4			5	1
U3		1	4	3	4	2		1	1	5
U5	4	2		1	3	2		1	2	3
U6	3			1		2		3	2	2

Note: Each row represents a unique user and each column is a unique item.

- b) Explain why we are usually interested in interactions more than ratings when working with recommendations generated by a GAN.

a) as it is a sparse matrix

↳ User 2 didn't rate any item

↳ Item 7 has no ratings by any user

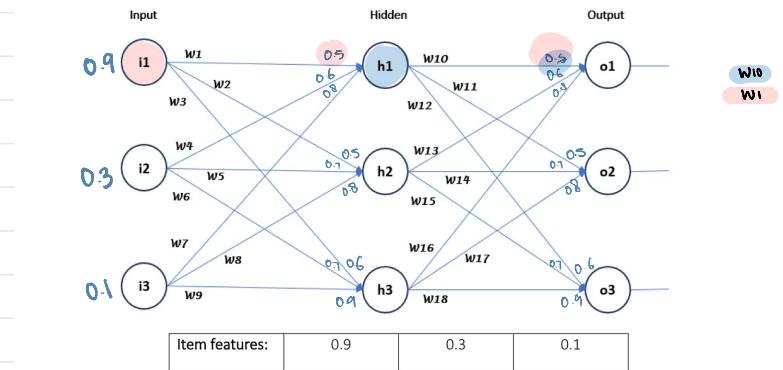
So we can't split such data

b) GANs are extremely sensitive

providing ratings is much more complex to converge

hence we use interactions (0 or 1)

Consider the neural network given below:



Initial weights are $w_1 = w_2 = w_{10} = w_{11} = 0.5$; $w_3 = w_4 = w_{12} = w_{13} = 0.6$; $w_5 = w_6 = w_{14} = w_{15} = 0.7$; $w_7 = w_8 = w_{16} = w_{17} = 0.8$; and $w_9 = w_{18} = 0.9$. Assume that the true output of o_1 was 0.7. Find the error at o_1 and run a backward pass to update w_{10} .

$$h_1 = (0.9)(0.5) + (0.3)(0.6) + (0.1)(0.8) = 0.71$$

$$h_2 = (0.9)(0.5) + (0.3)(0.7) + (0.1)(0.8) = 0.74$$

$$h_3 = (0.9)(0.6) + (0.3)(0.7) + (0.1)(0.9) = 0.84$$

$$O_1 = (0.71)(0.5) + (0.74)(0.6) + (0.84)(0.8) = 1.471$$

$$\text{error: } \frac{1}{2} (1.471 - 0.7)^2 = 0.297$$

$$\text{taking } \alpha = 0.05$$

$$\Delta = 1.471 - 0.7 = 0.77$$

Updating w_{10}

$$\begin{aligned} *w_{10} &= w_{10} - \alpha \Delta h_1 \\ &= 0.5 - 0.05(0.77)(0.71) \\ &= 0.41 \end{aligned}$$

Updating w_1

$$\begin{aligned} *w_1 &= w_1 - \alpha (\Delta \cdot w_{10} \cdot i_1) \\ &= 0.5 - 0.05(0.77)(0.5)(0.9) \\ &= 0.483 \end{aligned}$$

Question 5: Consider the given inputs and parametric values for a simple recommendation neural network. The following tables show data for three training samples. Each sample has distinct user and item features as the input. We can ignore the activations and hidden layer values to denote them simply as a , b , c and d . The output of the network is either 0 or 1 representing the predicted interaction of the particular user with the given item.

User Features	I	H1	H2	O
Item Features	1	a	c	1
	0.99	a	c	1
	0.005	b	d	1
	0.99	a	b	1
	0.8	c	d	1

Answer the following questions:

- Based on the observed values for input features. Will it be a good idea to omit user features from the input and use only item features for classification? Justify your answer with a reason.
- Given the variance in item features, is the model exhibiting over-fit or under-fit behavior?
- What should be the ideal activation function in the output layer of this network?
- For 1000 users, if we assume the variance in their features to be 0.08. What does it explain about the users present on the website/system?

