

Recommender Systems || Lecture 7

Summer 2022

Dr. Karthik Mohan

Univ. of Washington, Seattle

July 21, 2022

Today

- 1 SGD Review

Today

- ① SGD Review
- ② Auto differentiation ✓

Today

- ① SGD Review
- ② Auto differentiation

Today

- ① SGD Review
- ② Auto differentiation
- ③ Auto Encoders for Recommendation Systems

SGD in practice - mini-batch SGD!

mini-batch SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

- 1 **Initialize** $w = w_0$ (randomize)

(u^0, v^0)

SGD in practice - mini-batch SGD!

mini-batch SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

- 1 **Initialize** $w = w_0$ (randomize) Pick a batch of k data points at random between 1 and N : $i_1, i_2, \dots, i_k!$

$i_1, i_2, \dots, i_k!$
↳ k indices

SGD in practice - mini-batch SGD!

mini-batch SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

① **Initialize** $w = w_0$ (randomize) Pick a batch of k data points at random between 1 and N : i_1, i_2, \dots, i_k !

② **Gradient Descent** $w^{k+1} \leftarrow w^k - lr * \sum_{j=1}^k \nabla_w L_{i_j}(w^k)$

Assignment 2 part 2
Fill this in

SGD in practice - mini-batch SGD!

mini-batch SGD

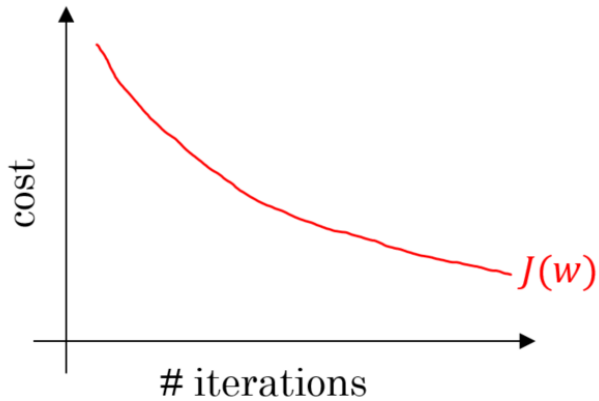
Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

- 1 **Initialize** $w = w_0$ (randomize) Pick a batch of k data points at random between 1 and N : i_1, i_2, \dots, i_k !
- 2 **Gradient Descent** $w^{k+1} \leftarrow w^k - lr * \sum_{j=1}^k \nabla_w L_{i_j}(w^k)$
- 3 **Iterate** Repeat step 2 and 3 until w converges, i.e.

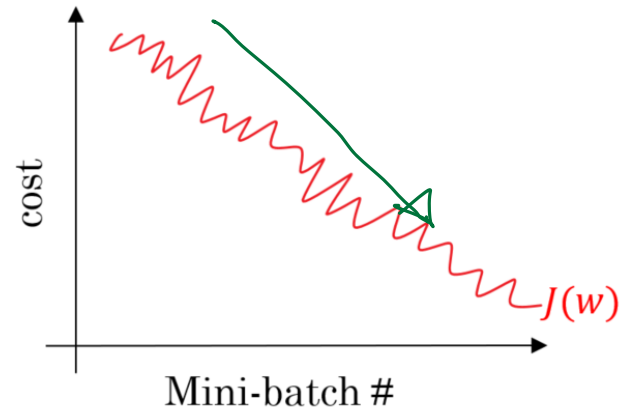
$$\|w^{k+1} - w^k\| / \|w^k\| \leq 10^{-3}$$

