



Discrete Collaborative Filtering

Hanwang Zhang¹, Fumin Shen², Wei Liu³,
Xiangnan He¹, Huanbo Luan⁴, Tat-Seng Chua¹

Presented by Xiangnan He

1. National University of Singapore
2. University of Electronic Science and Technology of China
3. Tencent Research
4. Tsinghua University

19 July 2016

Highlights of Discrete Collaborative Filtering

Online Recommendation

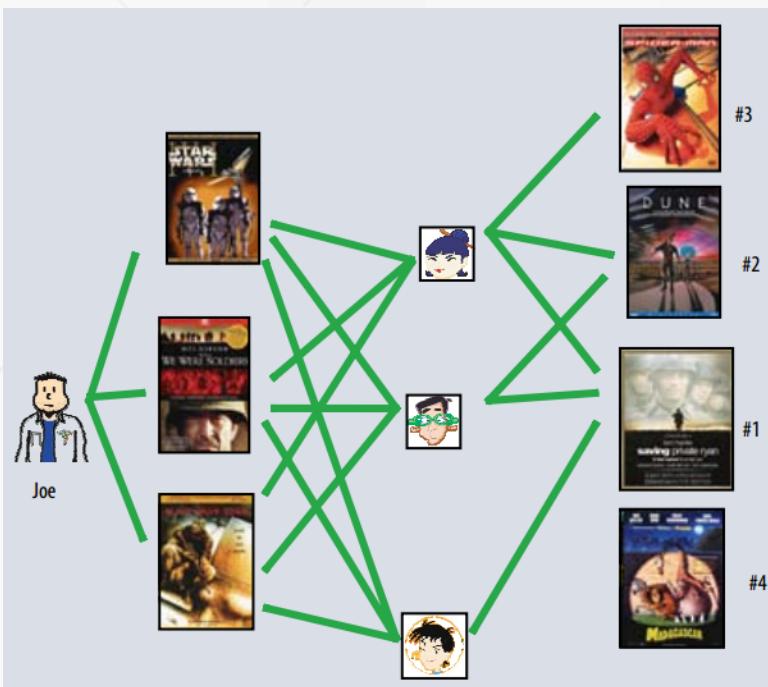
- An **Efficient** Recommender System
- Latent Model: **Binary** Representation for Users and Items
- Recommendation as **Search** with Binary Codes

Offline Training

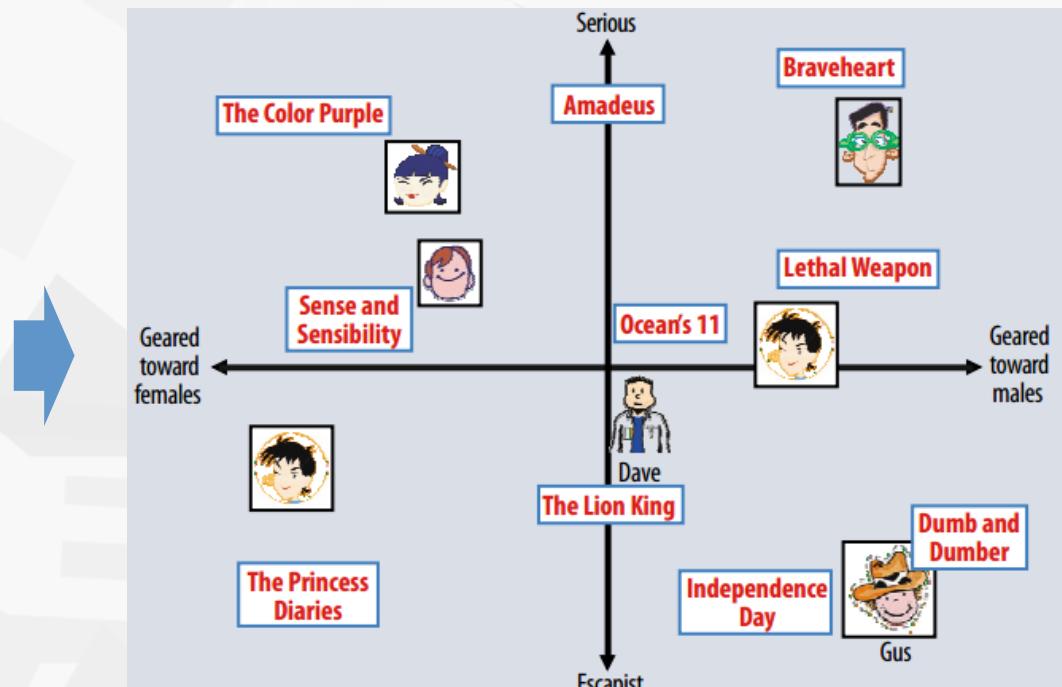
- **End-to-end** binary optimization
- **Balanced** and **Decorrelated** Constraint
- Small **SVD** + **Discrete** Coordinate Descent

