

CF-NADE: A Neural Autoregressive Approach to Collaborative Filtering

Yin Zheng, Bangsheng Tang, Wenkui Ding, Hanning Zhou

Hulu LLC.

Beijing, China

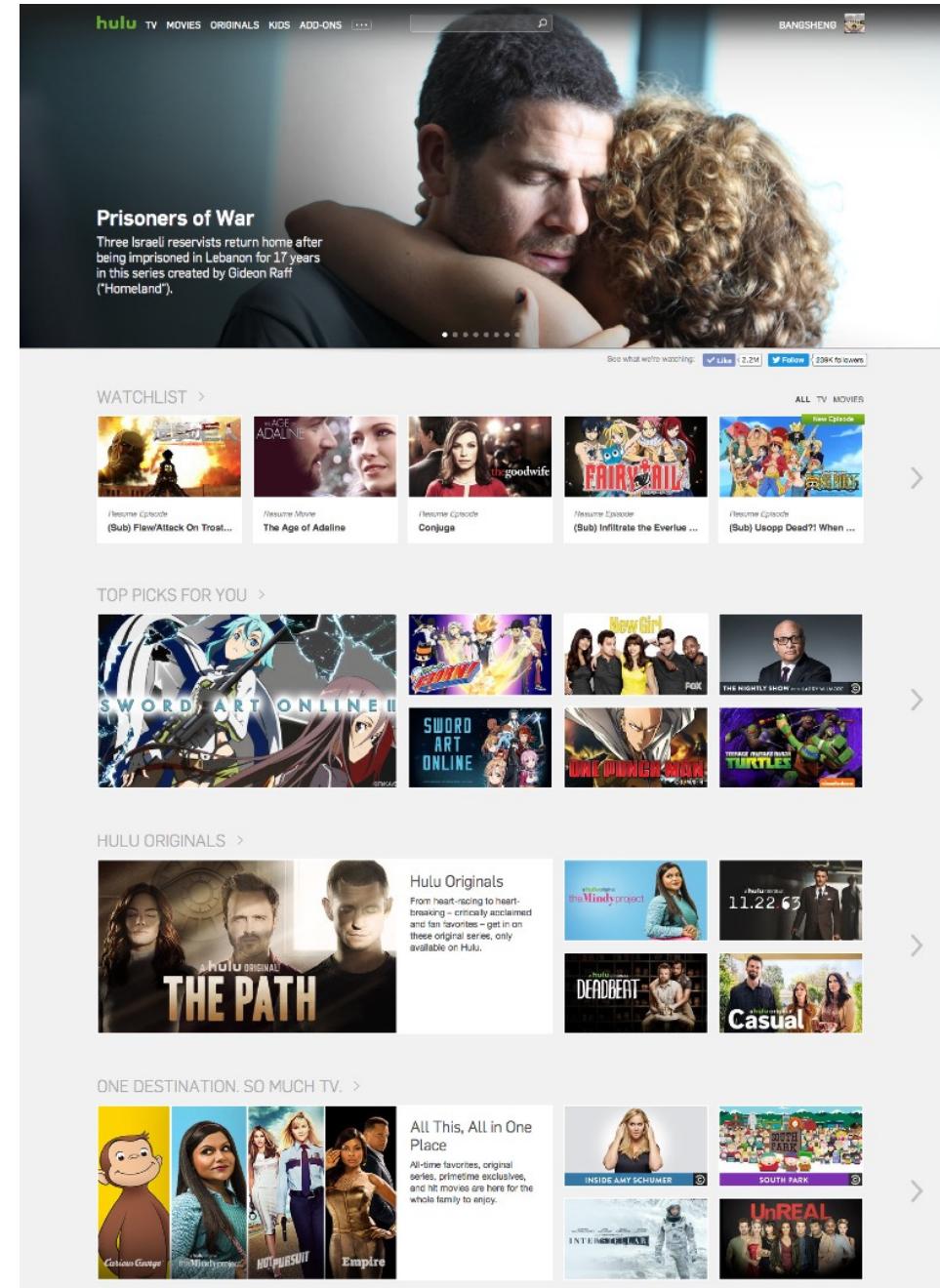
<https://sites.google.com/site/zhengyin1126/>

Outline

- Motivation
- Related Work
- CF-NADE
 - The basic model
 - Sharing parameters
 - Scalability
- Ordinal Cost in Collaborative Filtering
- Deep CF-NADE
- Experimental Results

Motivation

- Collaborative Filtering (CF):
 - Predict users' preferences
 - Previous behavior
 - Similar users
 - One's preferences do not change enormously
- Core of Recommender Systems
 - Hulu:
 - Personalized Homepage
 - Personalized Masthead
 - Watch list
 - Top picks for you
 - Related Shows
 - ...

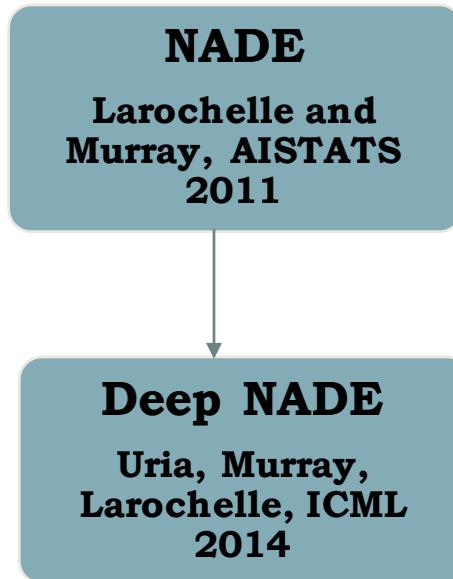


Related Work: CF methods

- Matrix Factorization based CF
 - SVD, SVD++
 - BiasMF (Koren et al., 2009)
 - PMF(Mnih & Salakhutdinov, 2007), BPMF(Salakhutdinov & Mnih,2008)
 - Poisson MF (Gopalan et al., 2013,2014)
 - LLORMA (Lee et al., 2015)
 - ...
- Neural Network based CF
 - RBM-CF (Salakhutdinov, Mnih and Hinton, 2007)
 - AutoRec (Sedhain et al., 2015)

