

Two-level monotonic multistage recommender systems

(Joint work with Xiaotong Shen and Wei Pan)

by **Ben Dai** (The Chinese University of Hong Kong)
on **EcoSta2022**

» Single-stage (binary) Recommender Systems

	Item 1	Item 2	...	Item m
User 1	✓	?	...	✗
User 2	?	✓	...	?
...				
User n	?	✗	...	✓

A single-stage (binary) Recommender System.

Input a user-item pair (i, j)

Goal predict its corresponding response $Y_{ij} = \pm 1$

Example predict if a user will **book** a item

» Multi-stage Recommender Systems

Deskdrop dataset

- * **Deskdrop** is an internal communications platform, which allows companies employees to share relevant articles with their peers, and collaborate around them

Size It contains about **73k logged users interactions** on more than **3k public articles** shared in the platform

User info [Id, agent, country]

Article info [Id, country, type, url]

Feedback **follow** > **like** > **view**

» Multi-stage Recommender Systems

Key characteristics to **Deskdrop** dataset

RS A typical RS dataset

Side info more features (continuous or categorical) for users/items

Feedback Monotonic property: **follow** > **like** > **view**

$$Y^t = -1 \text{ if } Y^{t-1} = -1, \text{ for } t = 1, \dots, T$$

* The feedback at stage- t' may already be observed.

Goal predict its responses Y_{ij}^t based on $Y_{ij}^{t'}$ for $1 \leq t' < t \leq T$.

- * obs a user **viewed** the item → if **like**?
- * obs a user **viewed** the item → if **follow**?
- * obs a user **liked** the item → if **follow**?

» Multi-stage Recommender Systems

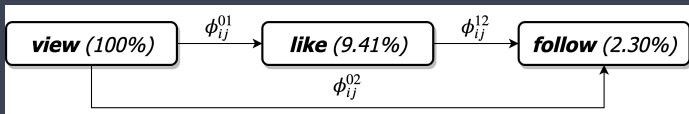
Input side info for a pair (i, j) + obs feedback $Y_{ij}^{t'}$ at stage- t'

Goal predict its response Y_{ij}^t at stage- t with $t > t'$

* A **pairwise** decision function based on the (t', t) -pair

$$\phi = (\phi^{t't}(\mathbf{x}, \mathbf{y}^{t'}))_{0 \leq t' < t \leq T}$$

Example In **Deskdrop** dataset



» Multi-stage RS: Evaluation

- * **Pairwise** mis-classification error:

$$e(\phi) = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \sum_{t' < t} w_{t't} \mathbf{E}(I(Y_{ij}^t \phi^{t't}(\mathbf{X}_{ij}, Y_{ij}^{t'}) \leq 0)) \quad (1)$$

- * $w_{t't}$ is a stagewise weight to quantify the **importance** of prediction for the t -th stage based on the t' -th stage
 - * **Next-stage-prediction**: $w_{t't} = I(t - t' = 1)$
 - * **Last-stage-prediction**: $w_{t't} = I(t = T)$
- * In general, $w_{t't} \geq 0$ can be a custom weight to indicate the relative importance among different stage-pairs

» Multi-stage RS: Bayes decision rule

Lemma 1 indicates two facts for the **optimal stagewise classifier** $\bar{\phi}$ with respect to the loss (1).

A. If $y_{ij}^{t'} = -1$ then the best prediction is -1.

$$\bar{\phi}^{t't}(\mathbf{x}_{ij}, y_{ij}^{t'}) = \begin{cases} -1, & \text{if } y_{ij}^{t'} = -1, \\ \text{sgn}(\bar{f}^{t't}(\mathbf{x}_{ij})), & \text{if } y_{ij}^{t'} = 1, \end{cases}$$

» Multi-stage RS: Bayes decision rule

Lemma 1 indicates two facts for the **optimal stagewise classifier** $\bar{\phi}$ with respect to the loss (1).

A. If $y_{ij}^t = -1$ then the best prediction is -1.

$$\bar{\phi}^{t,t}(\mathbf{x}_{ij}, y_{ij}^t) = \begin{cases} -1, & \text{if } y_{ij}^t = -1, \\ \text{sgn}(\bar{f}^{t,t}(\mathbf{x}_{ij})), & \text{if } y_{ij}^t = 1, \end{cases}$$