Collaborative Filtering

Latent Factor Approach [Koren et al. 2009]



User-Item Matrix



Latent Space

Efficient CF: Hashing Users & Items

Recommendation is *Search*

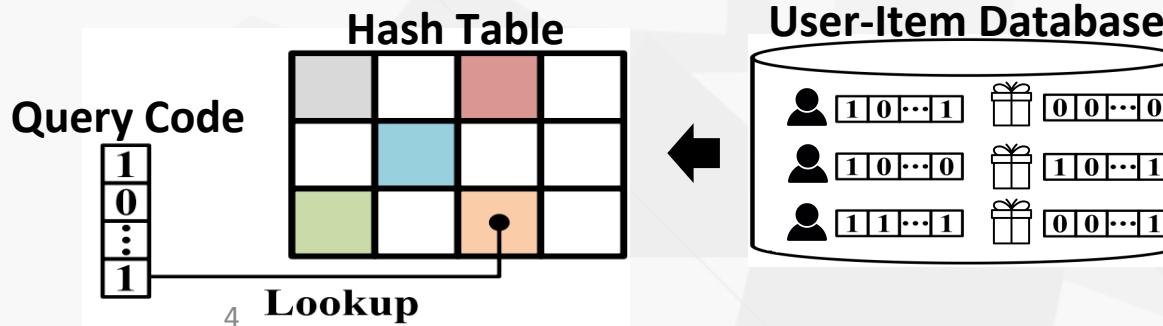
Ranking by <user vector, item vector>

Search in Euclidean space is *slow*

Requires **float** operations & **linear** scan of the data

Search in *Hamming Space* is *fast*.

Only requires **XOR** operation & **constant-time** lookup



Existing Hashing for CF: Two-Stage Approach

[Zhang et al, SIGIR'14; Zhou et al, KDD'12]

- Stage 1: Relaxed Real-Valued Problem
 $\{B, D\} \leftarrow \text{Continuous CF Methods}$
- Stage 2: Binarization
 $B \leftarrow \text{sgn}(B), \quad D \leftarrow \text{sgn}(D)$