Related Work: NADE based Models

- NADE (Larochelle and Murray 2011):
 - A tractable distribution estimator
 - For Fixed Length Binary Vectors
 - Tied parameters for networks
 - $O(DH)$
 - Practical for High dimensional data



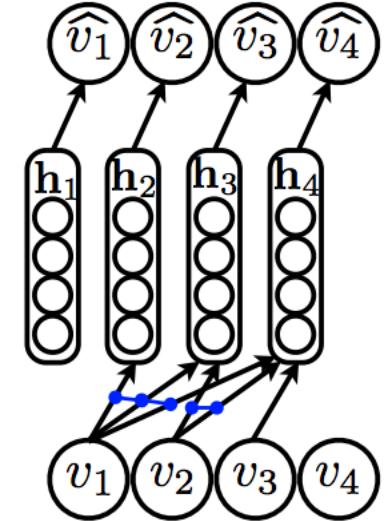
Real NADE
Uria, Murray, Larochelle, NIPS 2013

DocNADE
Larochelle and Lauly NIPS 2012

Fixation NADE
Zheng, Zemel, Zhang, Larochelle, IJCV 2014

Supervised DocNADE
Zheng, Zhang and Larochelle, CVPR, 2014

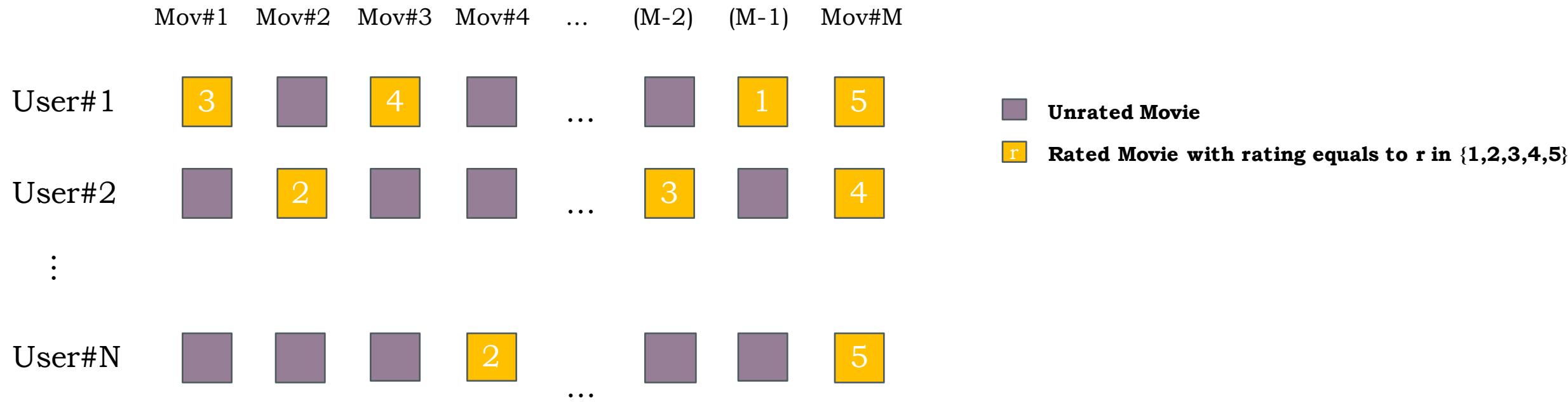
Deep DocNADE
Zheng, Zhang, Larochelle, tPAMI, 2015



NADE model

Source: Larochelle and Murray 2011

CF-NADE: The basic model



For User u we denote it as $\mathbf{r}^u = (r_{m_{o_1}}^u, r_{m_{o_2}}^u, \dots, r_{m_{o_D}}^u)$

CF-NADE tries to model $p(\mathbf{r}) = \prod_{i=1}^D p(r_{m_{o_i}} | \mathbf{r}_{m_{o_{<i}}})$

CF-NADE: The basic model

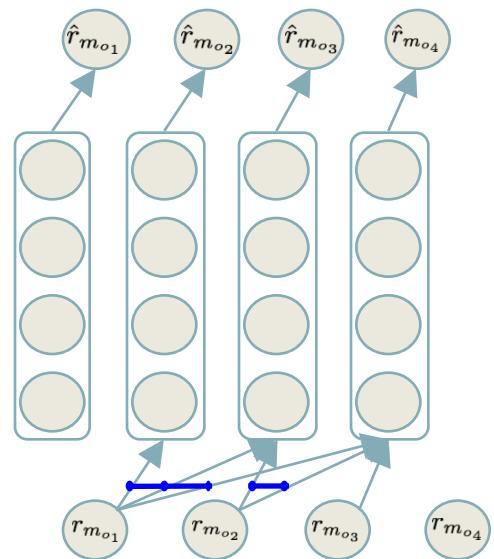
$$p(\mathbf{r}) = \prod_{i=1}^D p(r_{m_{o_i}} | \mathbf{r}_{m_{o_{<i}}})$$

One Option for the conditionals

$$p(r_{m_{o_i}} = k | \mathbf{r}_{m_{o_{<i}}}) = \frac{\exp(s_{m_{o_i}}^k(\mathbf{r}_{m_{o_{<i}}}))}{\sum_{k'=1}^K \exp(s_{m_{o_i}}^{k'}(\mathbf{r}_{m_{o_{<i}}}))}$$

$$s_{m_{o_i}}^k(\mathbf{r}_{m_{o_{<i}}}) = b_{m_{o_i}}^k + \mathbf{V}_{m_{o_i},:}^k \mathbf{h}(\mathbf{r}_{m_{o_{<i}}})$$

$$\mathbf{h}(\mathbf{r}_{m_{o_{<i}}}) = \mathbf{g}\left(\mathbf{c} + \sum_{j < i} \mathbf{W}_{:,m_{o_j}}^{r_{m_{o_j}}}\right)$$