B. The **optimal conditional classifier** $\bar{f}^{t,t}(\mathbf{x}_{ij})$ satisfies

$$\text{sgn}(\bar{f}_{ij}^{t,t}(\mathbf{x}_{ij})) = \text{sgn}(\mathbf{P}_{ij}(Y_{ij}^t = 1 | Y_{ij}^t = 1, \mathbf{X} = \mathbf{x}_{ij}) - 1/2),$$

iff. **two-level monotonic property**

$$\text{Forward: } \text{sgn}(\bar{f}_{ij}^{t,t+1}(\mathbf{x}_{ij})) = -1, \quad \text{if } \text{sgn}(\bar{f}_{ij}^{t,t}(\mathbf{x}_{ij})) = -1,$$

$$\text{Backward: } \text{sgn}(\bar{f}_{ij}^{t,t}(\mathbf{x}_{ij})) = -1, \quad \text{if } \text{sgn}(\bar{f}_{ij}^{t+1,t}(\mathbf{x}_{ij})) = -1. \quad (2)$$

» Multi-stage RS: Bayes decision rule

- C. **Additive** form. There exists $\bar{h}_{ij}^r(\mathbf{x}_{ij}) \geq 0$, (it suffices to model h_{ij}^r)

$$\bar{f}_{ij}^{t't}(\mathbf{x}_{ij}) = \bar{h}_{ij}^0(\mathbf{x}_{ij}) - \sum_{r=t'+1}^t \bar{h}_{ij}^r(\mathbf{x}_{ij}), \text{ with } \bar{h}_{ij}^r \geq 0; 0 \leq r \leq T. \quad (3)$$

This yields when $t' \rightarrow t$, $\bar{f}_{ij}^{t't}(\mathbf{x}_{ij})$ increases.

» Classifier: formulation

- * FunkSVD got **third place** in Netflix's competition (first and second places are resemble models)
- * Formulate the **interaction** btw the i -th user and j -th item.

$$h_{ij}^r = \mathbf{a}_i^\top \mathbf{b}_j$$

- * How to incorporate the **side information**? additional continuous and discrete features for users and items

$$\begin{aligned}
 * \underbrace{\text{sideInfo}}_{\mathbf{x}_{ij}}: & \left(\underbrace{\text{user_cont}}_{\mathbf{u}_i \in [0,1]^{p_1}} + \underbrace{\text{user_cate}}_{\mathbf{v}_j \in [n_1] \times \dots \times [n_{d_1}]} \right) \\
 & + \left(\underbrace{\text{item_cont}}_{\mathbf{s}_i \in [0,1]^{p_2}} + \underbrace{\text{item_cate}}_{\mathbf{o}_j \in [m_1] \times \dots \times [m_{d_2}]} \right)
 \end{aligned}$$

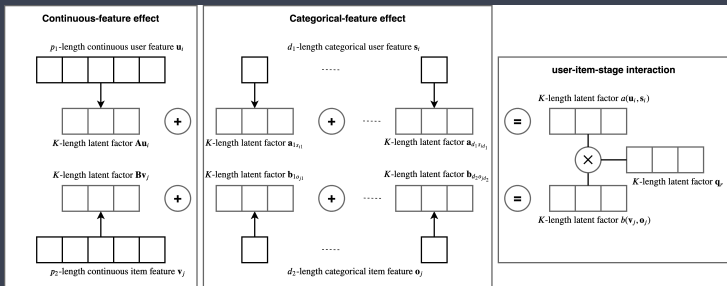
- * **User info:** [Id, agent, country]
- * **Article info:** [Id, country, type, url]

» Classifier: formulation

* Our formulation:

$$h_{ij}^r(\mathbf{x}_{ij}) = (a(\mathbf{u}_i, \mathbf{s}_i) \circ b(\mathbf{v}_j, \mathbf{o}_j))^T \mathbf{q}_r,$$

$$a(\mathbf{u}_i, \mathbf{s}_i) = \mathbf{A}\mathbf{u}_i + \sum_{l=1}^{d_1} \mathbf{a}_{l s_{il}}, \quad b(\mathbf{v}_j, \mathbf{o}_j) = \mathbf{B}\mathbf{v}_j + \sum_{l=1}^{d_2} \mathbf{b}_{l o_{jl}}, \quad (4)$$



* Additional constraints: $h_{ij}^r(\mathbf{x}_{ij}) \geq 0$

» Optimization: empirical loss

- * Given observations $(\mathbf{x}_{ij}, y_{ij}^t)_{(i,j) \in \Omega_0, 1 \leq t \leq T}$ and positive index set $\Omega_{t'} = \{(i, j) : y_{ij}^{t'} = 1\}$ at stage t' , the empirical loss is given as

$$\min_{\mathbf{A}, \mathbf{B}, \mathbf{a}, \mathbf{b}, \mathbf{q}} \sum_{t' < t} \sum_{(i,j) \in \Omega_{t'}} w_{t't} V \left(y_{ij}^t (h_{ij}^0(\mathbf{x}_{ij}) - \sum_{r=t'+1}^t h_{ij}^r(\mathbf{x}_{ij})) \right) + \lambda \text{Reg},$$

$$\mathbf{A} \geq \mathbf{0}, \mathbf{B} \geq \mathbf{0},$$

$$\mathbf{a}_{lh} \geq \mathbf{0}_K, h = 1, \dots, n_l; l = 1, \dots, d_1,$$

$$\mathbf{b}_{lh} \geq \mathbf{0}_K, h = 1, \dots, m_l; l = 1, \dots, d_2, \quad (5)$$

