

210430 Special Lectures on Database (RecSys)

Factorization Machines

Heeseung Yun

heeseung.yun@vision.snu.ac.kr

Contents

- Background
- Factorization Machines
- Comparison
- Conclusion

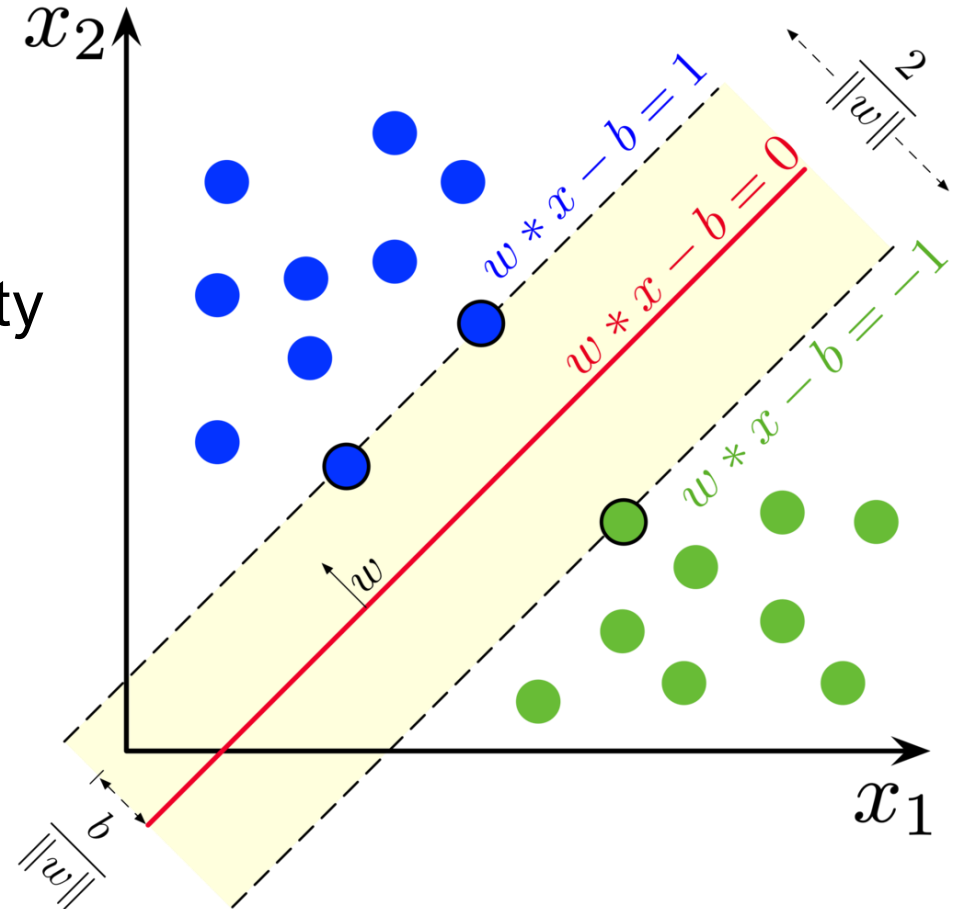
Background

Previous Recommender Systems

- Support Vector Machines
- Tensor Factorization

Previous Recommender Systems

- Support Vector Machines
 - (+) General predictor for any \mathbb{R}^n
 - (-) Susceptible to sparsity
 - (-) Susceptible to complex non-linearity

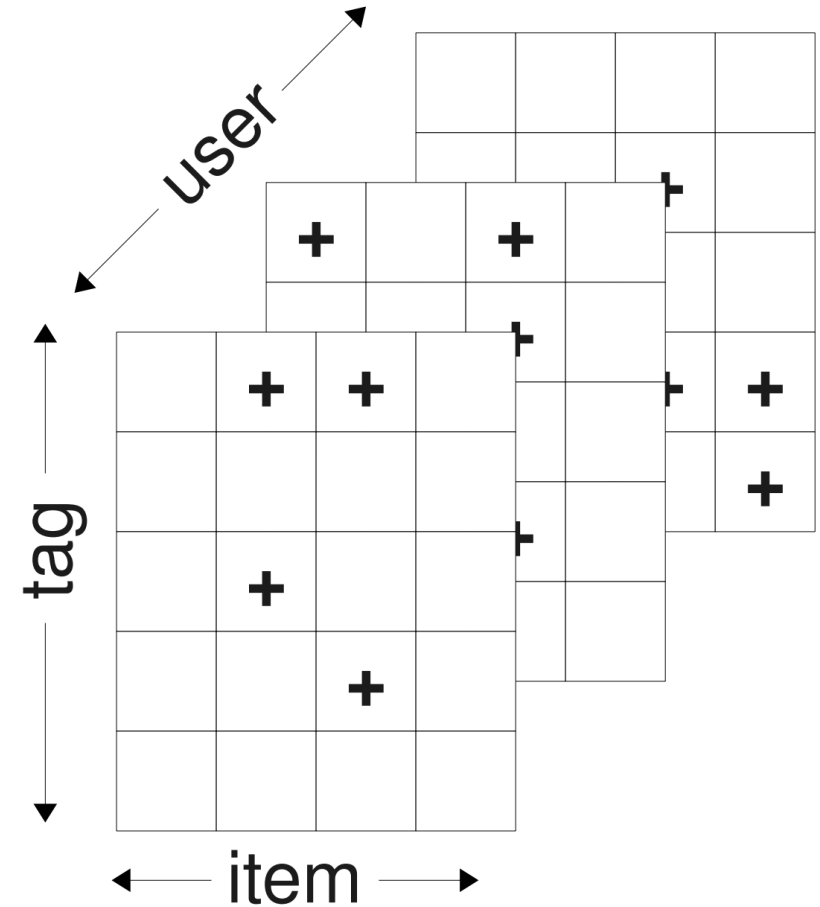


Previous Recommender Systems

- Tensor Factorization
 - (+) Sparsity-aware
 - (−) General features (\mathbb{R}^n) not applicable
 - Mostly categorical
 - (−) Task-specific modeling

Previous Recommender Systems

- Tensor Factorization
 - (ex1) SVD++ for movie rating prediction
 - Factorization with implicit feedback
 - (ex2) PITF for personalized tagging
 - Linearized Tucker decomposition



Koren. Factorization Meets the Neighborhood. In SIGKDD 2008.