About the order of the ratings

1. Should be predefined, e.g. *Timestamps*
2. A random order works well in practice
 - Different order is an different instantiation of CF-NADE for the same user
 - Key to extend CF-NADE to a deep model

CF-NADE: The basic model

hulu

Training Objective

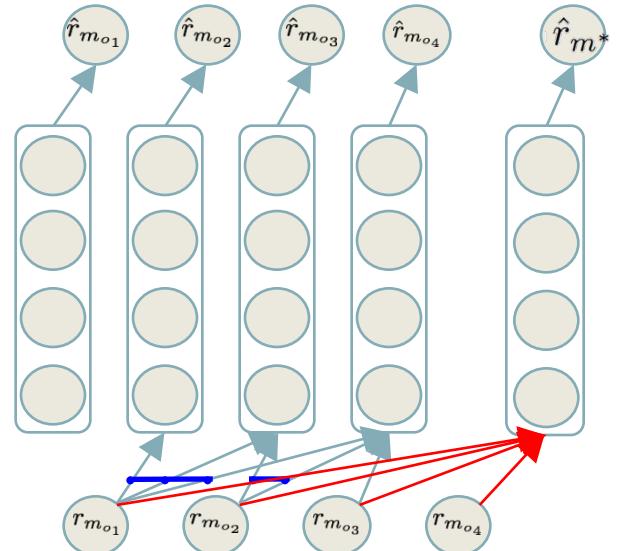
$$-\log p(\mathbf{r}) = -\sum_{i=1}^D \log p(r_{m_o_i} | \mathbf{r}_{m_o_{<i}})$$

Testing Phase

$$\hat{r}_{m^*} = \mathbb{E}_{p(r_{m^*} = k | \mathbf{r})} [k]$$

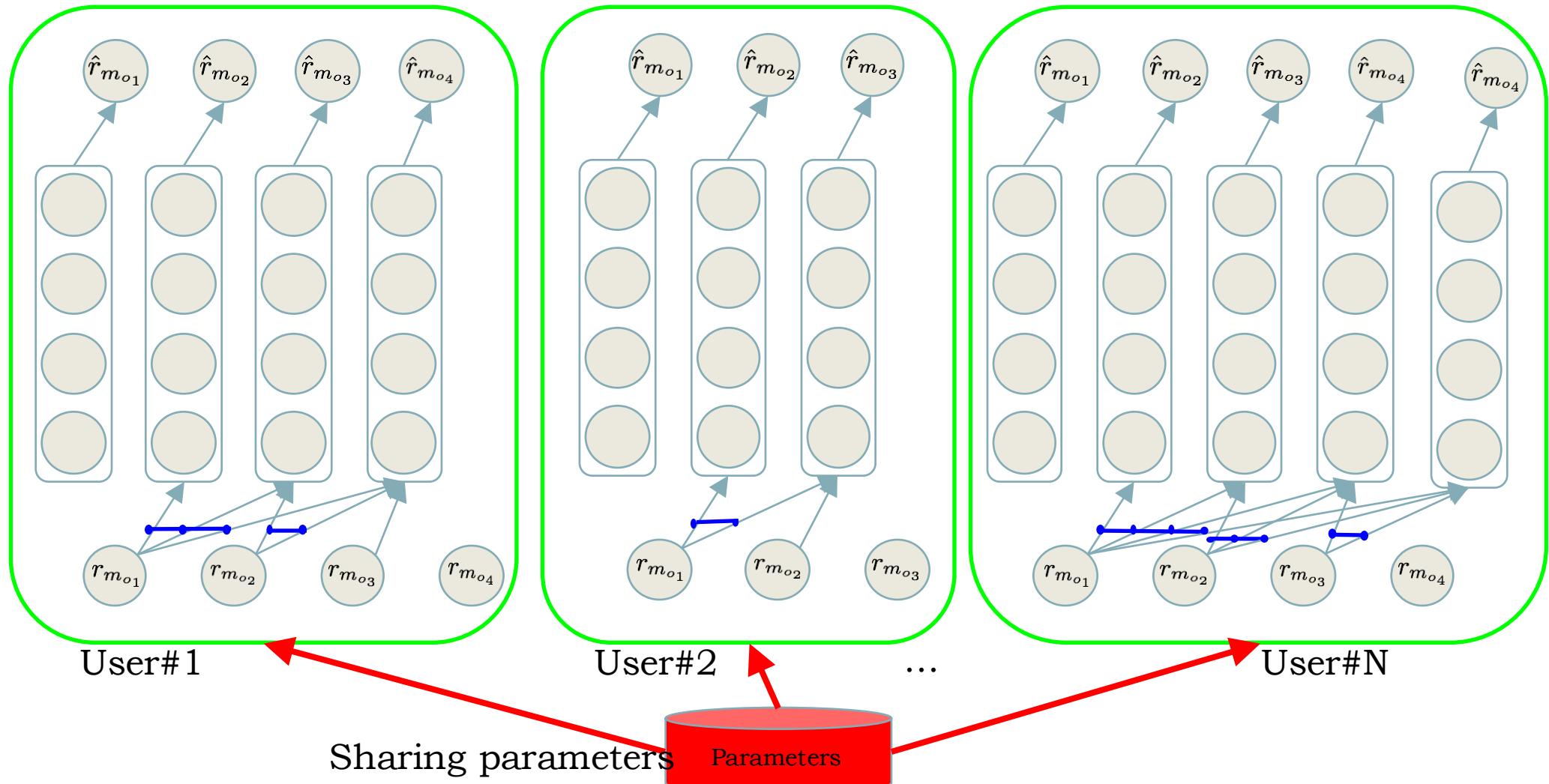
where $s_{m^*}^k(\mathbf{r}) = b_{m^*}^k + \mathbf{V}_{m^*,:}^k \mathbf{h}(\mathbf{r})$

$$\mathbf{h}(\mathbf{r}) = \mathbf{g} \left(\mathbf{c} + \sum_{j=1}^D \mathbf{W}_{:,m_{o_j}}^{r_{m_{o_j}}} \right)$$