Code learning and CF are isolated

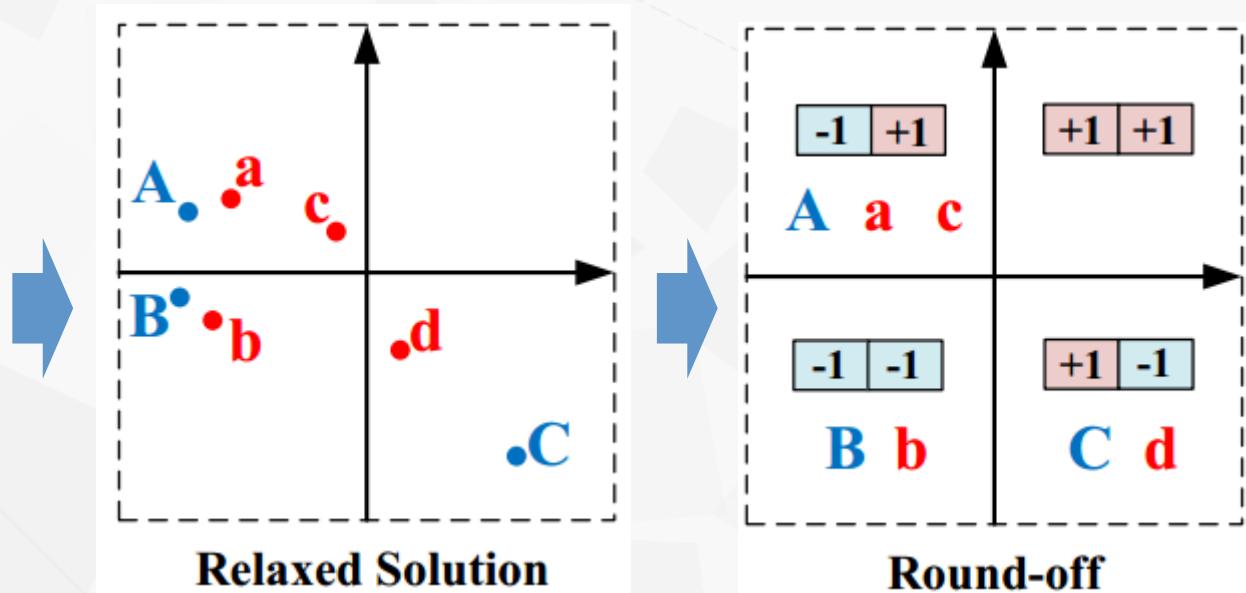


Quantization Loss

Quantization Loss

User-Item Matrix

	a	b	c	d
A	1	.8	.4	?
B	.8	.8	?	.4
C	.2	?	.2	.4

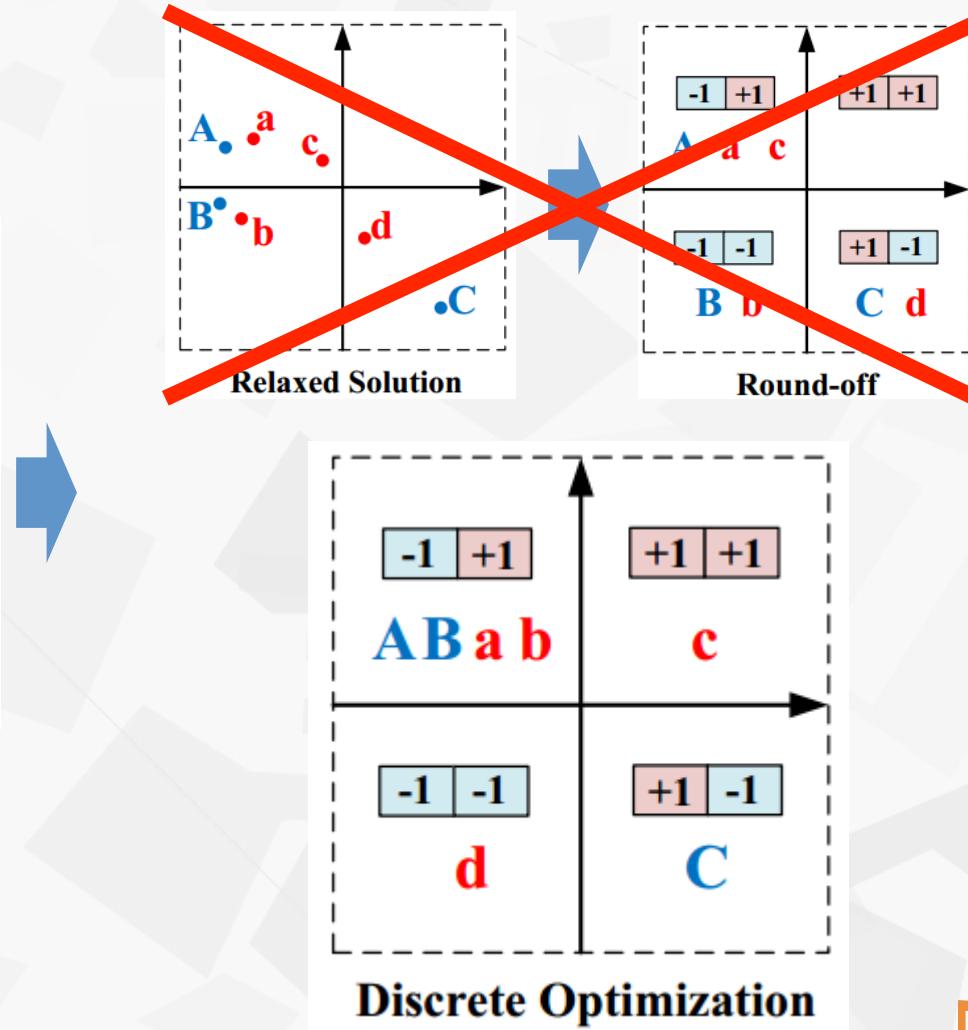


1. **A,B,a,b** are close but they are separated into different quadrants
2. **C, d** should be far but they are assigned to the same quadrant

Tackling Quantization Loss

User-Item Matrix

	a	b	c	d
A	1	.8	.4	?
B	.8	.8	?	.4
C	.2	?	.2	.4



Intuitive Solution: Binary Constraints

$$\operatorname{argmin}_{\mathbf{B}, \mathbf{D}} \sum_{i,j \in \mathcal{V}} \left(S_{ij} - \mathbf{b}_i^T \mathbf{d}_j \right)^2$$

