

Lecture 2 – Part 2
Research on the Netflix Problem
(Mostly)

(January 16, 2015)

Mu Zhu
University of Waterloo

Content-boosted Matrix Factorization

- Zhu M (2014), “Making personalized recommendations in e-commerce”, in *Statistics in Action: A Canadian Outlook*, J. F. Lawless, Ed., Chapman & Hall, pp. 259–268.
- Nguyen J, Zhu M (2013), “Content-boosted matrix factorization techniques for recommender systems”, *Statistical Analysis and Data Mining* **6**, pp. 286–301.
- Forbes P, Zhu M (2011), “Content-boosted matrix factorization for recommender systems: Experiments with recipe recommendation”, in *Proceedings of the 5th ACM Conference on Recommender Systems*, pp. 261–264.

Toy Example

	<i>The Interpretation</i>		<i>The Origin</i>	
	<u><i>Moby-Dick</i></u> (Melville)	<u><i>of Dreams</i></u> (Freud)	<u><i>of Species</i></u> (Darwin)	<u><i>Relativity</i></u> (Einstein)
Alice	90	70 ?	30	10
Bob	90	70	30 ?	10
Cathy	10	30 ?	70	90
David	10	30	70 ?	90

User-item ratings (matrix \mathbf{R}). A question mark (?) indicates that the corresponding rating is treated **as if** it were missing/unobserved and to be predicted.

Toy Example

	<i>Moby-Dick</i>	<i>Dreams</i>	<i>Species</i>	<i>Relativity</i>
Themes of conflict	✓	✓	✓	×
Elem. of moral phil.	✓	✓	×	×
Darkness of human nature	✓	✓	×	×
Other empirical evidence	×	×	✓	✓
Logical rigor	×	×	✓	✓
A grand new theory	×	✓	✓	✓

Content attributes (matrix \mathbf{A}). Sometimes, we may have content information for each item, such as these. Here, each $\mathbf{a}_{ij} \in \{0, 1\}$ is a binary indicator of whether item i has attribute j .

Boosting by Content

- intuition: content features may explain why some users prefer certain items to others
- toy example: any user's ratings of any two items are always $20 \times (4 - Z)$ points apart if the two items share Z content features in common ($Z = 0, 1, 2, \text{ or } 3$)
- reality: content information will never completely determine user preferences, but still more likely than not to be at least partially informative
- our work: **two** different ways to incorporate such content information into the **matrix factorization** approach

Two Ways to Boost by Content

- add extra penalty

$$\min_{\mathbf{p}_u, \mathbf{q}_i} \quad (\text{original objective}) - \lambda \left[\sum_{i=1}^m \sum_{i' \in \mathcal{S}_c(i)} \frac{\mathbf{q}_i^T \mathbf{q}_{i'}}{|\mathcal{S}_c(i)|} \right]$$

where

$$\mathcal{S}_c(i) \equiv \{i' < i : \mathbf{a}_i^T \mathbf{a}_{i'} \geq c\}$$

- enforce explicit constraint

$$\mathbf{q}_i = \mathbf{B} \mathbf{a}_i \quad \text{or} \quad \mathbf{q}_i(k) = \sum_j \mathbf{B}(k, j) \mathbf{a}_i(j)$$

(I) baseline MF; (II) extra penalty; (III) explicit constraint

		<i>Moby-Dick</i>	<i>Dreams</i>	<i>Species</i>	<i>Relativity</i>
(I)	Alice	89	48	30	10
	Bob	90	70	51	11
	Cathy	11	52	70	90
	David	10	30	49	89
(II)	Alice	89	63	30	10
	Bob	90	70	37	11
	Cathy	11	37	70	90
	David	10	30	63	89
(III)	Alice	90	75	30	10
	Bob	90	70	25	10
	Cathy	10	25	70	90
	David	10	30	75	90

using $K = 2$; $\lambda = 1$; $c = 3$

Expanded Alternating Optimization

Murdoch WJ, Zhu M (2014), “Expanded alternating optimization of nonconvex functions with applications to matrix factorization and penalized regression”, [arXiv:1412.4128](https://arxiv.org/abs/1412.4128).