A CF-NADE for each user

hulu



Sharing Parameters



The basic CF-NADE model

- different parameters for different ratings
- More parameters, less Optimization

$$\mathbf{h}(\mathbf{r}_{m_{o < i}}) = \mathbf{g}\left(\mathbf{c} + \sum_{j < i} \mathbf{W}_{:, m_{o_j}}^{r_{m_{o_j}}}\right)$$
$$s_{m_{o_i}}^k(\mathbf{r}_{m_{o < i}}) = b_{m_{o_i}}^k + \mathbf{V}_{m_{o_i}, :}^k \mathbf{h}(\mathbf{r}_{m_{o < i}})$$

Sharing Parameters between Ratings

$$\mathbf{h}(\mathbf{r}_{m_{o < i}}) = \mathbf{g}\left(\mathbf{c} + \sum_{j < i} \mathbf{W}_{:, m_{o_j}}^{r_{m_{o_j}}}\right)$$

Sharing 

$$\mathbf{h}(\mathbf{r}_{m_{o < i}}) = \mathbf{g}\left(\mathbf{c} + \sum_{j < i} \sum_{k=1}^{r_{m_{o_j}}} \mathbf{W}_{:, m_{o_j}}^k\right)$$

$$s_{m_{o_i}}^k(\mathbf{r}_{m_{o < i}}) = b_{m_{o_i}}^k + \mathbf{V}_{m_{o_i}, :}^k \mathbf{h}(\mathbf{r}_{m_{o < i}})$$

Sharing 

$$s_{m_{o_i}}^k(\mathbf{r}_{m_{o < i}}) = \sum_{j \leq k} \left(b_{m_{o_i}}^j + \mathbf{V}_{m_{o_i}, :}^j \mathbf{h}(\mathbf{r}_{m_{o < i}}) \right)$$

Better Scalability



Too many parameters...

$$\mathbf{W}^k \in \mathbb{R}^{H \times M} \quad \mathbf{V}^k \in \mathbb{R}^{M \times H} \quad k \in \{1, 2, \dots, K\}$$

Netflix dataset, $M=17700$, $H=500$, $K=5$, then **89** million free parameters

Weight Decay and Dropout is OK, but...

Solution: Factorize \mathbf{W} , \mathbf{V} by a product of 2 low rank matrices

$$W_{i,m}^k = \sum_{j=1}^J B_{i,j} A_{j,m}^k \quad V_{m,i}^k = \sum_{j=1}^J P_{m,j}^k Q_{j,i} \quad J \ll H \quad J \ll M$$

e.g. $J=50$, $M=17700$, $H=500$, $K=5$, then only **9** million free parameters for Netflix dataset

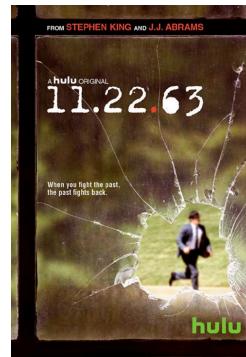
Ordinal Cost for CF



Ordinal Nature of Rating data:



I give this show 4-star score!