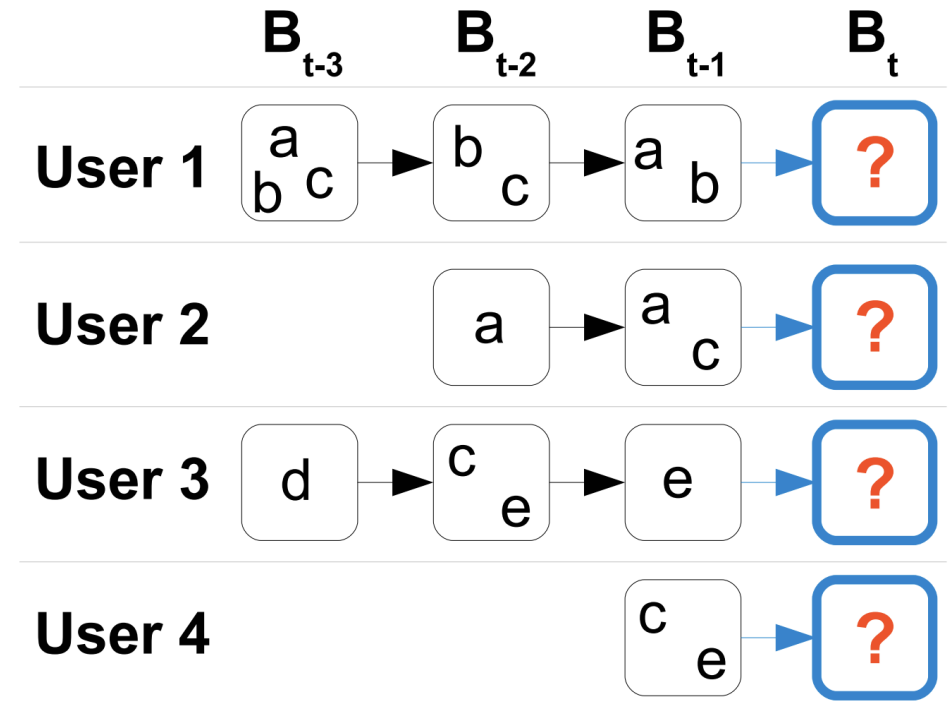
Rendle et al. Pairwise Interaction Tensor Factorization for Personalized Tag Recommendation. In WSDM 2010.

Previous Recommender Systems

- Tensor Factorization

- (ex3) FPMC for Next-Basket Recommendation

- Factorization + Markov Chain



- (ex4) PARAFAC

- Parallel Factor Analysis

Previous Recommender System

- Need for a unified framework
 - Take nested variable interaction into account
 - Efficient in terms of time & parameters
 - Work with various features and task-agnostic

Previous Recommender System

- Need for a unified framework
 - Take nested variable interaction into account → Sparsity-Robustness
 - Efficient in terms of time & parameters → Scalability
 - Work with various features and task-agnostic → Generalizability
- Solution: **Factorization Machines**

Factorization Machines

Definition

- Objective – predict $y : x \rightarrow T$ ($x \in \mathbb{R}^n$)
 - T depends on context
 - Real-valued, binary, ranking, etc.
 - x is generally categorical
 - BoW, transaction, etc.
 - Thereby highly sparse, i.e., $\bar{m}_D \ll n$
 - But not always categorical
 - Month, normalized indicator, \mathbb{R}^n

Definition

- Common tasks
 - Rating – regression with MSE
 - Binary – classification with hinge/logit loss (BCE)
 - Ranking – pairwise classification loss
 - Kendall's τ (Joachim, 2002)
 - Margin ranking loss $\max\left(0, f(x_{low}) + \alpha - f(x_{high})\right)$

Definition

- Objective – predict $y : x \rightarrow T$ ($x \in \mathbb{R}^n$)
 - T depends on context
 - Real-valued, binary, ranking, etc.
 - x is generally categorical
 - BoW, transaction, etc.
 - Thereby highly sparse, i.e., $\bar{m}_D \ll n$
↓ average # of nonzero entries
 - But not always categorical
 - Month, normalized indicator, \mathbb{R}^n