Alternating Optimization (AO)

- coordinate descent \in alternating optimization

$$\min_{z \in \mathbb{R}^d} f(z)$$

by partitioning z into blocks z_1, z_2, \dots and solving

$$\min_{z_b} f(z_1, z_2, \dots)$$

successively over $b = 1, 2, \dots, 1, 2, \dots$

- attractive if each individual problem (over z_b) is “easy”

Toy Example

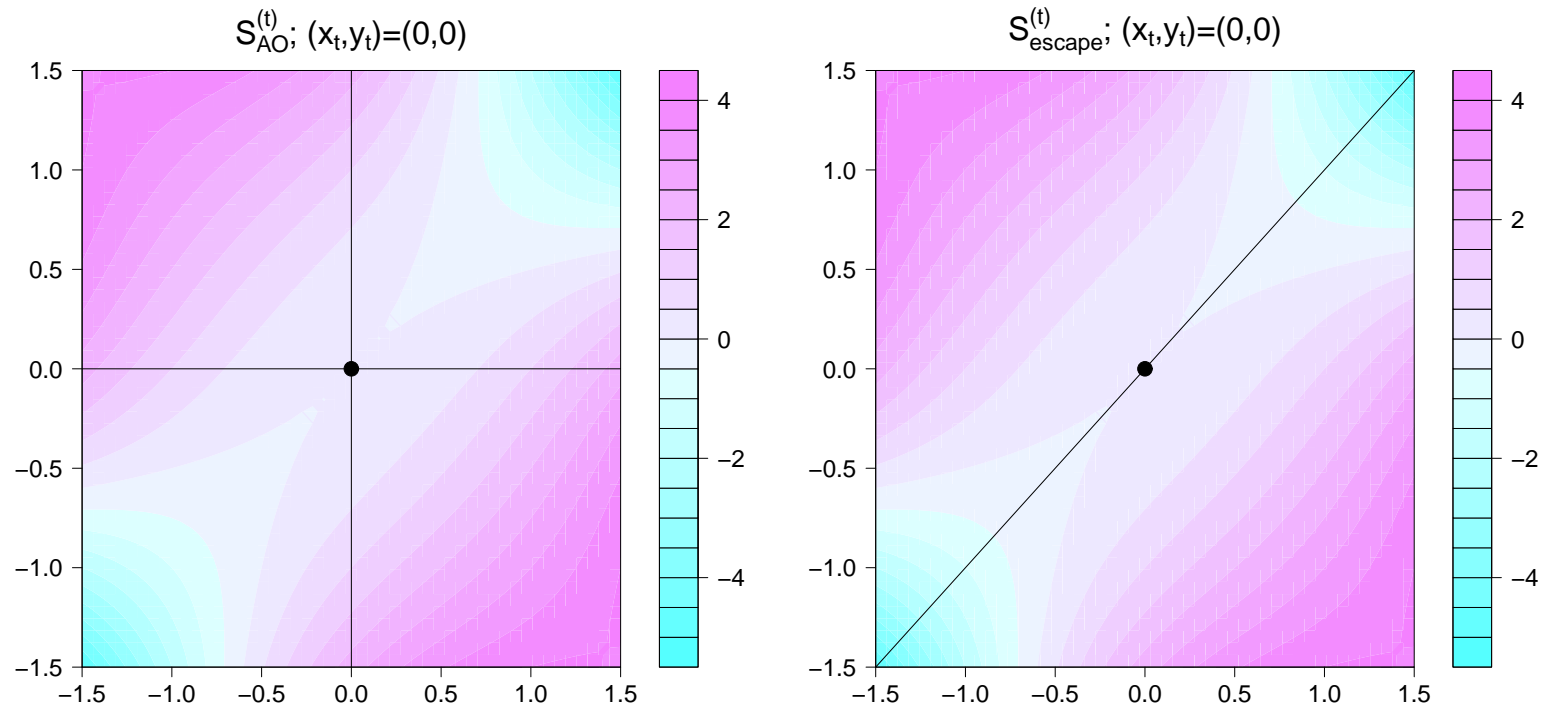
- minimize

$$f(x, y) = (x - y)^2 - x^2y^2$$

with coordinate descent

- suppose we reached $(x_t, y_t) = (0, 0)$ at iteration t
- fixing $x_t = 0$, $f(0, y) = y^2$ minimized at $y = 0$
- fixing $y_t = 0$, $f(x, 0) = x^2$ minimized at $x = 0$
- so, stuck at **saddle point** $(0, 0)$

Toy Example



Contours of $f(x, y) = (x - y)^2 - x^2 y^2$

Toy Example

$$\mathcal{S}_{AO}^{(t)} = \{(x, y) : x = 0\} \cup \{(x, y) : y = 0\}$$

$$\mathcal{S}_{escape}^{(t)} = \{(x, y) : x = y\}$$

Scaling

A. Tayal, T. F. Coleman & Y. Li (2014), “Primal explicit max margin feature selection for nonlinear support vector machines”, *Pattern Recognition* **47**, pp. 2153–2164.