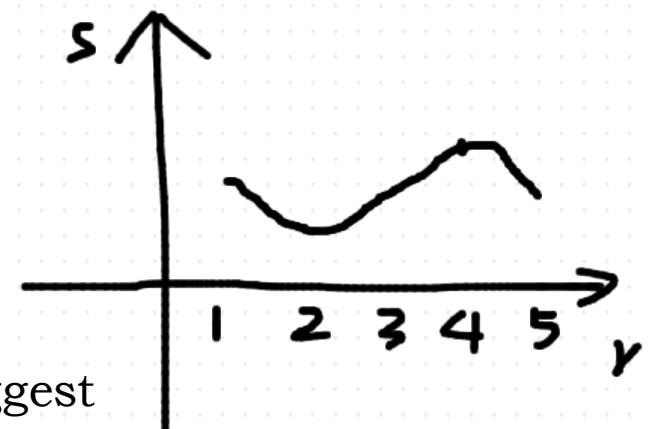
Rankings of Ratings

$$\mathbf{y}^{\text{down}} = (k, k-1, \dots, 1)$$

$$\mathbf{y}^{\text{up}} = (k, k+1, \dots, K)$$

\mathcal{C}_{reg} Regular Cost: Softmax

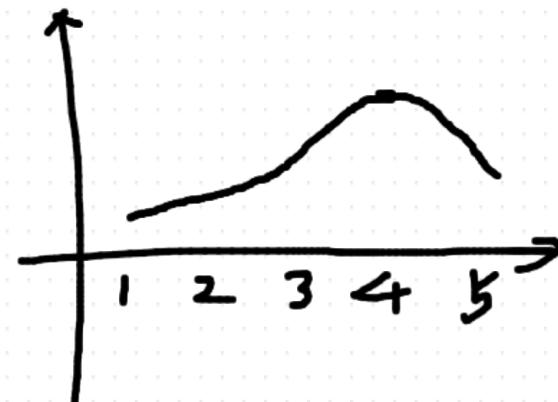
Score of rating k is the biggest



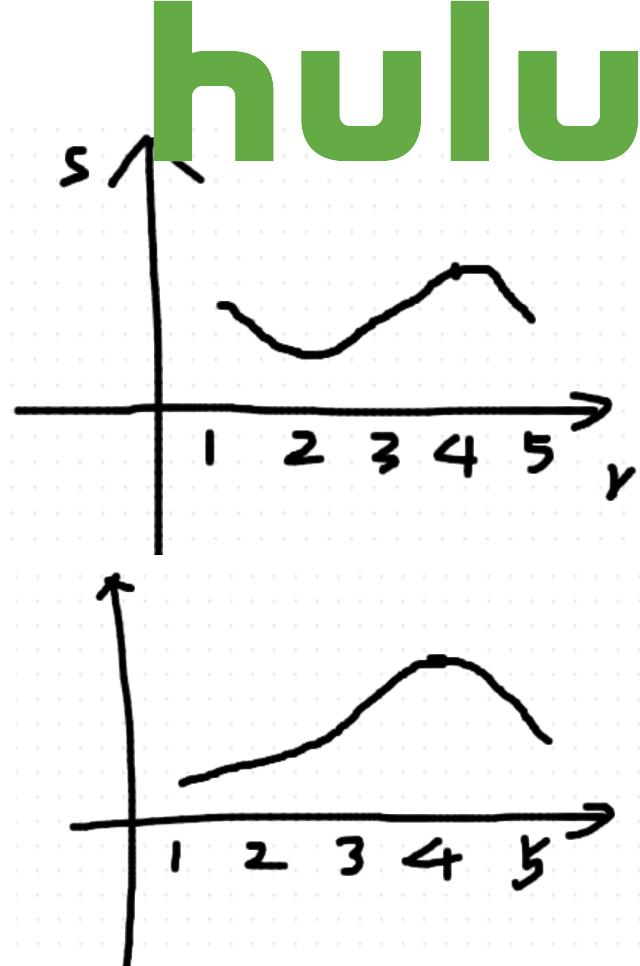
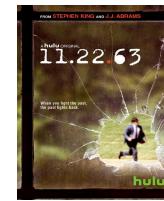
\mathcal{C}_{ord} Ordinal Cost: Ranking loss

$$k \succ k-1 \succ \dots \succ 1$$

$$k \succ k+1 \succ \dots \succ K$$

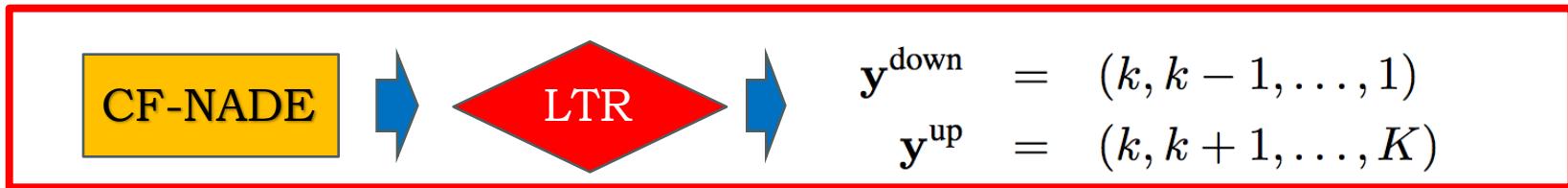


Ordinal Cost for CF



Ordinal Cost

$$p(r_{m_{o_i}} = k | \mathbf{r}_{m_{o_{<i}}}) = \prod_{j=k}^1 \frac{\exp(s_{m_{o_i}}^j)}{\sum_{t=1}^j \exp(s_{m_{o_i}}^t)} \prod_{j=k}^K \frac{\exp(s_{m_{o_i}}^j)}{\sum_{t=j}^K \exp(s_{m_{o_i}}^t)}$$



Hybrid Cost

$$\mathcal{C}_{\text{hybrid}} = (1 - \lambda)\mathcal{C}_{\text{reg}} + \lambda\mathcal{C}_{\text{ord}}$$

Deep CF-NADE

About the orderings of rated items

$$-\log p(\mathbf{r}) = -\sum_{i=1}^D \log p(r_{m_{o_i}} | \mathbf{r}_{m_{o_{<i}}})$$

Training over all possible orderings

$$\mathcal{C} = \mathbb{E}_{o \in \mathcal{O}} \sum_{i=1}^D -\log p(r_{m_{o_i}} | \mathbf{r}_{m_{o_{<i}}}, o)$$

- Given a context $\mathbf{r}_{m_{o_{<i}}}$, CF-NADE be equally good at modeling $\mathbf{r}_{m_{o_{\geq i}}}$

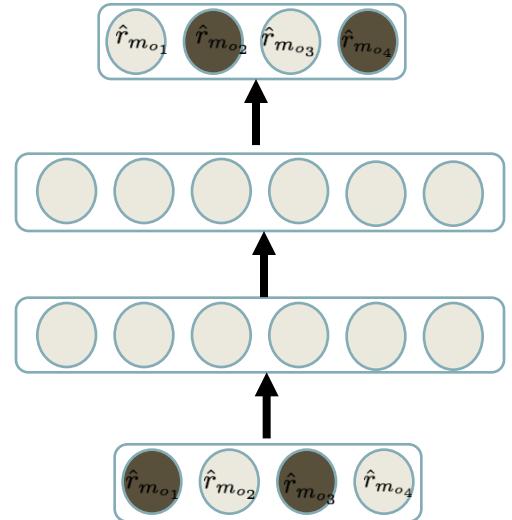
Cost Function is rewritten as

$$\mathcal{C} = \frac{D}{D-i+1} \sum_{j \geq i} -\log p(r_{m_{o_j}} | \mathbf{r}_{m_{o_{<i}}})$$

- only relies on $\mathbf{r}_{m_{o_{<i}}}$
- More hidden layers can be added as regular multiple layer neural networks

$$\mathbf{h}^{(l)}(\mathbf{r}_{m_{o_{<i}}}) = \mathbf{g}(\mathbf{c}^{(l)} + \mathbf{W}^{(l)} \mathbf{h}^{(l-1)}(\mathbf{r}_{m_{o_{<i}}}))$$