Definition

- Example – movie rating

| | Feature vector x | | | | | | | | | | | | | | | Target y | | | | | | |
|-----------|--------------------|---|---|-----|-------|----|----|----|-----|--------------------|-----|-----|-----|-----|------|------------------|----|----|----|-----|---|-----------|
| $x^{(1)}$ | 1 | 0 | 0 | ... | 1 | 0 | 0 | 0 | ... | 0.3 | 0.3 | 0.3 | 0 | ... | 13 | 0 | 0 | 0 | 0 | ... | 5 | $y^{(1)}$ |
| $x^{(2)}$ | 1 | 0 | 0 | ... | 0 | 1 | 0 | 0 | ... | 0.3 | 0.3 | 0.3 | 0 | ... | 14 | 1 | 0 | 0 | 0 | ... | 3 | $y^{(2)}$ |
| $x^{(3)}$ | 1 | 0 | 0 | ... | 0 | 0 | 1 | 0 | ... | 0.3 | 0.3 | 0.3 | 0 | ... | 16 | 0 | 1 | 0 | 0 | ... | 1 | $y^{(2)}$ |
| $x^{(4)}$ | 0 | 1 | 0 | ... | 0 | 0 | 1 | 0 | ... | 0 | 0 | 0.5 | 0.5 | ... | 5 | 0 | 0 | 0 | 0 | ... | 4 | $y^{(3)}$ |
| $x^{(5)}$ | 0 | 1 | 0 | ... | 0 | 0 | 0 | 1 | ... | 0 | 0 | 0.5 | 0.5 | ... | 8 | 0 | 0 | 1 | 0 | ... | 5 | $y^{(4)}$ |
| $x^{(6)}$ | 0 | 0 | 1 | ... | 1 | 0 | 0 | 0 | ... | 0.5 | 0 | 0.5 | 0 | ... | 9 | 0 | 0 | 0 | 0 | ... | 1 | $y^{(5)}$ |
| $x^{(7)}$ | 0 | 0 | 1 | ... | 0 | 0 | 1 | 0 | ... | 0.5 | 0 | 0.5 | 0 | ... | 12 | 1 | 0 | 0 | 0 | ... | 5 | $y^{(6)}$ |
| | A | B | C | ... | TI | NH | SW | ST | ... | TI | NH | SW | ST | ... | Time | TI | NH | SW | ST | ... | | |
| | User | | | | Movie | | | | | Other Movies rated | | | | | | Last Movie rated | | | | | | |

Definition

- Example – movie rating

| Feature vector x | | | | | | | | | | | | | | | Target y | | | | | | | |
|--------------------|-------------|---|---|-----|-------|----|----|----|-----|--------------------|-----|-----|-----|-----|------------|------------------|----|----|----|-----|---|-----------|
| $x^{(1)}$ | 1 | 0 | 0 | ... | 1 | 0 | 0 | 0 | ... | 0.3 | 0.3 | 0.3 | 0 | ... | 13 | 0 | 0 | 0 | 0 | ... | 5 | $y^{(1)}$ |
| $x^{(2)}$ | 1 | 0 | 0 | ... | 0 | 1 | 0 | 0 | ... | 0.3 | 0.3 | 0.3 | 0 | ... | 14 | 1 | 0 | 0 | 0 | ... | 3 | $y^{(2)}$ |
| $x^{(3)}$ | 1 | 0 | 0 | ... | 0 | 0 | 1 | 0 | ... | 0.3 | 0.3 | 0.3 | 0 | ... | 16 | 0 | 1 | 0 | 0 | ... | 1 | $y^{(2)}$ |
| $x^{(4)}$ | 0 | 1 | 0 | ... | 0 | 0 | 1 | 0 | ... | 0 | 0 | 0.5 | 0.5 | ... | 5 | 0 | 0 | 0 | 0 | ... | 4 | $y^{(3)}$ |
| $x^{(5)}$ | 0 | 1 | 0 | ... | 0 | 0 | 0 | 1 | ... | 0 | 0 | 0.5 | 0.5 | ... | 8 | 0 | 0 | 1 | 0 | ... | 5 | $y^{(4)}$ |
| $x^{(6)}$ | 0 | 0 | 1 | ... | 1 | 0 | 0 | 0 | ... | 0.5 | 0 | 0.5 | 0 | ... | 9 | 0 | 0 | 0 | 0 | ... | 1 | $y^{(5)}$ |
| $x^{(7)}$ | 0 | 0 | 1 | ... | 0 | 0 | 1 | 0 | ... | 0.5 | 0 | 0.5 | 0 | ... | 12 | 1 | 0 | 0 | 0 | ... | 5 | $y^{(6)}$ |
| | A | B | C | ... | TI | NH | SW | ST | ... | TI | NH | SW | ST | ... | Time | TI | NH | SW | ST | ... | | |
| | User | | | | Movie | | | | | Other Movies rated | | | | | | Last Movie rated | | | | | | |
| | Categorical | | | | | | | | | Real-valued | | | | | Int | | | | | | | |

Definition

- Prediction

$$\hat{y}(\mathbf{x}) := \overset{\text{Bias}}{w_0} + \sum_{i=1}^n \overset{\text{Signal Strength}}{w_i x_i} + \sum_{i=1}^n \sum_{j=i+1}^n \overset{\text{Interaction}}{\langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j}$$

- Parameters: $O(kn)$

$$w_0 \in \mathbb{R}, \quad \mathbf{w} \in \mathbb{R}^n, \quad \mathbf{V} \in \mathbb{R}^{n \times k}$$

- Hyperparameter: factor dimension $k \in \mathbb{N}_0^+$
 - Larger k theoretically reconstructs better
 - Smaller k leads to better generalization

Definition

- Factorization counts
 - Latent factors reflect context

| (Rating) | Titanics | Star Wars | Star Trek |
|----------|----------|-----------|-----------|
| Alice | 5 | 1 | ?? |
| Bob | | 4 | 5 |
| Charlie | 1 | | 5 |

→ Ratings are not independent

Complexity

- Definition

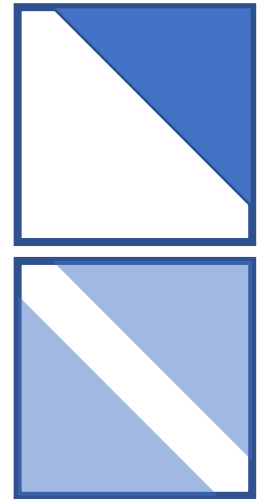
$$\hat{y}(\mathbf{x}) := w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j$$