Rating Prediction

Observed rating User code Item code

$\mathbf{B} \in \{\pm 1\}^{r \times m}, \mathbf{D} \in \{\pm 1\}^{r \times n}$

Binary Constraint

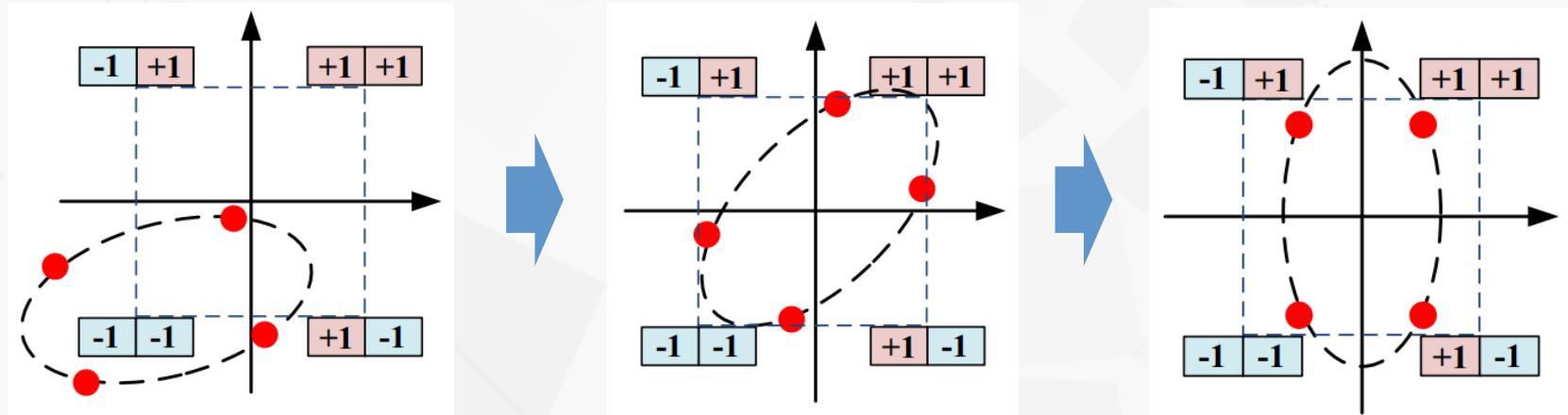
The diagram illustrates the components of the rating prediction formula. Observed rating, User code, and Item code are shown as boxes with arrows pointing to the terms S_{ij} , \mathbf{b}_i^T , and \mathbf{d}_j respectively in the formula. A bracket under the formula is labeled "Rating Prediction". To the right, a bracket under \mathbf{B} and \mathbf{D} is labeled "Binary Constraint".