Complexity: $O(K\hat{D}H + H^2L)$



Experiments

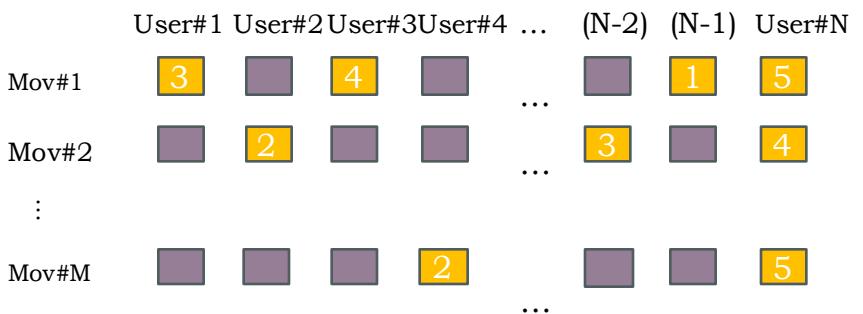
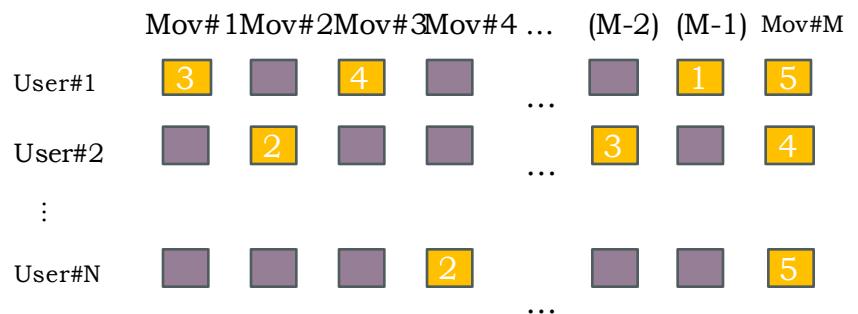
hulu

dataset	#Users	#Items	#Scales	#Ratings	type
MovieLen 1M	6040	3952	5	10^6	I-CF-NADE
MovieLen 10M	71567	10681	10	10^7	U-CF-NADE
Netflix	480189	17770	5	10^8	U-CF-NADE

Train: 90% ratings
Test: 10% ratings

Metric $RMSE = \sqrt{\frac{\sum_{i=1}^S (r_i - \tilde{r}_i)^2}{S}}$

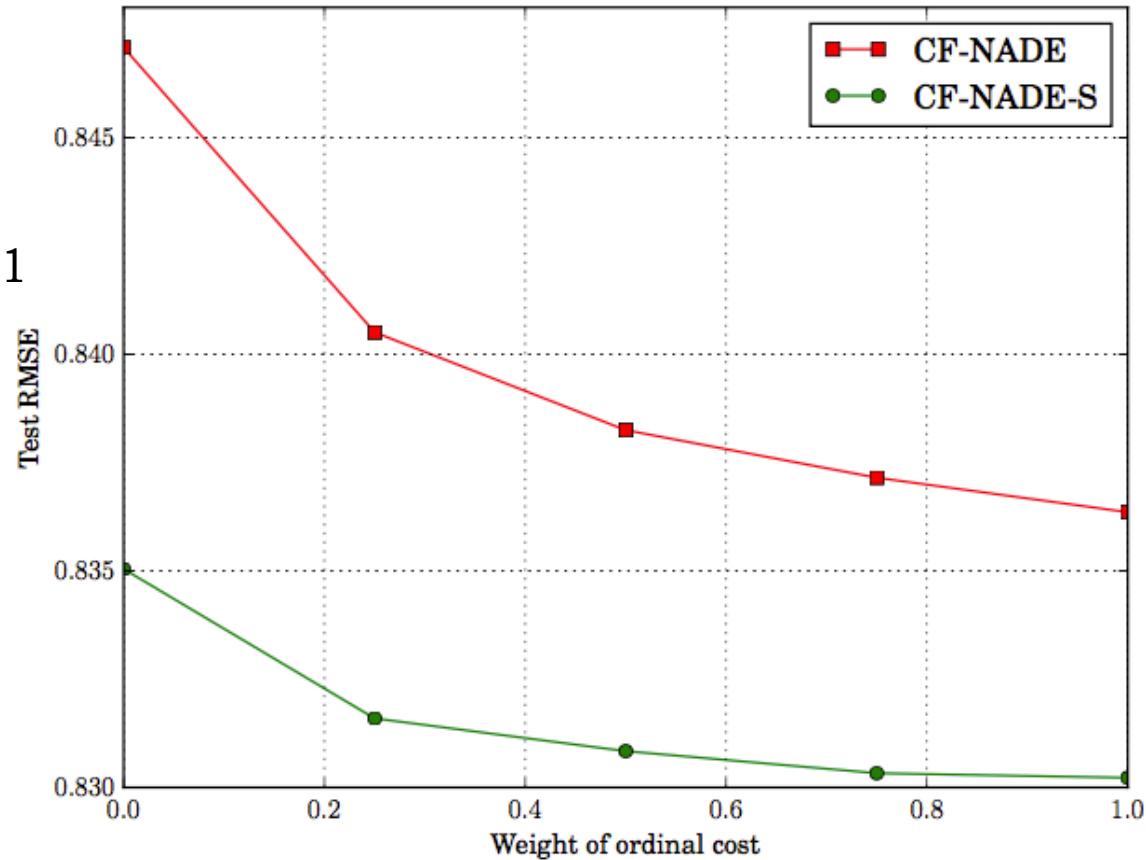
User based CF-NADE (U-CF-NADE) VS Item based CF-NADE (I-CF-NADE)



Experiments on MovieLens 1M

Impact of Ordinal Cost

- Hidden size = 500
- Weight decay = 0.015
- Adam Optimizer, learning rate 0.001



Experiments on MovieLens 1M

hulu

Compare with other baselines

- Hidden Size {250, 500, 750}
- Weight decay {0.015, 0.02}
- Learning rate {0.001, 0.0005, 0.0002}