a) both U and I features vary b/w samples

U features contributes unique info about users preference

which can enhance recommendations accuracy

so omitting U features is not a good idea AS

it will lead to loss of crucial info

b) all 'O' values are 1 → consistent output

item features have low variance → low variance features

this pattern suggests overfitting to a specific item

c) since output predicts 0 and 1

a sigmoid activation function is recommended

d) It implies low diversity in user feature values → as low variance of 0.08

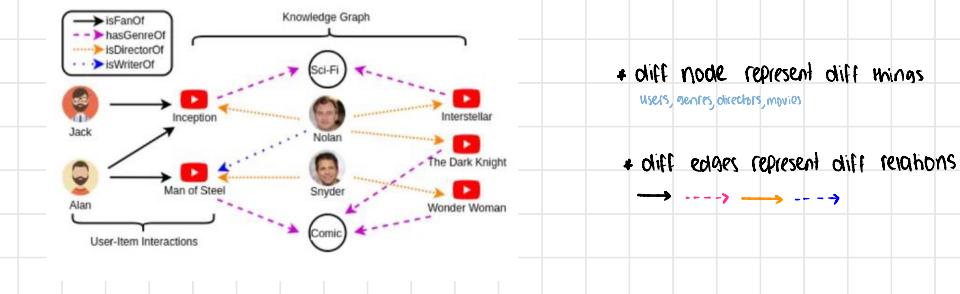
meaning users on the website mostly have same preferences

Consider the case where the data points are linearly separable. What will be the role of activations in a neural network architecture with an input layer, several hidden layers and an output layer?

If the data points in a given training set are already linearly separable, there isn't any need to induce non-linearity, hence there should be no role of activations in that case and a large regressor would suffice.

KNOWLEDGE GRAPHS based SYSTEMS

- ↳ uses graphs to represent users, items and their relationships
- ↳ uses a network of interconnected data
- ↳ nodes → entities
- ↳ how entities connected in social network



↳ Diversity

- ↳ KGs represents semantically rich relations b/w entities
- ↳ different aspects of user preference
 - ↳ director
 - ↳ genre

↳ Scalability

- ↳ doesn't need sparse matrix to predict similarity
- ↳ can't be vector based
 - ↳ as can define dimensions

↳ Users preference are accumulated and stored in the graph structure



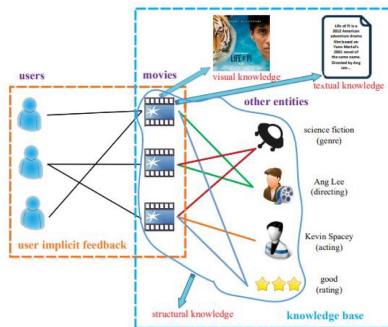
↳ ALGOS we can use for KG based recommendation

1. Embedding based → Neural Networks
2. Path based → Data Structures
3. RippleNet

1 Embedding based KG Recommendation

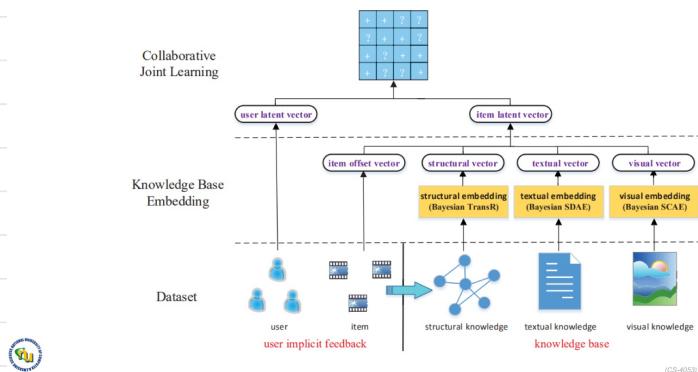
- ❑ Pre-process KG using knowledge graph embedding algorithms
 - ❑ For **structural** knowledge: TransR embedding
 - ❑ For **textual** knowledge: Word2Vec, Stacked Denoising Autoencoder
 - ❑ For **visual** knowledge: Stacked Convolutional Autoencoder
- ❑ movie posters
↳ movie is closer to which movie in KG
↳ movie synopsis
↳ actors
- ❑ The generated embeddings become input to joint collaborative (content) filtering
- ❑ Use cosine similarity (or some other measure) for finding distance of u to v_i ($i=1, 2, \dots, n$)