However, it may lead to non-informative codes, e.g.:

1. **Unbalanced Codes** → each bit should have split the dataset evenly
2. **Correlated Codes** → each bit should be as independent as possible

Balanced and Decorrelated Constraint

Illustration of the effectiveness of the two constraints in DCF



Without any constraints:
3 points are (-1, -1) and 1
point is (+1, -1), which is
not discriminative.

Balanced:
Separated in the
1st & 3rd quadrant

Decorrelated:
Well separated

Solution with the Two Constraints

$$\operatorname{argmin}_{\mathbf{B}, \mathbf{D}} \sum_{i,j \in \mathcal{V}} \left(S_{ij} - \mathbf{b}_i^T \mathbf{d}_j \right)^2,$$

$$s.t. \quad \mathbf{B} \in \{\pm 1\}^{r \times m}, \mathbf{D} \in \{\pm 1\}^{r \times n}$$

$$\underbrace{\mathbf{B}\mathbf{1} = 0, \mathbf{D}\mathbf{1} = 0}_{\text{Balanced Partition}}, \quad \underbrace{\mathbf{B}\mathbf{B}^T = m\mathbf{I}, \mathbf{D}\mathbf{D}^T = n\mathbf{I}}_{\text{Decorrelation}}$$



However, the hard constraints of **zero-mean** and **orthogonality** may not be satisfied in Hamming space!

Our DCF Formulation

Objective function:

$$\sum_{i,j \in \mathbb{V}} (S_{ij} - \mathbf{b}_i^T \mathbf{d}_j)^2 + 2\alpha \|\mathbf{B} - \mathbf{X}\|^2 + 2\beta \|\mathbf{D} - \mathbf{Y}\|^2$$

Binary Constraint

$$\mathbf{B} \in \{\pm 1\}^{r \times m}, \mathbf{D} \in \{\pm 1\}^{r \times n}$$

Delegate Code Quality Constraint

$$\mathbf{X}\mathbf{1} = \mathbf{0}, \mathbf{Y}\mathbf{1} = \mathbf{0}, \mathbf{X}\mathbf{X}^T = m\mathbf{I}, \mathbf{Y}\mathbf{Y}^T = n\mathbf{I}$$

Balanced Constraint

Decorrelated Constraint

Mixed-Integer Programming NP-Hard [Hastad 2001]

Our Solution: Alternating Optimization

Alternative Procedure

- B-Subproblem

$$\operatorname{argmin}_{\mathbf{B}} \sum_{i,j \in \mathcal{V}} (S_{ij} - \mathbf{b}_i^T \mathbf{d}_j)^2 - 2\alpha tr(\mathbf{B}^T \mathbf{X}) \text{ s.t., } \mathbf{B} \in \{\pm 1\}^{r \times m}$$

- D-Subproblem

$$\operatorname{argmin}_{\mathbf{D}} \sum_{i,j \in \mathcal{V}} (S_{ij} - \mathbf{b}_i^T \mathbf{d}_j)^2 - 2\beta tr(\mathbf{D}^T \mathbf{Y}) \text{ s.t., } \mathbf{D} \in \{\pm 1\}^{r \times n}$$

- X-Subproblem

$$\operatorname{argmin}_{\mathbf{X}} - 2\alpha tr(\mathbf{B}^T \mathbf{X}) \text{ s.t., } \mathbf{X}\mathbf{1} = 0, \mathbf{X}\mathbf{X}^T = m\mathbf{I}$$

- Y-Subproblem

$$\operatorname{argmin}_{\mathbf{Y}} - 2\beta tr(\mathbf{D}^T \mathbf{Y}) \text{ s.t., } \mathbf{Y}\mathbf{1} = 0, \mathbf{Y}\mathbf{Y}^T = n\mathbf{I}$$

B-Subproblem for Binary Codes

Objective Function

$$\underset{\mathbf{B}}{\operatorname{argmin}} \sum_{i,j \in \mathcal{V}} (S_{ij} - \mathbf{b}_i^T \mathbf{d}_j)^2 - 2\alpha t r(\mathbf{B}^T \mathbf{X}) \text{ s.t., } \mathbf{B} \in \{\pm 1\}^{r \times m}$$

