# Nonparametric Max-Margin Matrix Factorization for Collaborative Prediction



### Abstract

We present a **probabilistic** formulation of **max-margin matrix factori***zation* and build accordingly a *nonparametric* Bayesian model which automatically resolves the unknown number of latent factors.

Our work demonstrates a successful example that integrates *Bayesian nonparametrics* and *max-margin learning*, which are conventionally two separate paradigms and enjoy complementary advantages.

We develop an efficient *variational* algorithm for posterior inference, and our extensive empirical studies on large-scale MovieLens 1M and EachMovie data sets appear to justify the aforementioned dual advantages.

### Method Outline

1.  $PM^{3}F$ : The first key step is a general probabilistic formulation of the standard max-margin matrix factorization  $(M^{3}F)$ , which is based on the maximum entropy discrimination (MED) principle.

2.  $iPM^{3}F$ : We can then extend it to a nonparametric model, which in theory has an unbounded number of latent factors. To avoid over-fitting, we impose a sparsity-inducing Indian buffet process (IBP) prior on the latent coefficient matrix, selecting only an appropriate number of active factors.

3. **iBPM<sup>3</sup>F**: To make one step further towards a hierarchical





To solve backgro

```
For a tr
functior
```

and pre MED pr margin subsum

In the c

**From M<sup>3</sup>E to iBPM<sup>3</sup>E**  
preference matrix for N users and M items, which is only  
y observed: 
$$Y \in \mathcal{Y}^{\times \times \times}$$
,  $\mathcal{I} = \{i, j\}|Y_0$  is available),  $\mathbb{M}^{F}_{parameters}$   
approximates it with X under a  
trace norm regularizer and hinge  
loss error measure (binary case):  
 $m_{in}^{(m)} ||X||_{1} \in \sum_{i \neq i} h_{i} h_{i} h_{i} X_{i}$   
 $m_{in}^{(m)} ||X||_{1} \in \sum_{i \neq i} h_{i} h_{i} h_{i} X_{i}$   
 $m_{in}^{(m)} ||X||_{1} \in \sum_{i \neq i} h_{i} h_{i} h_{i} X_{i}$   
 $m_{in}^{(m)} ||X||_{1} \in \sum_{i \neq i} h_{i} h_{i} h_{i} X_{i}$   
 $m_{in}^{(m)} ||X||_{1} \in \sum_{i \neq i} h_{i} h_{i} h_{i} X_{i} h_{i} X_{i} h_{i} h_$ 

- obse
- pred
- subs
- still with

To auto **IBP** prio unboun Minjie Xu, Jun Zhu and Bo Zhang

Department of Computer Science and Technology, LITS, TNList, Tsinghua University

Tsinghua University



\_ \_ \_ \_ \_ \_ \_ \_

 $+\infty$ 

# **Learning and Inference**

We perform variational inference under the truncated mean-field assumption (K the truncation level):

 $p(\nu, Z, V, \theta) = \prod_{k=1}^{K} p(\nu_k) p(Z) p(V) p(\theta), \ \nu_k \sim \text{Beta}(\gamma_{k1}, \gamma_{k2})$ 

- For iPM<sup>3</sup>F (iBPM<sup>3</sup>F likewise), we have (details in Appendix C & D)
- $p(V) = \prod_{j=1}^{M} \mathcal{N}(V_j | \Lambda_j, \sigma^2 I)$  and each  $\Lambda_j$  is solution to a linear SVM
- $\gamma_{k1}, \gamma_{k2}$  follows the same approximate update rule as developed by (Doshi-Velez, et al., AISTATS'09)
- $p(Z) = \prod_{i=1}^{N} \prod_{k=1}^{K} \text{Bernoulli}(Z_{ik}|\psi_{ik})$  and each  $\psi_i$  can be easily updated via coordinate descent
- $p(\theta) = \prod_{i=1}^{N} \prod_{r=1}^{L-1} \mathcal{N}(\theta_{ir} | \varrho_{ir}, \varsigma^2)$  and each  $\varrho_{ir}$  is solution to a linear SVM We iterate through each component until convergence.

Note that the number of *active* factors *a posteriori* follows

$$\mathbb{E}_{p}[K_{+}] = \sum_{k=1}^{K} \left( 1 - \prod_{i=1}^{N} (1 - \psi_{ik}) \right)$$

# **Experimental Results**

We test on MovieLens 1M and EachMovie data sets and compare with M<sup>3</sup>F, PMF and BPMF under NMAE error measure. Performance of iPM<sup>3</sup>F on MovieLens 1M

Table 1: NMAE performance of different models on MovieLens and EachMovie

	MovieLens				EachMovie			
Algorithm	weak		strong		weak		strong	
M <sup>3</sup> F [11]	.4156	$5 \pm .0037$ .4203 $\pm$		E.0138	$.4397 \pm .0006$		$.4341 \pm .0025$	
PMF [13]	.4332	$2 \pm .0033$	.4413 =	$\pm .0074$	$.4466 \pm .000$	.0016	$.4579 \pm .0$	016
BPMF [12]	$.4235 \pm .0023$		$.4450 \pm .0085$		$.4352 \pm .0014$		$.4445 \pm .0005$	
$M^3F^*$	.4176	$6 \pm .0016$	.4227 =	E.0072	$.4348 \pm$	.0023	$.4301 \pm .0$	034
iPM <sup>3</sup> F	$.4031 \pm .0030$		$.4135 \pm .0109$		$.4211 \pm .0019$		$.4224 \pm .0051$	
iBPM <sup>3</sup> F	$.4050 \pm .0029$		$.4089 \pm .0146$		$.4268 \pm .0029$		$.4403\pm.0040$	
Table 3: Performance of iPM <sup>3</sup> F with and without probabilistic treatment of $\theta$								
Algor	Algorithm		MovieLens		EachMovie		pEachMovie	
w/ pro	w/ prob.		:.0030	.4211	$\pm .0019$	.3954	$\pm .0026$	
w/o pi	w/o prob.		$.4056\pm.0043$		$.4256\pm.0011$		$.4026\pm.0023$	
margi	margin		.0013	.0045	$\pm .0016$	.0072	$\pm .0045$	

Expected number of features per user a posteriori (3 a priori)





Influence of the truncation level



## Conclusion

We endow M<sup>3</sup>F with a probabilistic background, which further enables the integration of nonparametric Bayesian techniques to *automatically* infer the dimensionality of the latent feature space and thus bypass explicit model complexity control. Our extensive experimental studies verify this benefit and demonstrate that our methods have competitive performance on real word collaborative filtering data sets as well.