- Naïve approach: $O(kn^2)$

Complexity

- Optimization

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j - \frac{1}{2} \sum_{i=1}^n \langle \mathbf{v}_i, \mathbf{v}_i \rangle x_i x_i \\ &= \frac{1}{2} \left(\sum_{i=1}^n \sum_{j=1}^n \sum_{f=1}^k v_{i,f} v_{j,f} x_i x_j - \sum_{i=1}^n \sum_{f=1}^k v_{i,f} v_{i,f} x_i x_i \right) \end{aligned}$$



Complexity

- Optimization

$$= \frac{1}{2} \sum_{f=1}^k \left(\left(\sum_{i=1}^n v_{i,f} x_i \right) \left(\sum_{j=1}^n v_{j,f} x_j \right) - \sum_{i=1}^n v_{i,f}^2 x_i^2 \right)$$

$$= \frac{1}{2} \sum_{f=1}^k \left(\left(\sum_{i=1}^n v_{i,f} x_i \right)^2 - \sum_{i=1}^n v_{i,f}^2 x_i^2 \right)$$

- Time complexity for inference: $O(kn)$
 - $O(k\bar{m}_D)$ for practical usage

Complexity

- Training
 - Stochastic Gradient Descent

$$\frac{\partial}{\partial \theta} \hat{y}(\mathbf{x}) = \begin{cases} 1, & \text{if } \theta \text{ is } w_0 \\ x_i, & \text{if } \theta \text{ is } w_i \\ x_i \underbrace{\sum_{j=1}^n v_{j,f} x_j}_{\text{Constant}} - v_{i,f} x_i^2, & \text{if } \theta \text{ is } v_{i,f} \end{cases}$$

- Time complexity for learning step: $O(kn)$, i.e., $O(1)$ per parameter

Params & inference time & training time are asymptotically linear!

Extension

- d-way Factorization Machine

$$\hat{y}(x) := \overset{\text{Bias}}{w_0} + \sum_{i=1}^n \overset{\text{Signal Strength}}{w_i x_i} + \sum_{l=2}^d \sum_{i_1=1}^n \cdots \sum_{i_l=i_{l-1}+1}^n \left(\prod_{j=1}^l x_{i_j} \right) \left(\sum_{f=1}^{k_l} \prod_{j=1}^l v_{i_j, f}^{(l)} \right)$$

(PARAFAC)
Interaction

- ex. Interaction of five SNS hashtags
- Can be optimized as in 2-way FM
 - i.e., linear complexity

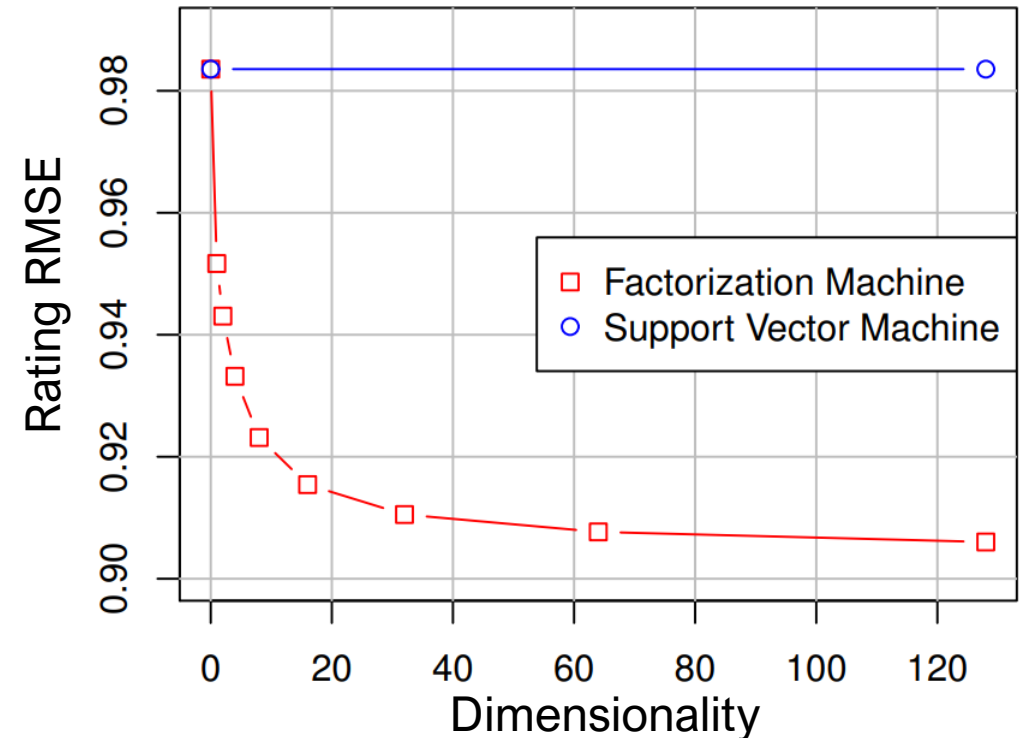
Comparison

SVM

- SVM with Linear Kernel $\phi(\mathbf{x}) := (1, x_1, \dots, x_n)$

$$\hat{y}(\mathbf{x}) = w_0 + \sum_{i=1}^n w_i x_i = w_0 + w_u + w_i$$

- Identical to degree=1 FM
- Netflix rating prediction
 - Interaction does count!