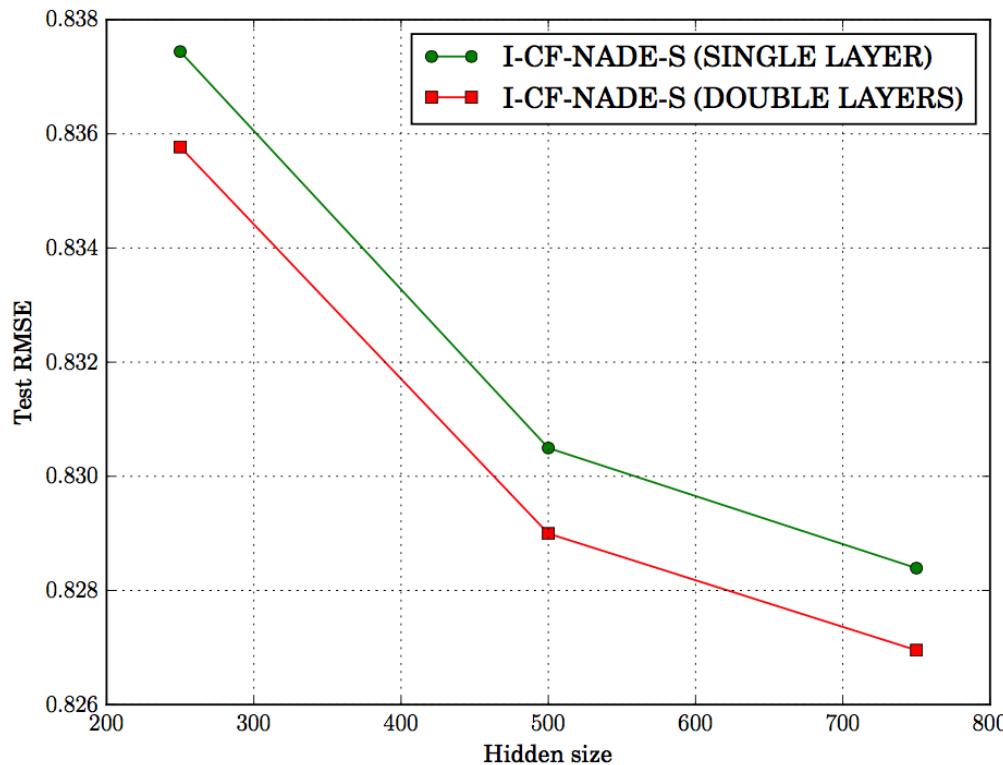


Table 1. Test RMSE of different models on MovieLens 1M.

METHOD	TEST RMSE
PMF†	0.883
U-RBM*	0.881
U-AUTOREC (SEDHAIN ET AL., 2015)	0.874
LLORMA-GLOBAL (LEE ET AL., 2013)	0.865
I-RBM*	0.854
BIASMF*	0.845
NNMF (DZIUGAITE & ROY, 2015)	0.843
LLORMA-LOCAL (LEE ET AL., 2013)	0.833
I-AUTOREC (SEDHAIN ET AL., 2015)	0.831
U-CF-NADE-S (SINGLE LAYER)	0.850
U-CF-NADE-S (2 LAYERS)	0.845
I-CF-NADE-S (SINGLE LAYER)	0.830
I-CF-NADE-S (2 LAYERS)	0.829

†: Taken from (Dziugaite & Roy, 2015).

*: Taken from (Sedhain et al., 2015).

Experiments on MovieLens 10M



Table 2. Test RMSE of different models on MovieLens 10M.

Configurations

- Hidden Size 500
- Factorized U-CF-NADE, J=50
- Learning rate 0.0005
- Weight decay 0.015

METHOD	TEST RMSE
U-AUTOREC (SEDHAIN ET AL., 2015)	0.867
I-RBM†	0.825
U-RBM†	0.823
LLORMA-GLOBAL (LEE ET AL., 2013)	0.822
BIASMF†	0.803
LLORMA-LOCAL (LEE ET AL., 2013)	0.782
I-AUTOREC (SEDHAIN ET AL., 2015)	0.782
U-CF-NADE-S (SINGLE LAYER)	0.772
U-CF-NADE-S (2 LAYERS)	0.771

†: Taken from (Sedhain et al., 2015).

Experiments on Netflix Dataset

Configurations

- Hidden Size 500
- Factorized U-CF-NADE, J=50
- Learning rate 0.0005
- Weight decay 0.001

Table 3. Test RMSE of different models on Netflix dataset.

METHODS	TEST RMSE
LLORMA-GLOBAL (LEE ET AL., 2013)	0.874
U-RBM†	0.845
BIASMF†	0.844
LLORMA-LOCAL (LEE ET AL., 2013)	0.834
I-AUTOREC (SEDHAIN ET AL., 2015)	0.823
U-CF-NADE-S (SINGLE LAYER)	0.804
U-CF-NADE-S (2 LAYERS)	0.803

†: Taken from (Sedhain et al., 2015).