- at each step,

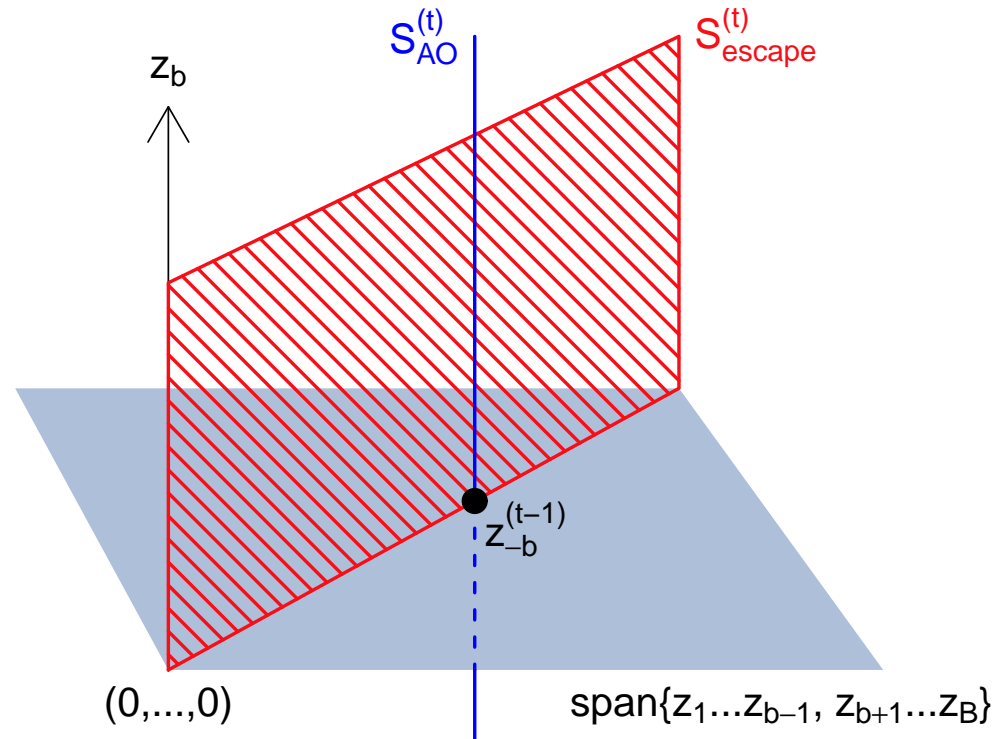
$$\min_{z_b, v_b} f(v_b z_1, \dots, v_b z_{b-1}, z_b, v_b z_{b+1}, \dots, v_b z_B)$$

- let v_b^* denote the solution of v_b ; all other components z_{-b} adjusted accordingly, i.e.,