SVM

- SVM with Polynomial Kernel $(1, x_1, \dots, x_n)^d$

- For $d = 2$

$$\hat{y}(\mathbf{x}) = w_0 + \sqrt{2} \sum_{i=1}^n w_i x_i + \sum_{i=1}^n w_{i,i}^{(2)} x_i^2 + \sqrt{2} \sum_{i=1}^n \sum_{j=i+1}^n \overset{\text{Independence}}{w_{i,j}^{(2)}} x_i x_j$$

- FM without factorization

- cf. Original FM

$$\hat{y}(\mathbf{x}) := w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \overset{\text{Interaction}}{\langle \mathbf{v}_i, \mathbf{v}_j \rangle} x_i x_j$$

SVM

- SVM with Polynomial Kernel $(1, x_1, \dots, x_n)^d$

- For degree=2 FM with $x_u = x_i = 1, x_{else} = 0$

$$\begin{aligned}\hat{y}(\mathbf{x}) &= w_0 + \sqrt{2} \sum_{i=1}^n w_i x_i + \sum_{i=1}^n w_{i,i}^{(2)} x_i^2 + \sqrt{2} \sum_{i=1}^n \sum_{j=i+1}^n w_{i,j}^{(2)} x_i x_j \\ &= w_0 + \sqrt{2}(w_u + w_i) + w_{u,u}^{(2)} + w_{i,i}^{(2)} + \sqrt{2}w_{u,i}^{(2)}\end{aligned}$$

SVM

- SVM with Polynomial Kernel $(1, x_1, \dots, x_n)^d$

- For degree=2 FM with $x_u = x_i = 1, x_{else} = 0$

$$\hat{y}(\mathbf{x}) = w_0 + \sqrt{2} \sum_{i=1}^n w_i x_i + \sum_{i=1}^n w_{i,i}^{(2)} x_i^2 + \sqrt{2} \sum_{i=1}^n \sum_{j=i+1}^n w_{i,j}^{(2)} x_i x_j$$

$$= w_0 + \sqrt{2}(w_u + w_i) + w_{u,u}^{(2)} + w_{i,i}^{(2)} + \sqrt{2}w_{u,i}^{(2)}$$

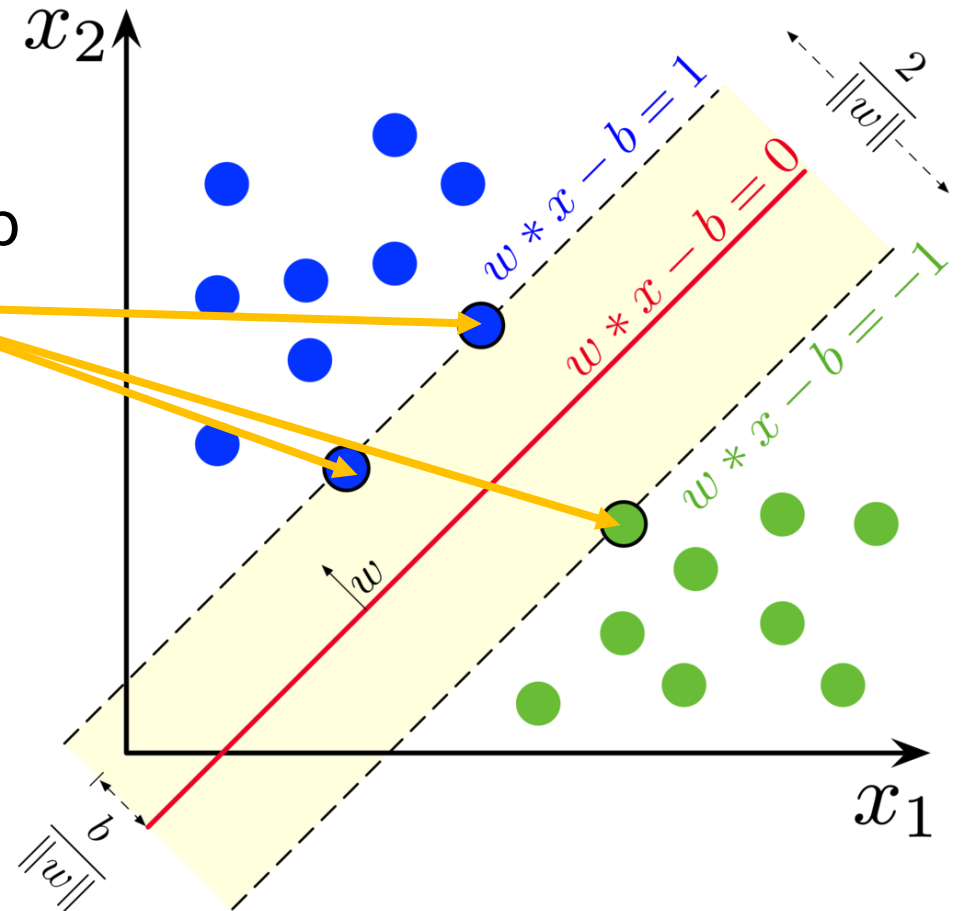
Identical to
 $w_u + w_i$

Zero for most
combination
 $(\bar{m}_D \ll n)$

- i.e., discards latent factors

SVM

- (Slide 5 revisited)
 - (–) Susceptible to sparsity
 - (–) Susceptible to complex relationship
 - (–) Dependent on support vectors
 - (·) Solved in dual, not primal



Tensor Factorization

- Matrix Factorization (MF, SVD)

$$\hat{y}(\mathbf{x}) = w_0 + w_u + w_i + \langle \mathbf{v}_u, \mathbf{v}_i \rangle$$

- Identical to FM with $x_u = x_i = 1, x_{else} = 0$

$$\hat{y}(\mathbf{x}) := w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j$$

- PARAFAC

- Features with multiple categorical variables, i.e., d-degree FM
- Both are limited to categorical variables