Embedding-based KG Recommendation



(CS-4053) Course Instructor: Syed Zain

Embedding-based KG Recommendation

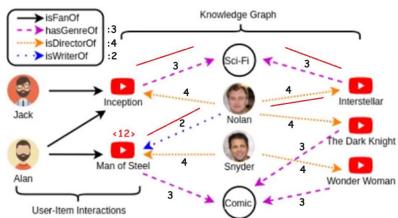


(CS-4053) Course Instructor: Syed Zain (

2. Path based KG Recommendation

- For an active user, traverse a path from each of the k top-rated items of the user
- Each edge is assigned a feature score
- The path follows edges representing knowledge about the item
- Each edge is visited once and its feature score is added to total path score
- For each path of an unrated item:
 - It ends on an unrated item and its path score is recorded
 - The maximum path score across all paths become the item utility value
- The top- k items with maximum utility values are recommended

Path-based KG Recommendation



- Domain knowledge is used to assign feature scores
- The path score in this example is 12

Inception → Sci-Fi → Interstellar → Nolan → Man of steel

$$3 + 3 + 4 + 2 = 12$$

↳ each edge is visited once

item utility value

↳ for each path of an unrated item

↳ ends on unrated item

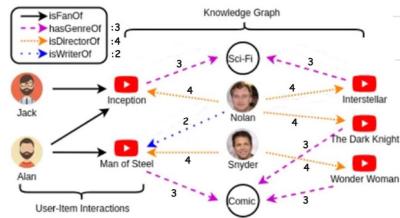
↳ max path score



top k items with
max utility is

RECOMMENDED

Path-based KG Recommendation



Inception → Nolan]

wonder ← snyder ← man of
woman

↓
recommended

?

Path

3. RippleNet

↳ inputs

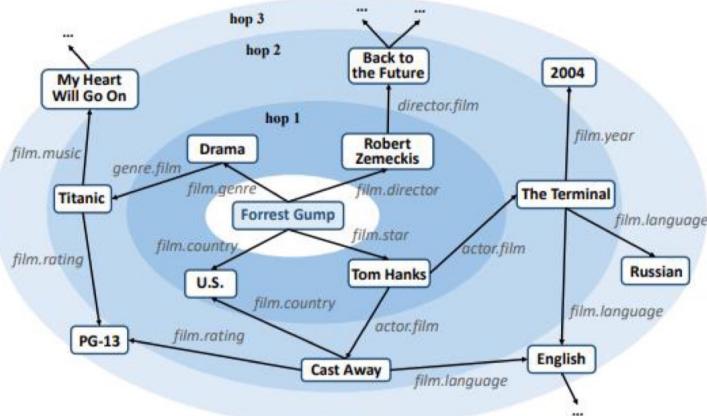
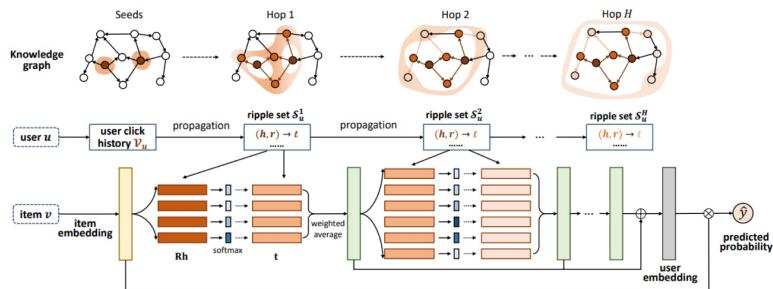
↳ user u

↳ item v

↳ outputs · predicted probability user u will click item v

- ◻ For the input user u , his historical set of interests V_u is treated as seeds in the KG, then extended along links to form multiple ripple sets
- ◻ A ripple set is the set of knowledge triples that are k -hop(s) away from the seed set V_u . These ripple sets are used to interact with the item embedding iteratively for obtaining the responses of user u with respect to item v , which are then combined to form the final user embedding
- ◻ Lastly, we use the embeddings of user u and item v together to compute the predicted probability

RippleNet



Generative Adversarial Networks (GANs)