GD vs Mini-batch convergence behavior

Batch gradient descent



Mini-batch gradient descent



GD vs mini-batch SGD

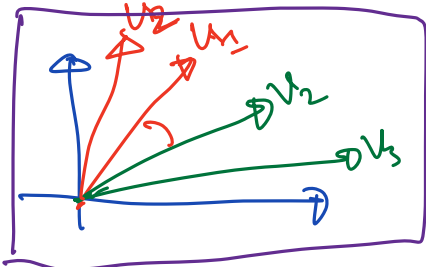


Factor	GD	Mini-batch SGD
Data	All per iteration	Mini-batch (usually 128 or 256)
Randomness	Deterministic	Stochastic
Error reduction	Monotonic	Stochastic
Computation	High	<u>Low</u>
Memory big data	Intractable	<u>Tractable</u>
Convergence	Low relative error	Few “passes” on data
Local Minima traps	Yes	<u>No</u>

ICE #1

The correct loss function based on user factors (V) and movie factors (U) for this problem is given by:

- ① $L(X; U, V) = (u_1^T v_1 - 1)^2 + (u_1^T v_3)^2 + (u_2^T v_2 - 1)^2 + (u_2^T v_4)^2 + (u_3^T v_1)^2 + (u_3^T v_2 - 1)^2 + (u_4^T v_1)^2 + (u_4^T v_2)^2 + (u_4^T v_3 - 1)^2$
- ② $L(X; U, V) = (v_1^T u_1 - 1)^2 + (v_1^T u_3)^2 + (v_2^T u_2 - 1)^2 + (v_2^T u_4)^2 + (v_3^T u_1)^2 + (v_3^T u_2 - 1)^2 + (v_4^T u_1)^2 + (v_4^T u_2)^2 + (v_4^T u_3 - 1)^2$
- ③ None of the above
- ④ Both are valid



Movies |

	Tom	Henry	Arnold	Shafiq
Sun	<u>1</u>	?	0	?
Moon	?	<u>1</u>	?	0
Stars	<u>0</u>	<u>1</u>	?	?
Galaxies	0	<u>0</u>	<u>1</u>	?

Automatic Differentiation in Torch

Notebook Example



Programming 2, Part 2

Notebook Walkthrough



Project Proposal (1 pager)

Auto Encoders for Recommendation Systems

$$L(X; U, V) = \frac{1}{N} \sum_{(i,j) \in O} (X_{ij} - \underbrace{U_i^T V_j}_{\substack{\downarrow \\ \text{Bilinear fn.}}})^2$$

Handwritten annotations in pink:
- X : Data
- U, V : Parameters
- N : # of observed entries
- $(i,j) \in O$: Observed entries
- $U_i^T V_j$: Bilinear function

- 1 Natural extension of matrix factorization
- 2 Matrix Factorization methods are “bi-linear” (linear in U , V) and linear models have limitations in learning power
- 3 Auto Encoders can have non-linearity