Tensor Factorization

- Movie rating (SVD++)

$$\hat{y}(\mathbf{x}) = w_0 + w_u + w_i + \langle \mathbf{v}_u, \mathbf{v}_i \rangle + \frac{1}{\sqrt{|N_u|}} \sum_{l \in N_u} \langle \mathbf{v}_i, \mathbf{v}_l \rangle$$

user-movie
+
movie-rated

- Partial model of FM with $x_u = x_i = 1, x_l = \frac{1}{\sqrt{|N_u|}}, x_{else} = 0$

$$\hat{y}(\mathbf{x}) := w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j$$

Tensor Factorization

- Movie rating (FM)

$$\hat{y}(\mathbf{x}) = w_0 + w_u + w_i + \langle \mathbf{v}_u, \mathbf{v}_i \rangle + \frac{1}{\sqrt{|N_u|}} \sum_{l \in N_u} \langle \mathbf{v}_i, \mathbf{v}_l \rangle$$

user-movie
+
movie-rated

$$+ \frac{1}{\sqrt{|N_u|}} \sum_{l \in N_u} \left(w_l + \langle \mathbf{v}_u, \mathbf{v}_l \rangle + \frac{1}{\sqrt{|N_u|}} \sum_{l' \in N_u, l' > l} \langle \mathbf{v}_l, \mathbf{v}'_{l'} \rangle \right)$$

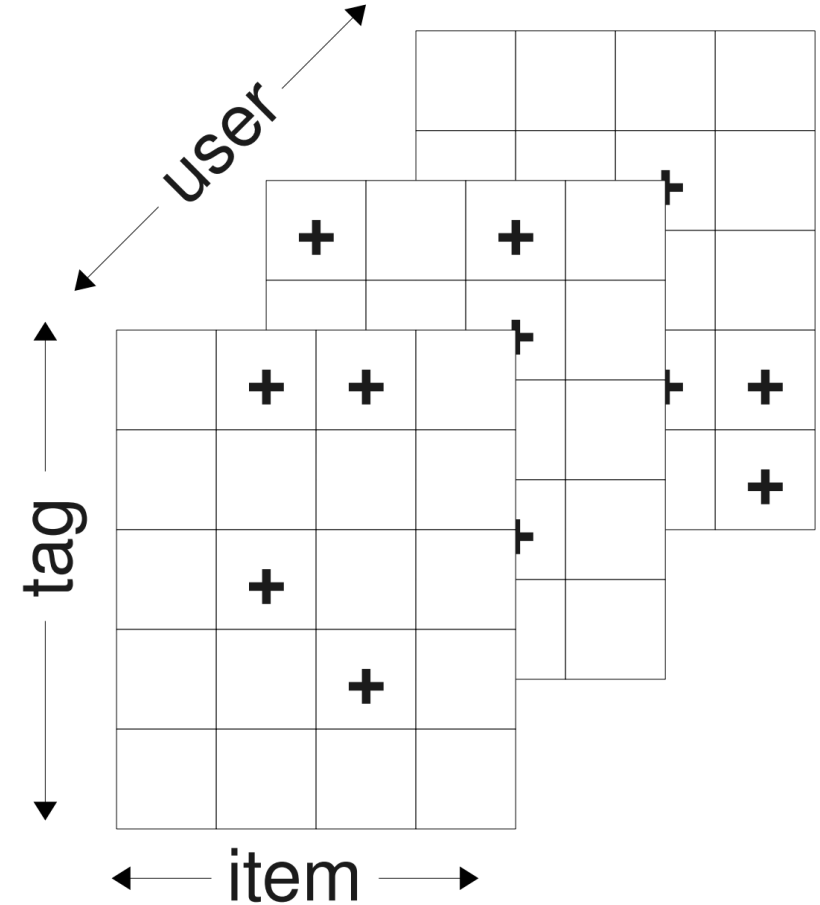
+
user-rated
+
rated-rated

Tensor Factorization

- Personal Tag Recommendation (PITF)

$$\hat{y}_{u,i,t} = \sum_f \hat{u}_{u,f} \cdot \hat{t}_{t,f}^U + \sum_f \hat{i}_{i,f} \cdot \hat{t}_{t,f}^I$$

user-tag MF item-tag MF



Tensor Factorization

- Personal Tag Recommendation (PITF)

$$\hat{y}_{u,i,t} = \sum_f \hat{u}_{u,f} \cdot \hat{t}_{t,f}^U + \sum_f \hat{i}_{i,f} \cdot \hat{t}_{t,f}^I$$

- Near-identical to FM optimized for ranking task with $x_u = x_i = x_t = 1$

$$\begin{aligned} \hat{y}(\mathbf{x}) &:= w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j \\ &= w_0 + w_u + w_i + w_t + \langle \mathbf{v}_u, \mathbf{v}_i \rangle + \langle \mathbf{v}_u, \mathbf{v}_t \rangle + \langle \mathbf{v}_i, \mathbf{v}_t \rangle \end{aligned}$$

Tensor Factorization

- Personal Tag Recommendation (PITF)

$$\hat{y}_{u,i,t} = \sum_f \hat{u}_{u,f} \cdot \hat{t}_{t,f}^U + \sum_f \hat{i}_{i,f} \cdot \hat{t}_{t,f}^I$$

- Near-identical to FM optimized for ranking task with $x_u = x_i = x_t = 1$

$$\begin{aligned} \hat{y}(\mathbf{x}) &:= w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j \\ &= w_0 + w_u + w_i + w_t + \langle \mathbf{v}_u, \mathbf{v}_i \rangle + \langle \mathbf{v}_u, \mathbf{v}_t \rangle + \langle \mathbf{v}_i, \mathbf{v}_t \rangle \end{aligned}$$

