

Case Study 4: Collaborative Filtering

Review: Matrix Completion Alternating Least Squares

Machine Learning for Big Data
CSE547/STAT548, University of Washington

Emily Fox

February 13th, 2014

©Emily Fox 2014

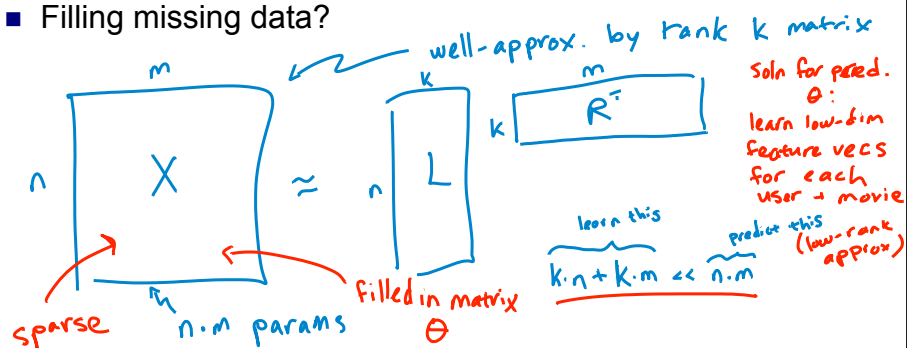
1

Matrix Completion Problem



X_{ij} known for black cells
 X_{ij} unknown for white cells
Rows index users
Columns index movies

- Filling missing data?



©Emily Fox 2014

2

Matrix Completion via Rank Minimization

- Given observed values: $(u, v, r_{uv}) \in X$ some $r_{uv} = ?$
- Find matrix $\hat{\Theta} \leftarrow$ filled in (not sparse) ↑
unobs. rating
- Such that: $\Theta_{uv} = r_{uv} \quad \forall r_{uv} \neq ?$ ← all obs. ratings
fit $r_{uv} \neq ?$ perfectly ← match ratings for all obs. values of X
- But... want Θ to be low-rank
- Introduce bias: one possible objective $\min \text{rank}(\Theta)$
 Θ s.t. $\Theta_{uv} = r_{uv} \quad \forall r_{uv} \neq ?$ ← for $k \leq \min(n, m)$
- Two issues: NP-hard
you can't hope to get exact matching

©Emily Fox 2014

3

Approximate Matrix Completion

- Minimize squared error:
 - (Other loss functions are possible) instead:
 $\min_{\Theta} \sum_{(u,v): r_{uv} \neq ?} (\Theta_{uv} - r_{uv})^2$ ← allow for some error
- Choose rank k :
 $\hat{\Theta} = L \hat{R}^T$ ← fix k
- Optimization problem:
 $\min_{L, R} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2$
 non-convex opt. problem ... local optima only

©Emily Fox 2014

4

Coordinate Descent for Matrix Factorization

$$\min_{L,R} \sum_{(u,v): r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2$$

- Fix movie factors R , optimize for user factors L
- First observation:

$$\min_{L_u, L_n} \sum_{(u,v): r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2 = \quad V_u = \text{set of movies user } u \text{ rated}$$

$$\min_{L_u, L_n} \sum_u \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2 = \quad \leftarrow \text{ind. opt. problem for each user}$$

$$\sum_u \min_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2 \quad \leftarrow \text{data parallel problem}$$

next slide

©Emily Fox 2014

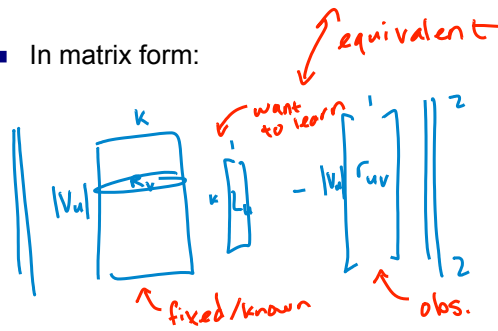
5

Minimizing Over User Factors

- For each user u : $\min_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2$

just consider one user

- In matrix form:



Think of as $\|X\beta - y\|_2^2$ normal LS problem

- Second observation: Solve by
 - matrix inversion
 - gradient methods
 - ...