For each user code \mathbf{b}_i , optimize **bit by bit**

for $i=1$ **to** m **do** Parallel **for** loop over m users

repeat Usually converges in 5 iterations

for $k=1$ **to** r **do** **for** loop over r bits

$$\hat{b}_{ik} \leftarrow \sum_{j \in \mathcal{V}_i} (S_{ij} - \mathbf{d}_{j\bar{k}}^T \mathbf{b}_{ik}) d_{jk} + \alpha x_{ik};$$

$$b_{ik} \leftarrow \operatorname{sgn} (K(\hat{b}_{ik}, b_{ik}));$$

end

until converge;

end

D-Subproblem can be solved in a similar way

B-Subproblem Complexity

$$\mathcal{O}(r^2 T_s |\mathcal{V}| / p)$$

#bits #bit-by-bit iterations #computing threads
#training ratings

The diagram illustrates the complexity formula $\mathcal{O}(r^2 T_s |\mathcal{V}| / p)$. Four red arrows point from labels above the formula to its components: '#bits' points to r^2 , '#bit-by-bit iterations' points to T_s , '#computing threads' points to p , and '#training ratings' points to $|\mathcal{V}|$.

Linear to data size!

X-Subproblem for Code Delegate

Objective Function

$$\underset{\mathbf{X}}{\operatorname{argmin}} - 2\alpha \operatorname{tr}(\mathbf{B}^T \mathbf{X}) \text{ s.t., } \mathbf{X}\mathbf{1} = 0, \mathbf{X}\mathbf{X}^T = m\mathbf{I}$$

Small SVD $r \times m$

$$([\mathbf{P}_b \hat{\mathbf{P}}_b], \mathbf{Q}_b) \leftarrow \text{SVD}(\bar{\mathbf{B}})$$

$r \times m$ row-centered user code matrix

Orthogonalization

$$\hat{\mathbf{Q}}_b \leftarrow \text{GramSchmidt}([\mathbf{Q}_b \mathbf{1}])$$

$$\mathbf{X} \leftarrow \sqrt{m} [\mathbf{P}_b \hat{\mathbf{P}}_b] [\mathbf{Q}_b \hat{\mathbf{Q}}_b]^T$$

Y-Subproblem can be solved in a similar way

X-Subproblem Complexity

$$\mathcal{O}(r^2 m)$$

#bits #users

The diagram features a mathematical expression $\mathcal{O}(r^2 m)$ enclosed in a white rectangular box. Two red arrows originate from the terms r^2 and m in the expression and point towards the labels "#bits" and "#users" respectively, indicating that the complexity is proportional to the product of the number of bits and the number of users.

Linear to data size!

Summary

- Recommendation is *search*
- We can accelerate search by *hashing*
- Unlike previous erroneous *two-stage* hashing, *DCF* is an *end-to-end* hashing method
- Fast *O(n)* *discrete optimization* for DCF

Evaluations

- Dataset (filtering threshold at 10):

Table 1: Statistics of datasets in evaluation.

Dataset	Rating#	User#	Item#	Density
Yelp	696,865	25,677	25,815	0.11%
Amazon	5,057,936	146,469	189,474	0.02%
Netflix	100,480,507	480,189	17,770	1.18%

- Random split: 50% training and 50% testing.
- Metric: **NDCG@K**
- Search Protocol:
Hamming ranking or hash table lookup

Evaluation 1: Compared to state-of-the-art

- **MF:** Matrix Factorization [Koren et al 2009]
Classific MF, Euclidean space baseline
- **BCCF:** Binary Code learning for Collaborative Filtering
[Zhou&Zha, KDD 2012]
MF+balance+binarization
- **PPH:** Preference Preserving Hashing [Zhang et al. SIGIR 2014]
Cosine MF + norm&phase binarization
- **CH:** Collaborative Hashing [Liu et al. CVPR 2014]
Full SVD MF + balance + binarization