» Optimization: BCD

- * Jointly non-convex problem
- * Blockwise convex with respect to \mathbf{A} , \mathbf{B} , \mathbf{a}_l , \mathbf{b}_l , \mathbf{q}_r
- * Motivate to **blockwise coordinate descent** (BCD)
 - * Fixing the others, update one parameter
 - * Each subproblem is a **non-negative + drifted + weighted SVM**.
 - * solved by our Python library `varsvm`

» Optimization: BCD

User-effect block \mathbf{A} . This convex optimization solves for \mathbf{A} by fixing $(\mathbf{B}, \mathbf{a}, \mathbf{b}, \mathbf{q})$:

$$\min_{\text{vec}(\mathbf{A}) \geq 0} \sum_{t' < t} \sum_{(i,j) \in \Omega_{t'}} \underbrace{w_{t't}}_{\text{weights}} V \left(\underbrace{y_{ij}^t}_{\text{label}} \left(\underbrace{\text{vec}(\mathbf{A})^\top}_{\text{coeff.}} \underbrace{(\mathbf{u}_i \otimes (\mathbf{b}(\mathbf{v}_j, \mathbf{o}_j) \circ \mathbf{q}_{t't}))}_{\text{input feats.}} \right) \right. \\ \left. + \underbrace{\left(\sum_{l=1}^{d_1} \mathbf{a}_{lS_{il}} \right)^\top (\mathbf{b}(\mathbf{v}_j, \mathbf{o}_j) \circ \mathbf{q}_{t't}) \right)}_{\text{fixed intercept (drift)}} \right) + \lambda_1 \sum_{h=1}^{n_l} \|\text{vec}(\mathbf{A})\|_2^2,$$

where $\mathbf{q}_{t't} = \mathbf{1}_K - \sum_{r=t'+1}^t \mathbf{q}_r$, \otimes is the Kronecker product and $\text{vec}(\mathbf{A})$ is the column vectorization of matrix \mathbf{A} .

» Python library varsvm

- * Python *scikit-learn* estimators module
- * weighted SVMs, drifted SVMs, non-negative SVMs
- * Much **faster** than `sklearn.svm.SVC` in Python for weighted SVMs

Doc <https://variant-svm.readthedocs.io>

» Experiments: real application

- * **Deskdrop** article sharing
- * Features:
 - * User **cate feats**: [**user_Id**, **user_agent**, **user_region**]
 - * Item **cont feats**: Doc2Vec embedding of “plain text” of an item
 - * Item **cate feats**: [**article_Id**, **author_Id**, **language**]
- * Evaluation:
 - * Pairwise Mis-classification Error
 - * Equal stage weight: $w_{t'} = 1$
- * **Comeptitors**: DeepNN, GradBoost, SVM, OSVM, SVD++

» Experiments: real application

(observe t' , predict t)	Proposed	SVD ⁺⁺	GradBoost	SVM	DeepNN	OSVM
(Stage 0, Stage 1)	0.136(.001)	0.499(.003)	0.166(.000)	0.139(.000)	0.154(.001)	0.140(.000)
(Stage 0, Stage 2)	0.191(.001)	0.499(.007)	0.218(.001)	0.218(.001)	0.223(.006)	0.216(.001)
(Stage 1, Stage 2)	0.043(.000)	0.497(.007)	0.049(.000)	0.045(.000)	0.139(.011)	0.045(.000)
Overall	0.123(.001)	0.498(.006)	0.144(.000)	0.134(.000)	0.172(.005)	0.134(.000)
%Inconsistent	0.00%(.000)	7.79%(.002)	1.40%(.000)	1.90%(.000)	13.7%(.005)	0.98%(.002)

Class-balanced zero-one losses of six competitors on the *Deskdrop* dataset, in addition to %Inconsistent denoting a proportion of inconsistent instances violating the monotonicity (2).

- * Overall performance with the amount of improvement over 8.20%
- * With respect to inconsistency prediction, all other methods have inconsistent cases ranging from 0.98% to 13.7%.

» Summary

- * The proposed method can produce a recommendation for any subsequent stage given observations at any present stage, which is highly demanded in real applications. Yet, most recommender systems focus on a fixed present stage.
- * The **two-level monotonic property** is fully accounted for by our nonnegative additive latent factor model. As a result, it ensures prediction **consistency** across different stages.
- * Asymptotics and more numerical examples are included in our paper published in EJS.

Background
○○○

Framework
○○○○

Formulation
○○

Optimization
○○○○

Experiments
○○

Summary
○○●

Thank you!