- ↳ An unsupervised architecture
- ↳ Uses 2 neural networks, competing against each other

thus the adversarial

that can pass
for real data

- ↳ They can learn to mimic any distribution of data

- ↳ Can be used for
 - ↳ image generation
 - ↳ voice generation
 - ↳ video generation

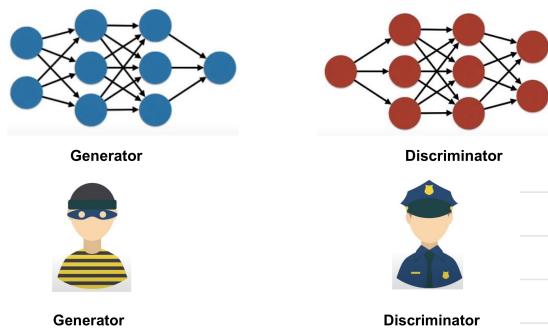
Image Generation

 [This Person Does Not Exist - Random Face Generator \(this-person-does-not-exist.com\)](https://thispersondoesnotexist.com/)



GAN architecture

- ↳ has 2 networks
 - ↳ generator
 - ↳ discriminator



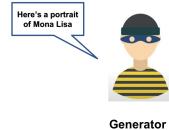
Adversarial Learning

↳ It's a 2 Player game b/w the generator and discriminator

The Generator

- ↳ Starts with a random noise
- ↳ tries to mimic the target distribution
 - e.g. image of oil rig
 - image of mona lisa

↳ (resonable) random noise



Generator

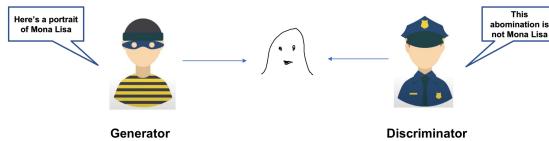
The Discriminator

- ↳ has already seen the training data
- ↳ so knows the target distribution
- ↳ so acts like a trainer and judge
 - No! It's obviously not



Discriminator

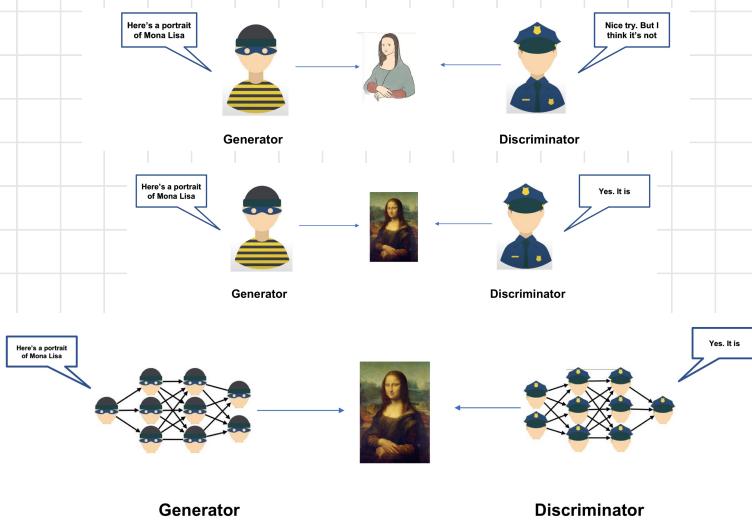
- ↳ It repeatedly tries to fool the discriminator
 - ↳ image
- ↳ that the generated distribution is a real one



Generator

Discriminator

- ↳ based on feedback of the discriminator
- ↳ the quality of the generated output is gradually improved