Evaluation 1

DCF is a new state-of-the-art

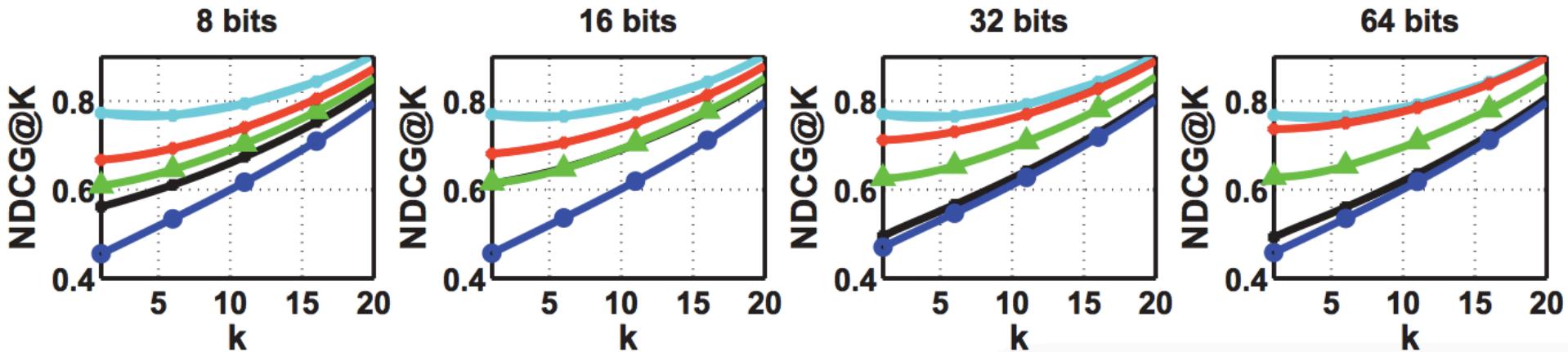
Performance of NDCG@10.

Protocol	Hamming Ranking		
Dataset	Yelp	Amazon	Netflix
BCCF(128)	0.623	0.827	0.633
PPH(128)	0.643	0.841	0.587
CH(128)	0.655	0.922*	0.693
DCF(8)	0.684*	0.890	0.730*

1. DCF learns compact and informative codes.
2. DCF's performance is most close to the real-valued MF.
3. End-to-end > Two stage

MF — BCCF — PPH — CH — DCF

Netflix



Evaluation 2

DCF generalizes well to unseen users

Training: full histories of 50% users

Testing: the other 50% users that have no histories in training

Evaluation: simulate online learning scenario.