$$z_{-b} \longleftarrow v_b^* z_{-b},$$

before the next step

Scaling



$S_{escape}^{(t)}$ larger than $S_{AO}^{(t)}$, but still not whole space

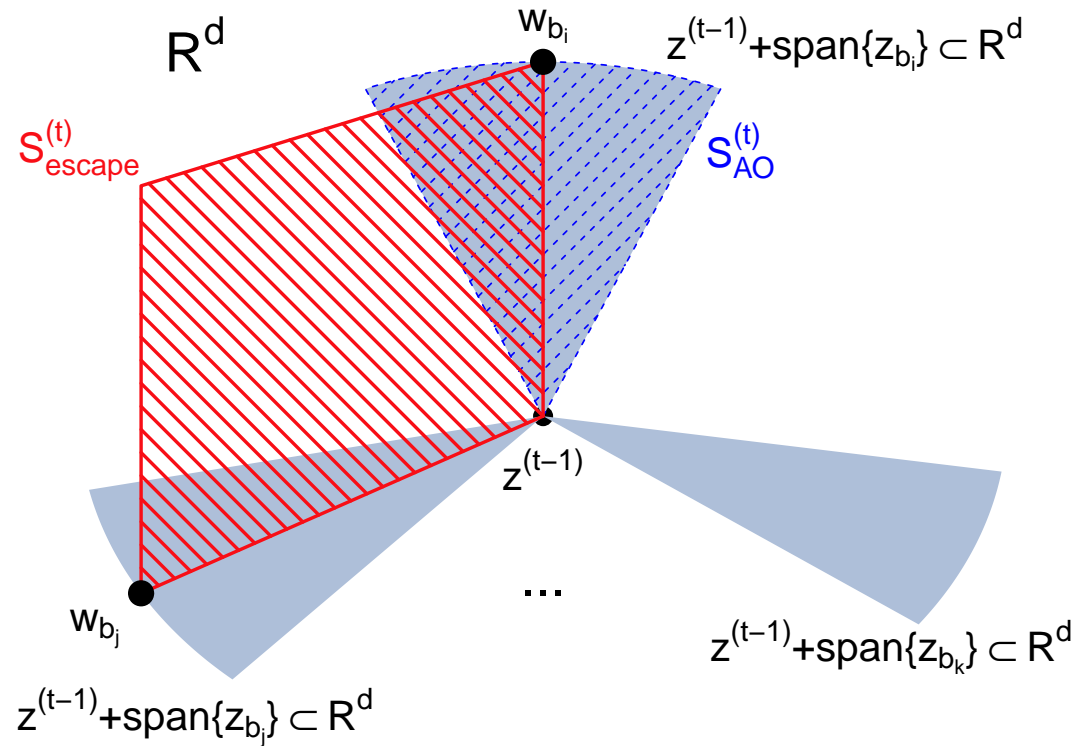
Restricted Joint Search (RJS)

$$\min_{\alpha_1, \dots, \alpha_B} f(z_1 + \alpha_1 w_1 I_1, z_2 + \alpha_2 w_2 I_2, \dots, z_B + \alpha_B w_B I_B)$$

$$w_b \in \text{span}\{z_b\}$$

$$I_b = \begin{cases} 1, & \text{if component } b \text{ chosen to "participate"}; \\ 0, & \text{otherwise} \end{cases}$$

Restricted Joint Search (RJS)



$S_{\text{escape}}^{(t)}$ never search by AO

RJS for Matrix Factorization

$$\min_{\alpha \in \mathbb{R}^n, \beta \in \mathbb{R}^m} \sum_{u,i \in \mathbb{T}} \left[r_{ui} - (\mathbf{p}_u + \alpha_u \mathbf{w}_p^u)^\top (\mathbf{q}_i + \beta_i \mathbf{w}_q^i) \right]^2 + \lambda \left[\sum_u \|\mathbf{p}_u + \alpha_u \mathbf{w}_p^u\|^2 + \sum_i \|\mathbf{q}_i + \beta_i \mathbf{w}_q^i\|^2 \right]$$

- search directions $\mathbf{w}_p^u, \mathbf{w}_q^i$ **mostly zero**
- nonzero ones selected ...
 - either randomly, e.g., $\mathbf{w}_p^u, \mathbf{w}_q^i \sim \mathbf{N}(\mathbf{0}, \mathbf{I})$
 - or greedily (next slide)

Greedy Choices of $\mathbf{w}_p^u, \mathbf{w}_q^i$

- for \mathbf{p}_u alone (fixing others), consider solving for α given \mathbf{w} :

$$\min_{\alpha} L(\alpha) = \sum_{i \in \mathbb{T}_u} [r_{ui} - (\mathbf{p}_u + \alpha \mathbf{w})^T \mathbf{q}_i]^2 + \lambda \|\mathbf{p}_u + \alpha \mathbf{w}\|^2$$