Model 1: AutoRec

Loss: $\frac{1}{|U|} \sum_u (x_u - \hat{x}_u)^2$

MSE \rightarrow mean squared Error

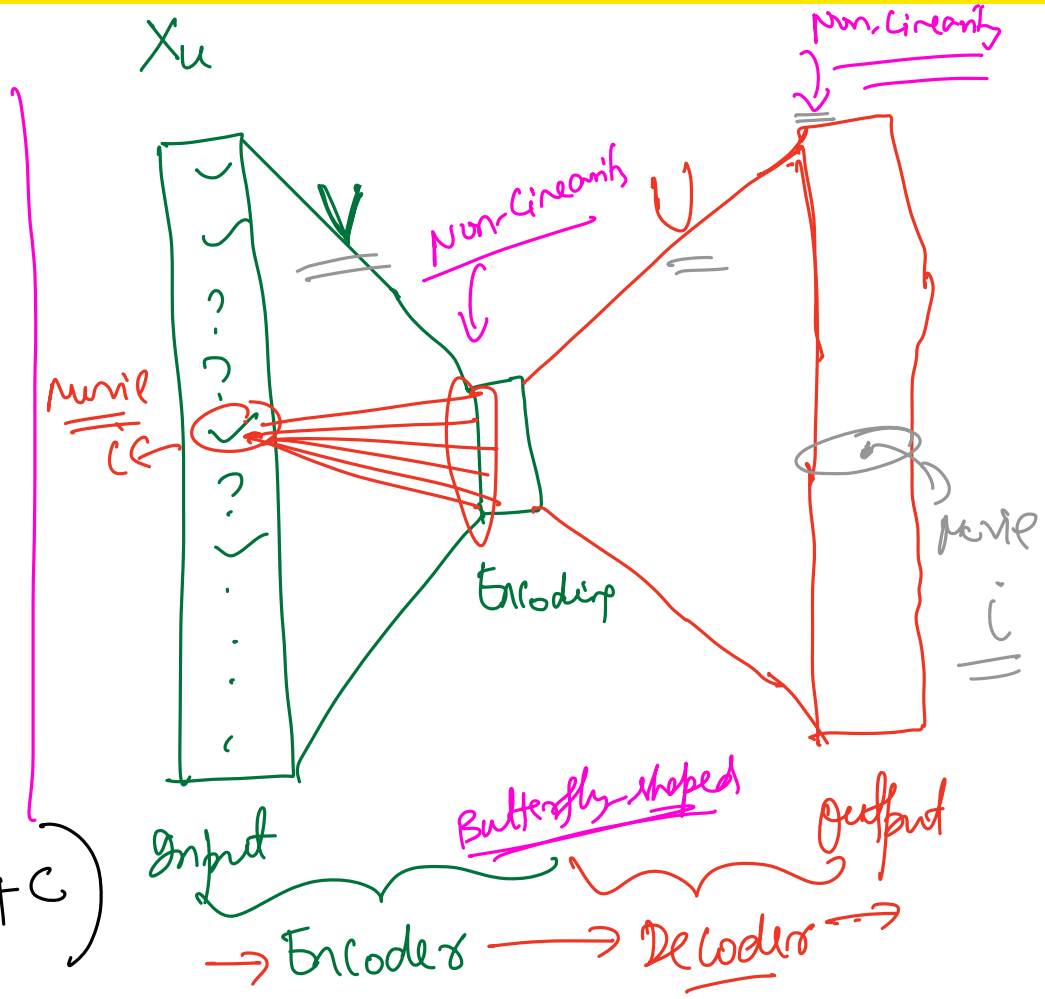
Reference: AutoRec

AutoEncoder for Rec Sys

Estimated \hat{x}_u

non-linearity

$$\hat{x}_u = g(U f(Vx^u + b) + c)$$



Model 1: AutoRec

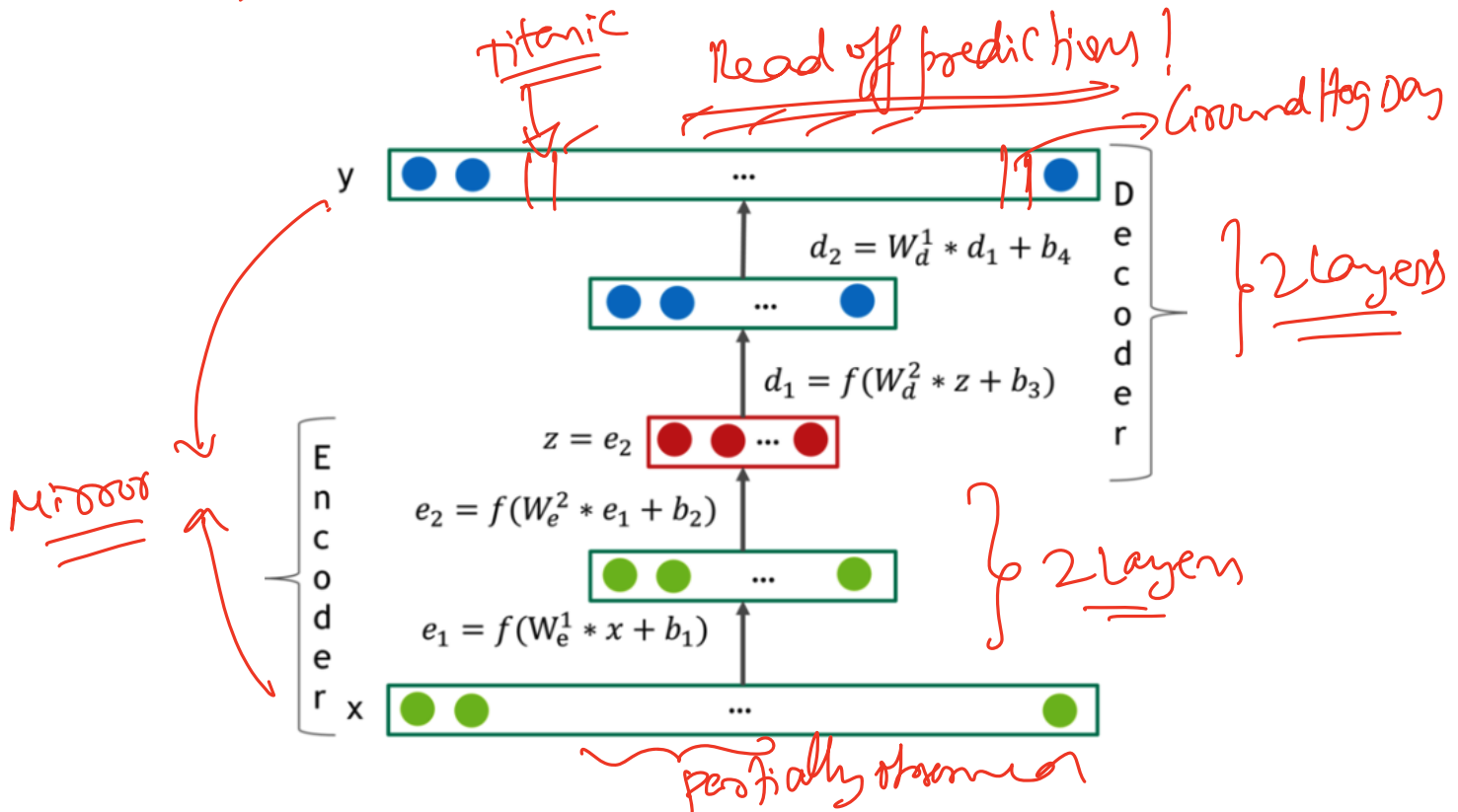
$f(\cdot)$	$g(\cdot)$	RMSE				
Identity	Identity	0.872	BiasedMF	0.845	0.803	0.844
Sigmoid	Identity	0.852	I-RBM	0.854	0.825	-
Identity	Sigmoid	0.831	U-RBM	0.881	0.823	0.845
Sigmoid	Sigmoid	0.836	LORMA	0.833	0.782	0.834
			<u>I-AutoRec</u>	0.831	0.782	0.823

movie lens (handwritten arrow pointing to ML-1M and ML-10M columns)

netflix Prize Dataset (handwritten arrow pointing to Netfix column)