Independent of tag t

Tensor Factorization

- Personal Tag Recommendation (PITF)

$$\hat{y}_{u,i,t} = \sum_f \hat{u}_{u,f} \cdot \hat{t}_{t,f}^U + \sum_f \hat{i}_{i,f} \cdot \hat{t}_{t,f}^I \quad \text{Distinct tag factorization}$$

- Near-identical to FM optimized for ranking task with $x_u = x_i = x_t = 1$

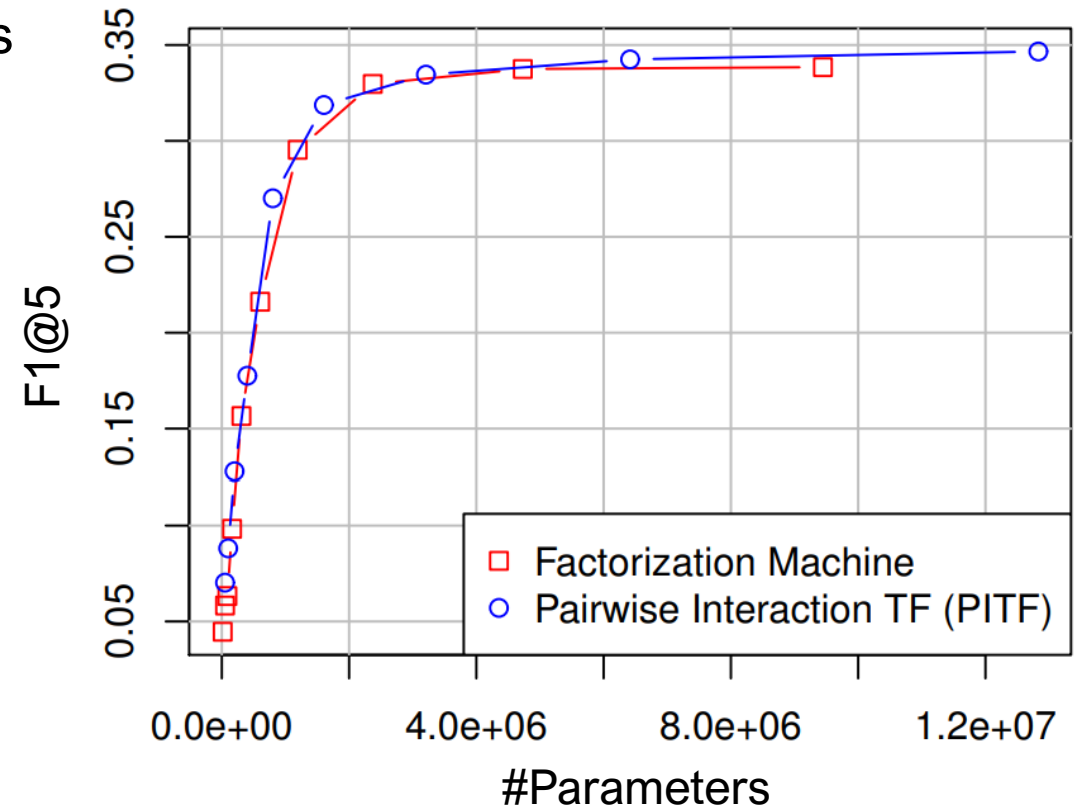
$$\begin{aligned} \hat{y}(\mathbf{x}) &:= w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j \\ &= w_0 + w_u + w_i + w_t + \langle \mathbf{v}_u, \mathbf{v}_i \rangle + \langle \mathbf{v}_u, \mathbf{v}_t \rangle + \langle \mathbf{v}_i, \mathbf{v}_t \rangle \end{aligned}$$

$$\hat{y}(\mathbf{x}) := w_t + \langle \mathbf{v}_u, \mathbf{v}_t \rangle + \langle \mathbf{v}_i, \mathbf{v}_t \rangle \quad \text{Shared tag factorization}$$

Bias

Tensor Factorization

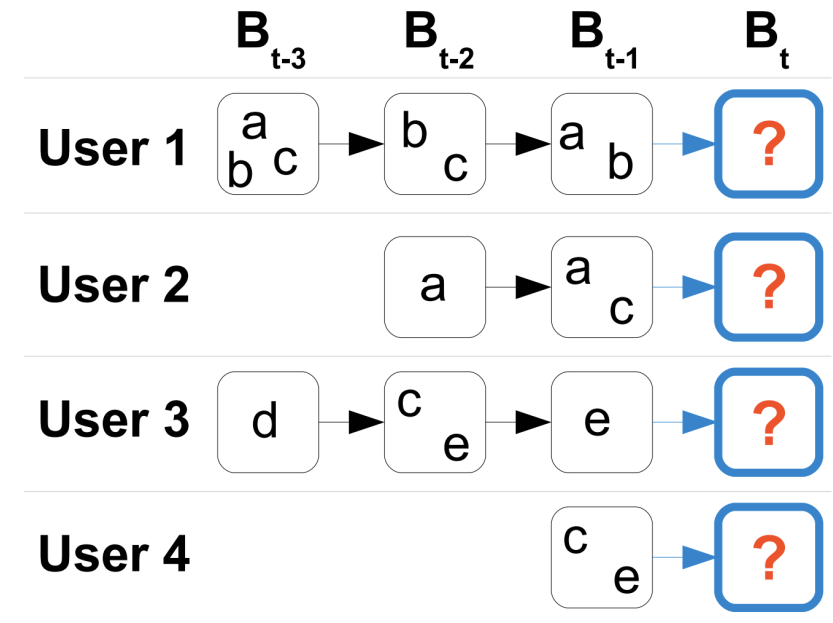
- Personal Tag Recommendation (PITF)
 - ECML/PKDD Discovery Challenge
 - Nearly identical behavior w.r.t params