STEPS TAKEN BY GANS

- ↳ The generator takes in random numbers and returns an image

- ↳ The discriminator is fed

- ↳ this generated image

- ↳ a stream of images

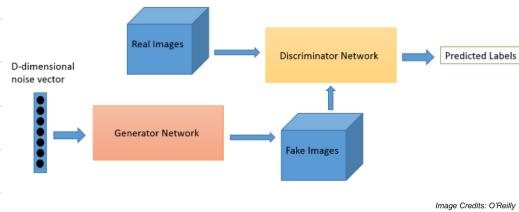
- ↳ It takes both real and fake images

- ↳ taken from the actual ground truth dataset

- ↳ and returns probabilities representing

- ↳ 1 → a prediction of authenticity

- ↳ 0 → fake



$$\min_G \max_D V(D, G)$$

$$V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]$$

EXAMPLE: Generating MNIST digits

Generator

- ↳ An inverse CNN

- ↳ Input → vector of random noise

- ↳ Output → upsamples it to an image

Steps

1. Show MNIST dataset to discriminator

- ↳ never train discriminator full dataset

- ↳ on 30 to 35% dataset shown

2. Train discriminator as a classifier

3. Generator will take noise and produce image

- ↳ It will have some loss that discriminator will use

Discriminator

- ↳ A standard CNN

- ↳ categorises images fed to it

- ↳ always a binary classifier, labels images

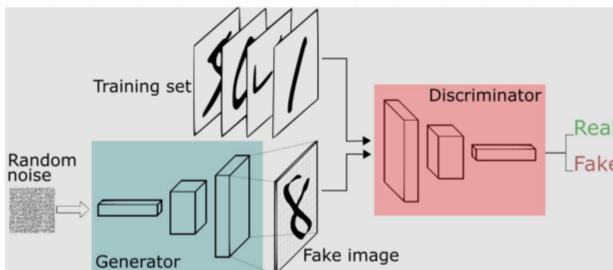
- ↳ real

- ↳ fake

MNIST dataset

- ↳ handwritten dataset of digits

- 0 – 9



HOW TO train a GAN

- ❑ Build your generator and discriminator models
- ❑ Sample data from true distribution and train discriminator on this data
 - ❑ Discriminator works as a binary classifier
- ❑ Create input data from latent space (noise) and feed it to the generator
- ❑ The generator output serves as input to the discriminator
- ❑ The discriminator uses Sigmoid activation to classify the generated output as real or fake
- ❑ The discriminator loss is fed-back to the generator
- ❑ The generator is trained until it reaches convergence

CON

1. GANs are hard to train → as convergence hard to attain
2. Hard to meet convergence criteria → how realistic do I want it to be
doesn't take time into account
3. Vanilla GANs suffer from vanishing gradient problem → gradients start becoming 0 so then training can't be done → as if user feedback of loss
4. Overwhelming adversary → discriminator shouldn't be overpowered → unless very realistic work except it
discriminator shouldn't be weak → unrealistic images accepted
5. Can not count objects → can't say make 5 Mona Lisa Paintings → Mona Lisa can't draw things perfectly
6. Sensitive to choice of hyperparameters → no optimal hyperparameter
7. ... and many other challenges

Guidelines for training

- ↳ use Tanh activation for generator
- ↳ use Sigmoid for discriminator → as 0 fake 1 real
- ↳ freeze generator for every other training iteration of discriminator
- ↳ USE LEAKYRELU in dense layers
- ↳ Use a Functional API → for keras
- ↳ Use drop out
 - ↳ randomly switches off some neurons during training → so no overfitting

TYPES OF GANS

- ↳ Vanilla GANs → OG
- ↳ CGAN ↳ Conditional GAN ↳ Uses class label as input
- ↳ Wasserstein GAN
- ↳ Style GAN ↳ Input product ↳ Output related product ↳ Pink shirt ↳ Pink hat
- ↳ Info GAN
- ↳ DCGAN
- ↳ DualGAN
- ↳ KDGAN
- ⋮

GANS in RS

- ↳ used in RS to generate ratings
- ↳ solves cold start problem
- ↳ slow training but prediction are fast
- ↳ CGAN
- ↳ Style GAN

BASIC DESIGN PRINCIPLES

- ↳ for recommendation output is always a 1D vector
- ↳ GANs pick one

- ◻ We never work on **ratings** i.e., these ratings are to be transformed to interactions
- ◻ We always use a supervised variant of GAN for recommender system e.g., cGAN, cWGAN, IRGAN
- ◻ We can do either user-to-item generation or item-to-user generation
- ◻ Always watch out for a trivial solution (*network learning to output all zeros or ones*)