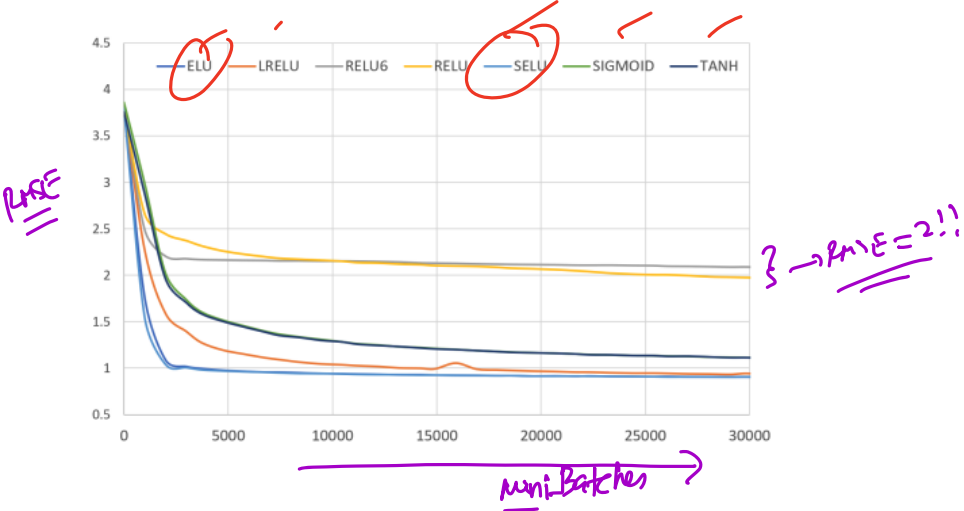
Reference: AutoRec

Model 2: Deep Rec



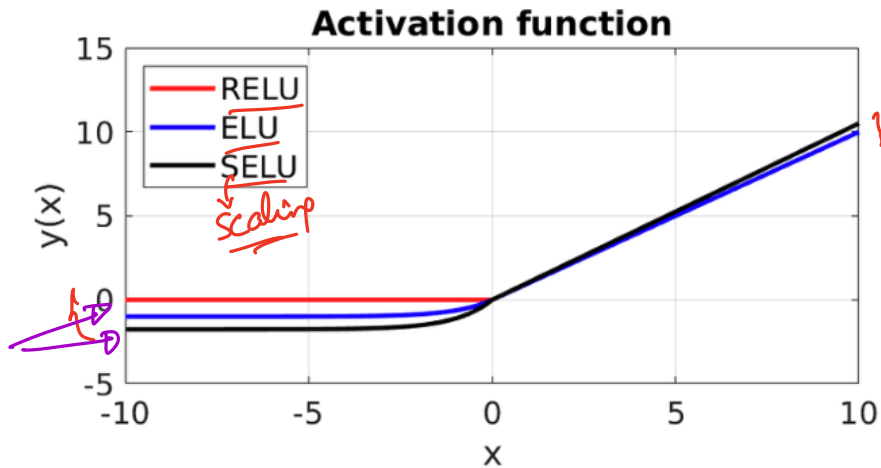
Training Deep Auto Encoders for collaborative filtering

Model 2: Deep Rec



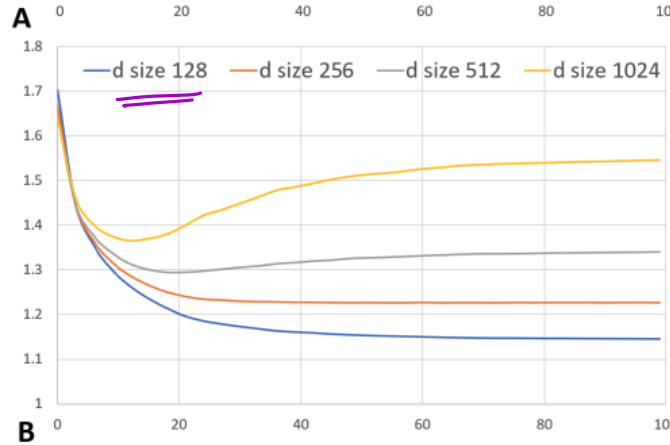
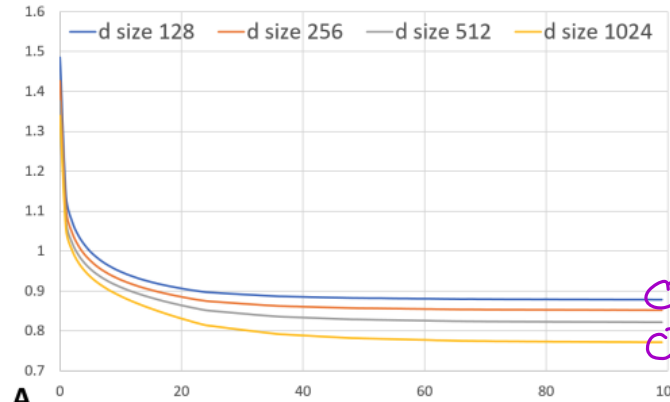
Training Deep Auto Encoders for collaborative filtering

