Discrete Factorization Machines for Fast Feature-based Recommendation

Han Liu¹, Xiangnan He², Fuli Feng², Liqiang Nie¹, Rui Liu³, Hanwang Zhang⁴

Shandong University
 National University of Singapore
 University of Electronic Science and Technology of China
 Nanyang Technological University

Motivation

Accurate Recommender System

Quality of Service & Profit of the Service Provider

side information



content-based : e.g.,
item descriptions



context-based : *e.g.*, when and where a purchase is made



session-based : e.g.,
recent browsing
history of users

user



item

Factorization Machines (FM)

FM is a score prediction function for a (user, item) pair feature x.



Motivation



Existing FM framework is not suitable for fast recommendation, especially for mobile users.

Discrete Factorization Machines





Balance Constraint: each bit should split the dataset evenly De-Correlation Constraint: each bit should be as independent as possible

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

Our DFM Formulation



Our Solution: Alternating Optimization



B-Subproblem for Binary Codes

Objective Function

 $\underset{\mathbf{B}}{\operatorname{arg\,min}} \sum_{(\mathbf{x},y)\in\mathcal{V}} (y-w_0 - \sum_{i=1}^n w_i x_i - \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{b}_i, \mathbf{b}_j \rangle x_i x_j)^2 - 2\beta tr(\mathbf{B}^T \mathbf{D}), \text{ s.t. } \mathbf{B} \in \{\pm 1\}^{k \times n}$

For each feature code **b***li*, optimize **bit by bit**

for
$$r = 1$$
 to n do for loop over n features
repeat
for $t = 1$ to k do for loop over k bits
 $\begin{vmatrix} \hat{b}_{rt} = \sum_{\mathcal{V}_r} (x_r \psi - x_r^2 \hat{\mathbf{x}}^T \mathbf{Z}_{\bar{t}} \mathbf{b}_{r\bar{t}}) \hat{\mathbf{x}}^T \mathbf{z}_t + \beta d_{rt}; \\ \hat{b}_{rt} \leftarrow \operatorname{sgn}(K(\hat{b}_{rt}, b_{rt})) \end{vmatrix}$
end
until converge;

D-Subproblem for Code Delegate



w-Subproblem for Bias



It is the standard multivariate linear regression problem, use Coordinate Descent algorithm



Experiment Settings

• Datasets:

Datasets	#users	#items	#ratings	Density
Yelp	13,679	12,922	640,143	0.36%
Amazon	35,151	33,195	1,732,060	0.15%

- Split: randomly split 50% training and 50% testing move items in the testing set that haven't occurred in the training set to the training set.
- Evaluation Protocol: rank the testing items of a user and evaluate the ranked list with NDCG@K

Compared to the state-of-the-art

- **libFM**: Factorization Machines with <u>libFM</u> [Rendle et al.,TIST'12] original implementation of FM
- DCF: <u>D</u>iscrete <u>C</u>ollaborative <u>F</u>iltering [Zhang et al., SIGIR'16] <u>CF+binarization+direct optimization</u>
- DCMF: <u>D</u>iscrete <u>Content-aware</u> <u>Matrix</u> Factorization [Lian et al.,KDD'17]

CF+binarization+direct optimization+constraint

 BCCF: <u>B</u>inary <u>C</u>ode learning for <u>C</u>ollaborative <u>F</u>iltering [Zhou&Zha,KDD'12]

MF+binarization+two-stage optimization

Performance Comparison



In figure, we show the recommendation performance (NDCG@1 to NDCG@10) of DFM and the baseline methods on the two datasets. The code length varies from 8 to 64.

Efficiency Study

Efficiency comparison between DFM and libFM regarding Testing Time Cost (TTC) on the two datasets.

тар							
Code Length	8	16	32	64			
libFM (TTC)	27.18	56.77	114.10	217.64			
DFM (TTC)	2.06	3.56	6.60	12.43			
Acceleration Ratio	13.19	15.95	17.29	17.51			

Voln

Amazon

Code Length	8	16	32	64
libFM (TTC)	177.03	357.46	716.83	1,414.67
DFM (TTC)	12.67	22.50	42.56	81.04
Acceleration Ratio	13.97	15.89	16.84	17.46



DFM is an operable solution for many large-scale Web service to reduce the computation cost of their recommender systems.

Conclusion & Future Work

- We propose DFM to enable fast feature-based recommendation.
- We develop an efficient algorithm to address the challenging optimization problem of DFM.
- We will extend binary technique to neural recommender models such as Neural FM.

Q&A

Thank you.



https://github.com/hanliu95/DFM