Complexity of CF-NADE

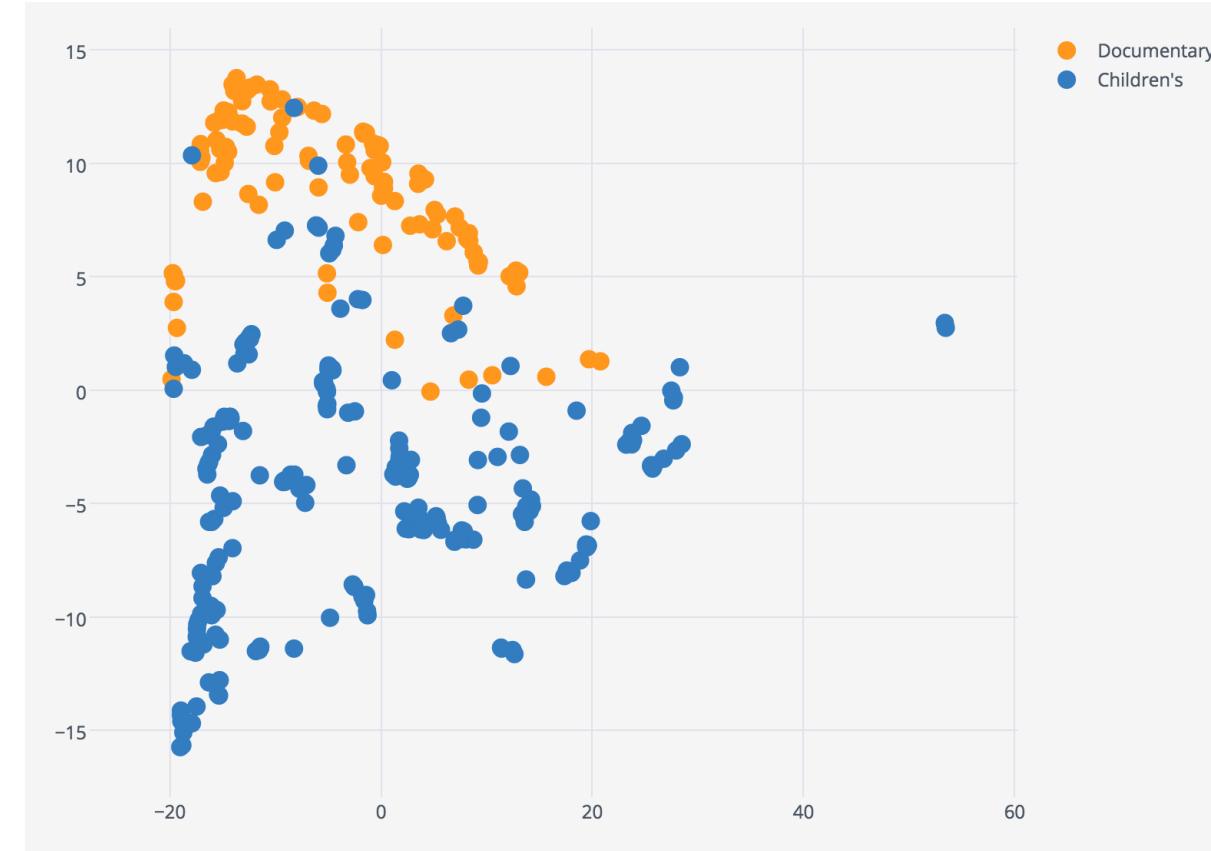
Table 4. Complexity of CF-NADE on different benchmarks

Datasets	# Layers	Train Time (second)	Test Time (second)	#Params (million)
MLens1M	1 layer	3.09	0.65	30.2
	2 layers	3.11	0.68	30.48
MLens10M	1 layer	134.76	31.72	10.78
	2 layers	135.62	32.73	10.98
Netflix	1 layer	1057.81	239.78	9.02
	2 layers	1064.33	243.79	9.19

Note that there are some rooms to speed up the implementation

The complexity is $O(K\hat{D}H + H^2L)$

Visualization of the Learned Model

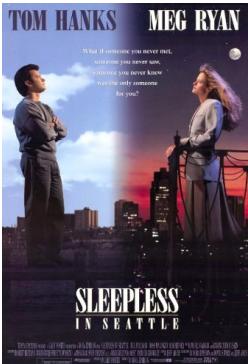


tSNE of the representations of movies on MovieLens 1M

Input Movies



Star Trek VI:
The Undiscovered Country



Sleepless in Seattle

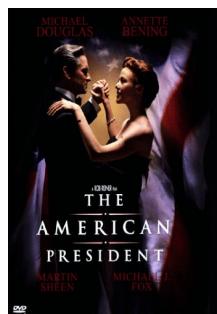


The Lion King

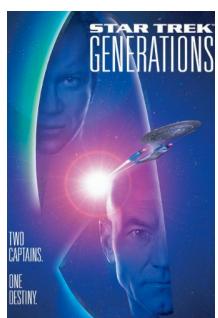
Top 6 most similar movies by CF-NADE on MovieLens 1M



Star Trek III:
The Search for Spock



The American President



Star Trek:
Generations



Pretty Woman



Star Trek:
Insurrection



Notting Hill



Star Trek:
First Contact



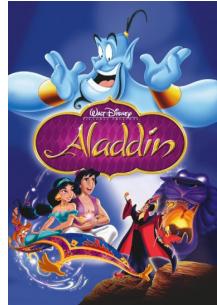
A League of Their Own



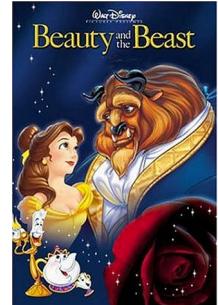
Star Trek IV:
The Voyage Home



Ghost



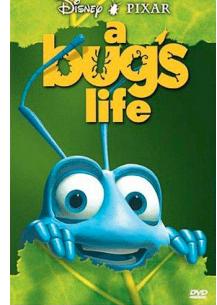
Aladdin



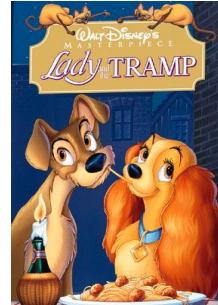
Beauty and
the Beast



The Little
Mermaid



A Bug's Life



Lady and
the Tramp

Conclusion

- CF-NADE:
 - An Efficient and Powerful architecture for CF tasks
 - Sharing parameters is a good way to improve the performance
 - Better scalability by factorization
 - Can be extended to a deep version
- Ordinal cost for rating data
 - Learning to rank
- Meaningful representations
- Source code at: <https://github.com/Ian09/CF-NADE>



Thanks! Q&A