Model 2: Deep Rec



Training Deep Auto Encoders for collaborative filtering

Model 2: Deep Rec



Training Deep Auto Encoders for collaborative filtering

Model 2: Deep Rec

<u>Number of layers</u>	<u>Evaluation RMSE</u>	<u>params</u>
2	<u>1.146</u>	4,566,504
4	<u>0.9615</u>	4,599,528
6	<u>0.9378</u>	4,632,552
8	0.9364	4,665,576
10	0.9340	4,698,600
12	0.9328	<u>4,731,624</u>

Handwritten notes:

- A purple arrow points down from 2 to 4 layers.
- A purple arrow points down from 4 to 6 layers.
- A purple arrow points down from 6 to 8 layers.
- A purple arrow points down from 8 to 10 layers.
- A purple arrow points down from 10 to 12 layers.
- The text "Point of Diminishing Returns" is written in purple, with a double underline, and a purple arrow pointing to the 8-layer row.
- The text "57M parameters" is written in purple, with a double underline, and a purple arrow pointing to the 12-layer row.

Training Deep Auto Encoders for collaborative filtering

Model 2: Deep Rec

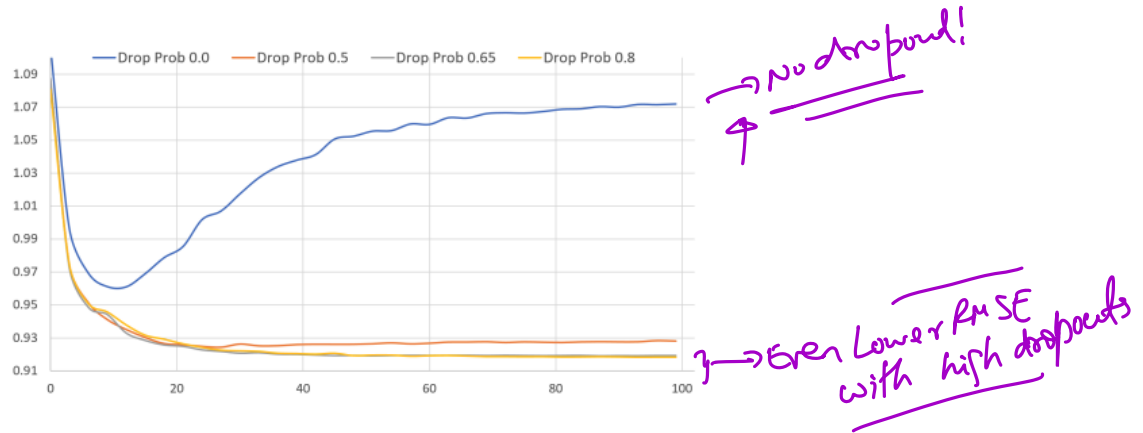


Figure 4: Effects of dropout. Y-axis: evaluation RMSE, X-axis: epoch number. Model with no dropout (Drop Prob 0.0) clearly over-fits. Model with drop probability of 0.5 over-fits as well (but much slowly). Models with drop probabilities of 0.65 and 0.8 result in RMSEs of 0.9192 and 0.9183 correspondingly.