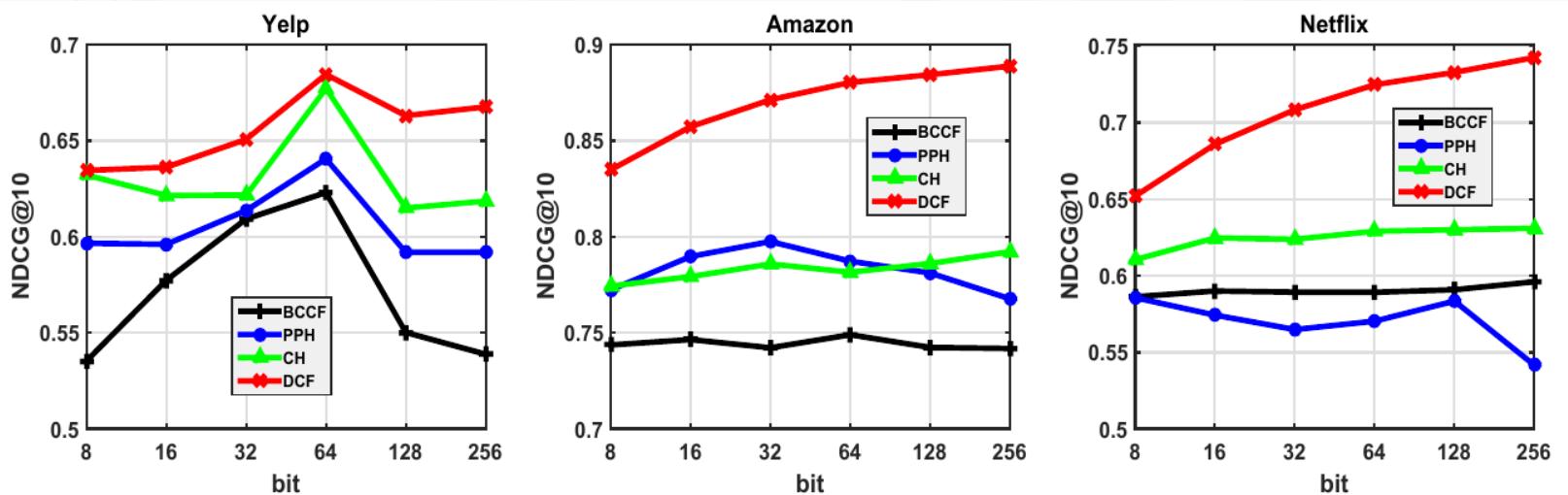


Figure 4: Recommendation performance (NDCG@10) on 50% held-out “new” users (RQ 2).

— MF — BCCF — PPH — CH — DCF

Evaluation 3

Balanced and Decorrelated Constraints are necessary

MF: original MF

MFB: MF+Binarization

DCFinit: the variant of DCF that discards the two constraints.

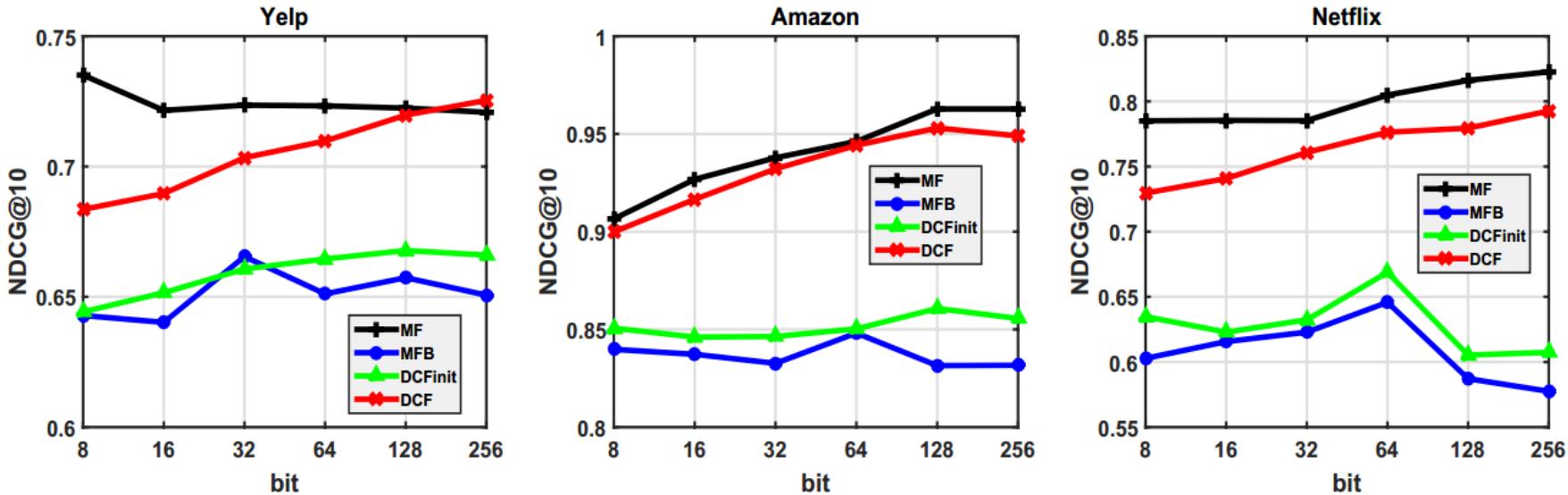


Figure 5: Recommendation performance (NDCG@10) of CF and DCF variants (RQ 3).

Conclusions

- Discrete Collaborative Filtering: an end-to-end hashing method for efficient CF
- A fast **algorithm** for DCF
- DCF is a general **framework**. It can be extended to any popular CF variants, such as SVD++ and factorization machines.

Thank you

Code available: <https://github.com/hanwangzhang/Discrete-Collaborative-Filtering>