Tensor Factorization

- Next-Basket Recommendation (FPMC)

$$\hat{x}_{u,t,i} := \underbrace{\langle v_u^{U,I}, v_i^{I,U} \rangle}_{\text{user-item MF}} + \underbrace{\frac{1}{|B_{t-1}^u|}}_{\text{Markov}} \sum_{l \in B_{t-1}^u} \underbrace{\langle v_i^{I,L}, v_l^{L,I} \rangle}_{\text{Item-bag MF}}$$



Tensor Factorization

- Next-Basket Recommendation (FPMC)

$$\hat{x}_{u,t,i} := \langle v_u^{U,I}, v_i^{I,U} \rangle + \frac{1}{|B_{t-1}^u|} \sum_{l \in B_{t-1}^u} \langle v_i^{I,L}, v_l^{L,I} \rangle$$

- Near-identical to FM optimized for ranking ($x_u = x_i = 1, x_l = \frac{1}{\sqrt{|B_{t-1}^u|}}$)

$$\hat{y}(\mathbf{x}) = w_0 + w_u + w_i + \langle \mathbf{v}_u, \mathbf{v}_i \rangle + \frac{1}{|B_{t-1}^u|} \sum_{l \in B_{t-1}^u} \langle \mathbf{v}_i, \mathbf{v}_l \rangle + \frac{1}{|B_{t-1}^u|} \sum_{l \in B_{t-1}^u} \left(w_l + \langle \mathbf{v}_u, \mathbf{v}_l \rangle + \frac{1}{|B_{t-1}^u|} \sum_{l' \in B_{t-1}^u, l' > l} \langle \mathbf{v}_l, \mathbf{v}_{l'} \rangle \right)$$

Tensor Factorization

- Next-Basket Recommendation (FPMC)

$$\hat{x}_{u,t,i} := \langle v_u^{U,I}, v_i^{I,U} \rangle + \frac{1}{|B_{t-1}^u|} \sum_{l \in B_{t-1}^u} \langle v_i^{I,L}, v_l^{L,I} \rangle$$

- Near-identical to FM optimized for ranking ($x_u = x_i = 1, x_l = \frac{1}{\sqrt{|B_{t-1}^u|}}$)

$$\hat{y}(\mathbf{x}) = w_0 + w_u + w_i + \langle \mathbf{v}_u, \mathbf{v}_i \rangle + \frac{1}{|B_{t-1}^u|} \sum_{l \in B_{t-1}^u} \langle \mathbf{v}_i, \mathbf{v}_l \rangle$$

$$+ \frac{1}{|B_{t-1}^u|} \sum_{l \in B_{t-1}^u} \left(w_l + \langle \mathbf{v}_u, \mathbf{v}_l \rangle + \frac{1}{|B_{t-1}^u|} \sum_{l' \in B_{t-1}^u, l' > l} \langle \mathbf{v}_l, \mathbf{v}_{l'} \rangle \right)$$

Independent of item i

Tensor Factorization

- Next-Basket Recommendation (FPMC)

$$\hat{x}_{u,t,i} := \langle v_u^{U,I}, v_i^{I,U} \rangle + \frac{1}{|B_{t-1}^u|} \sum_{l \in B_{t-1}^u} \langle v_i^{I,L}, v_l^{L,I} \rangle \quad \text{Distinct item factorization}$$

- Near-identical to FM optimized for ranking ($x_u = x_i = 1, x_l = \frac{1}{\sqrt{|B_{t-1}^u|}}$)

$$\begin{aligned} \hat{y}(\mathbf{x}) &= w_0 + w_u + w_i + \langle \mathbf{v}_u, \mathbf{v}_i \rangle + \frac{1}{|B_{t-1}^u|} \sum_{l \in B_{t-1}^u} \langle \mathbf{v}_i, \mathbf{v}_l \rangle \\ &\quad + \frac{1}{|B_{t-1}^u|} \sum_{l \in B_{t-1}^u} \left(w_l + \langle \mathbf{v}_u, \mathbf{v}_l \rangle + \frac{1}{|B_{t-1}^u|} \sum_{l' \in B_{t-1}^u, l' > l} \langle \mathbf{v}_l, \mathbf{v}_{l'} \rangle \right) \\ &= w_i + \langle \mathbf{v}_u, \mathbf{v}_i \rangle + \frac{1}{|B_{t-1}^u|} \sum_{l \in B_{t-1}^u} \langle \mathbf{v}_i, \mathbf{v}_l \rangle \quad \text{Shared item factorization} \\ &\quad \text{Bias} \end{aligned}$$

Conclusion

Conclusion

- Factorization Machines
 - Factorization-oriented
 - Linear parameters & prediction time & learning time
 - Subsumes existing (task-specific) State-of-the-Art models
- Additional remarks
 - (+) Unified RS framework with clear theoretical explanation
 - (−) Lacks experimental evidence (SVD++, FPMC, time complexity)

Thank You

Steffen Rendle. Factorization Machines. In ICDM 2010.

Appendix

Lagrangian Duality

- Primal SVM

$$\zeta_i = \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i - b))$$

$$\text{minimize } \frac{1}{n} \sum_{i=1}^n \zeta_i + \lambda \|\mathbf{w}\|^2$$

subject to $y_i(\mathbf{w}^T \mathbf{x}_i - b) \geq 1 - \zeta_i$ and $\zeta_i \geq 0$, for all i

- Dual SVM

$$\text{maximize } f(c_1 \dots c_n) = \sum_{i=1}^n c_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i c_i (\mathbf{x}_i^T \mathbf{x}_j) y_j c_j$$

subject to $\sum_{i=1}^n c_i y_i = 0$, and $0 \leq c_i \leq \frac{1}{2n\lambda}$ for all i