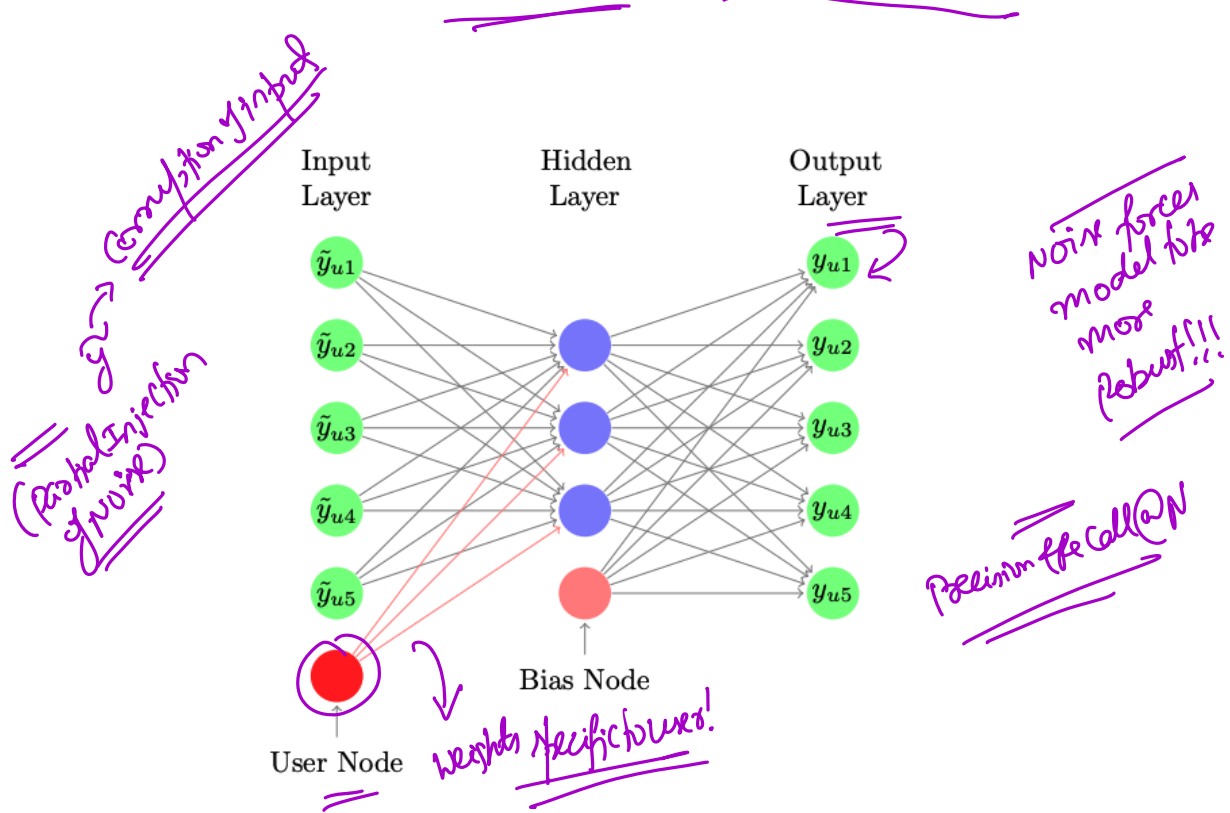
Training Deep Auto Encoders for collaborative filtering

Model 2: Deep Rec

DataSet	I-AR	U-AR	RRN	DeepRec
Netflix 3 months	0.9778	0.9836	0.9427	0.9373
Netflix Full	0.9364	0.9647	0.9224	0.9099

Training Deep Auto Encoders for collaborative filtering

Model 3: Collaborative De-noising AutoEncoder



CDAE

Summarize

1. MF models - Bi-Linear
2. Bi-linear models Limitations (RMSE limitation)
3. Auto-Encoders work well for RecSys
4. Non-Linear & they can be deep
5. Attain SOTA on many benchmark datasets
6. Scalability issues! (E.g. output layer becomes too big).

DSSM | Siamese Networks } In next Lecture!
(NLP)