Example: MIND News Recommendation

- ↳ a news dataset with

↳ impression logs

- ↳ historical choices for users and news items

ImpressionID	UserID	Time	History	Impressions
0	1	11/11/2019 9:05:58 AM	N55189 N42782 N34694 N45794 N18445 N63302 N104...	N55689-1 N35729-0
1	2	11/12/2019 6:11:30 PM	N31739 N6072 N63045 N23979 N35656 N43353 N8129...	N20678-0 N39317-0 N58114-0 N20495-0 N42977-0 N...
2	3	11/14/2019 7:01:48 AM	N10732 N25792 N7563 N21087 N41087 N5445 N60384...	N50014-0 N23877-0 N35389-0 N49712-0 N16844-0 N...
3	4	11/11/2019 5:28:05 AM	N45729 N2203 N871 N53880 N41375 N43142 N33013...	N35729-0 N33632-0 N49685-1 N27581-0
4	5	11/12/2019 4:11:21 PM	N10078 N56514 N14904 N33740	N39985-0 N36050-0 N16096-0 N8400-1 N22407-0 N6...

User saw but didn't click → ignored
 User saw and clicked → historical

CONSTRUCT TRAINING DATA

- ↳ create training data by treating

↳ historical choice → positive interaction → 1

↳ ignored news → negative interaction → -1

change 0 to -1

UserID	NewsID	Interaction
78833	U46778	1

UserID	NewsID	Interaction
78833	U46778	-1
78833	U46778	-1
78833	U46778	1
78833	U46778	-1
78833	U46778	-1

→ change data so it's suitable for GANs

↳ row → user

↳ column → news item

↳ intersection → interaction

news item
 user
 interaction

$$\begin{bmatrix} [1, 0, 0, \dots, 0, 0, 0] \\ [1, 0, 0, \dots, 0, 0, 0] \\ [1, 0, 0, \dots, 0, 0, 0] \\ \dots \\ [0, 0, 0, \dots, 0, 0, 0] \\ [0, 0, 0, \dots, 0, 0, 0] \\ [0, 0, 0, \dots, 0, 0, 0] \end{bmatrix}$$

Training Data Shape: (26040, 6838)

Interaction Table

-1 → saw, not interested
 0 → didn't see
 1 → saw, interested

↓ NO labels

since user features are not given and we need class labels

so do clustering on interest of user → news features → what type of news does user like

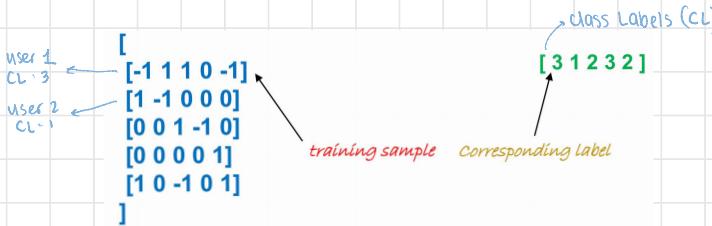
↳ every cluster is a class label

■ We often do not have any labelled recommendation data hence we can use clustering or some other annotation technique to create our label set

NewsID	Category	Subcategory	Title	Abstract	Processed_Text	Embeddings	Label	
2	N61837	news	newsworld	The Cost of Trump's Aid Freeze in the Trenches...	Lt. Ivan Molchanov peeked over a parapet of s...	news newsworld cost trump aid freeze trench u...	[-0.05749018, -0.05521024, -0.016900085, 0.030...	5
3	N53526	health	voices	I Was An NBA Wife. Here's How It Affected My M...	I felt like I was a fraud, and being an NBA wi...	health voice nba wife here affected mental hea...	[-0.03127052, -0.029868228, -0.028364867, 0.05...	2
4	N38324	health	medical	How to Get Rid of Skin Tags, According to a Doc...	They seem harmless, but there's a very good re...	health medical get rid skin tag according derm...	[-0.07352022, -0.0478042, -0.007855178, 0.0447...	1
5	N2073	sports	football_nfl	Should NFL be able to fine players for critic...	Several fines came down against NFL players fo...	sport footballnfl nfl able fine player critici...	[0.018641314, -0.004955198, -0.01916305, 0.091...	3
10	N9721	health	nutrition	50 Foods You Should Never Eat, According to He...	This is so depressing.	health nutrition food never eat according heal...	[-0.05535877, -0.08080809906, -0.044505723, 0.01...	9

GAN TRAINING

↳ we can train it in the following manner



Example: Generation Recommendation

Conditional input



Synthetic interactions

ishma hafeez
notes
represent

* if cold start user

↳ forward it to news predictor