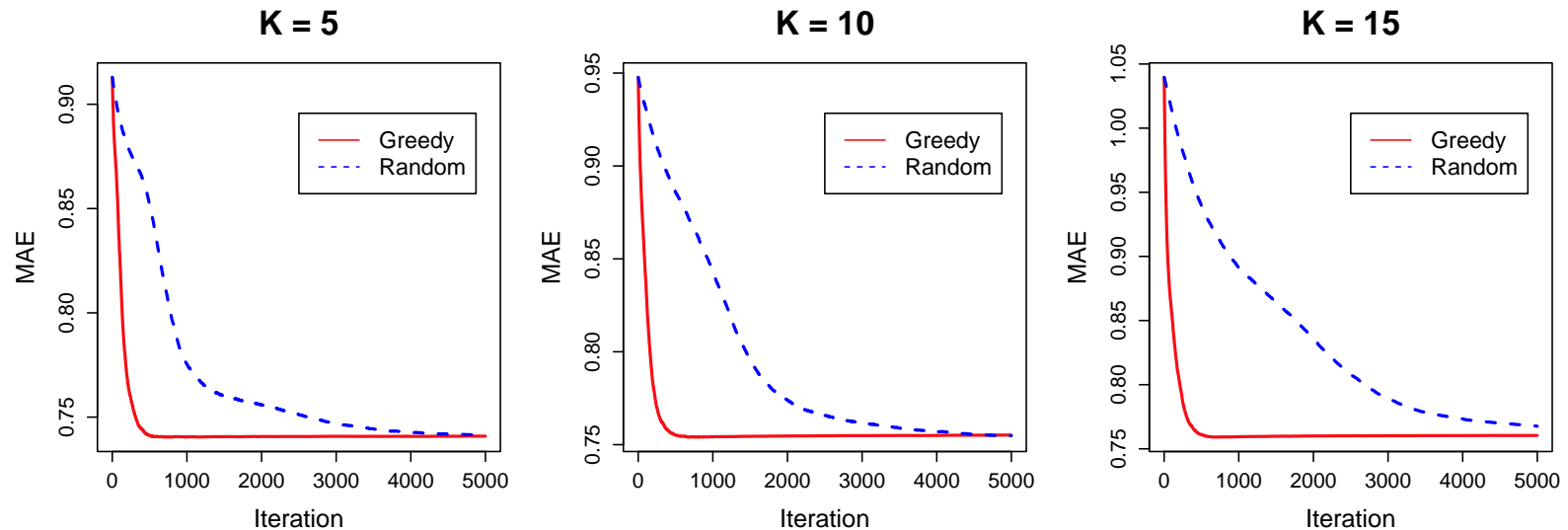
where $\mathbb{T}_u = \{i : r_{ui} \text{ is known}\}$

- optimal α must satisfy

$$\hat{\alpha}(\mathbf{w}) = \frac{\sum_{i \in \mathbb{T}_u} \mathbf{w}^T \mathbf{q}_i (r_{ui} - \mathbf{p}_u^T \mathbf{q}_i) - \lambda \mathbf{w}^T \mathbf{p}_u}{\sum_{i \in \mathbb{T}_u} (\mathbf{w}^T \mathbf{q}_i)^2 + \lambda \|\mathbf{w}\|^2}$$

- can plug $\hat{\alpha}(\mathbf{w})$ back into $L(\alpha)$ above and solve for “best” \mathbf{w} , e.g., with [quasi-Newton](#)

Example: Amazon Data



Mean absolute error (MAE) on the hold-out test set (average of 10 runs). Iterations start at convergence of baseline AO. Data source: McAuley & Leskovec (*RecSys 2013*). Original data contained ~ 35.3 M reviews from `amazon.com` (1995–2013). We took a dense subset consisting of users having rated ≥ 55 items & items with ≥ 24 ratings (~ 5.5 M reviews).

Selective Scaling: Regression w/ MCP

- for given j , let

$$E_j = \{k \neq j \text{ such that } |\text{corr}(\mathbf{x}_k, \mathbf{x}_j)| > \rho_{\min}\}$$

- solve

$$\min_{\beta_j, v} \left\| \mathbf{y} - \beta_j \mathbf{x}_j - \sum_{\ell \in E_j^C \setminus \{j\}} \beta_\ell \mathbf{x}_\ell - \sum_{k \in E_j} (v \beta_k) \mathbf{x}_k \right\|^2 +$$
$$J(\beta_j) + \sum_{\ell \in E_j^C \setminus \{j\}} J(\beta_\ell) + \sum_{k \in E_j} J(v \beta_k)$$

- can solve explicitly — first order condition piecewise linear in v on $< \infty$ intervals, so just check if each interval crosses 0

Summary

- main message:
 - optimization of nonconvex functions is hard
- specific content:
 - content-boosted matrix factorization
 - expanded alternating optimization