©Emily Fox 2014

6

Coordinate Descent for Matrix Factorization: Alternating Least-Squares

$$\min_{L,R} \sum_{(u,v):r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\| + \lambda_v \|R\|$$

- Fix movie factors, optimize for user factors
 - ★ □ Independent least-squares over users $\min_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\|$
- Fix user factors, optimize for movie factors
 - ★ □ Independent least-squares over movies $\min_{R_v} \sum_{u \in U_v} (L_u \cdot R_v - r_{uv})^2 + \lambda_v \|R\|$
- System may be underdetermined: *use regularization*
- Converges to *local optima*

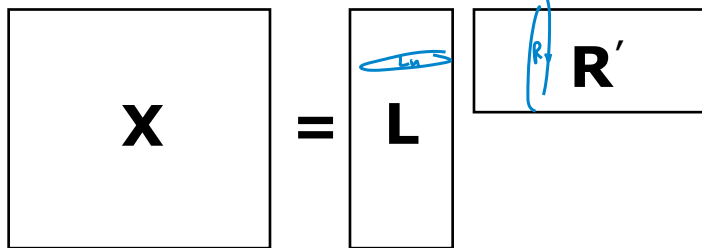
©Emily Fox 2014

7

Effect of Regularization

$$\|A\|_F = \sqrt{\sum_{i,j} A_{ij}^2}$$

$$\min_{L,R} \sum_{(u,v):r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\| + \lambda_v \|R\|$$



$\| \cdot \| = \| \cdot \|_F^2$
 each sub problem
 uses $\|L_u\|_2^2 \rightarrow$ ridge
 regression

$\| \cdot \| = \| \cdot \|_1$
 each subproblem $\|L_u\|_1$
 \rightarrow solved by LASSO
 methods

©Emily Fox 2014

8

What you need to know...

- Matrix completion problem for collaborative filtering
- Over-determined \rightarrow low-rank approximation
- Rank minimization is NP-hard
- Minimize least-squares prediction for known values for given rank of matrix
 - Must use regularization
- Coordinate descent algorithm = “Alternating Least Squares”

©Emily Fox 2014

9

Case Study 4: Collaborative Filtering

SGD for Matrix Completion
Matrix-norm Minimization

Machine Learning for Big Data
CSE547/STAT548, University of Washington

Emily Fox

February 13th, 2014

©Emily Fox 2014

10

Stochastic Gradient Descent

$$\min_{L,R} \frac{1}{2} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \frac{\lambda_u}{2} \|L\|_F^2 + \frac{\lambda_v}{2} \|R\|_F^2$$

- Observe one rating at a time $r_{uv}^{(t)}$ $\epsilon_t = L_u^{(t)} \cdot R_v^{(t)} - r_{uv}^{(t)}$

- Gradient observing r_{uv} :

$$\left. \begin{aligned} \frac{\partial F}{\partial L_u} &= \epsilon_t R_v + \lambda_u L_u \\ \frac{\partial F}{\partial R_v} &= \epsilon_t L_u + \lambda_v R_v \end{aligned} \right\} \nabla F_t = \begin{bmatrix} \epsilon_t R_v + \lambda_u L_u \\ \epsilon_t L_u + \lambda_v R_v \end{bmatrix}$$

- Updates: step size η_t : $\begin{bmatrix} L \\ R \end{bmatrix} \leftarrow \begin{bmatrix} L \\ R \end{bmatrix} - \eta_t \nabla F_t$

$$\begin{bmatrix} L_u^{(t+1)} \\ R_v^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} (1 - \eta_t \lambda_u) L_u^{(t)} - \eta_t \epsilon_t R_v^{(t)} \\ (1 - \eta_t \lambda_v) R_v^{(t)} - \eta_t \epsilon_t L_u^{(t)} \end{bmatrix} \leftarrow \text{fast + easy to implement}$$

©Emily Fox 2014

11

Local Optima v. Global Optima

- We are solving:

$$\min_{L,R} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\|_F^2 + \lambda_v \|R\|_F^2$$

- We (kind of) wanted to solve:

$$\min_{\theta} \text{rank}(\theta) \quad \theta_{uv} = r_{uv} \quad \forall (u,v, r_{uv}) \in X \text{ s.t. } r_{uv} \neq ?$$

- Which is NP-hard...

- How do these things relate???

©Emily Fox 2014

12

Eigenvalue Decompositions for PSD Matrices

- Given a (square) symmetric positive semidefinite matrix: $\theta \succeq 0$

□ Eigenvalues: $\lambda_1, \dots, \lambda_d \geq 0$ $\lambda = (\lambda_1, \dots, \lambda_d)$

- Thus rank is:

$$|\{\lambda_i : \lambda_i > 0\}| \equiv \text{rank}(\theta) \equiv \|\lambda\|_0$$

- Approximation:

$$\|\lambda\|_0 \approx \|\lambda\|_1 = \sum_{i=1}^d |\lambda_i| \stackrel{\text{PSD}}{=} \sum_{i=1}^d \lambda_i \quad \leftarrow \text{L}_1 \text{ norm is sum of eigvals}$$

- Property of trace:

$$\text{trace}(\theta) = \sum_{i=1}^d \lambda_i$$

- Thus, approximate rank minimization by:

$$\begin{array}{l} \min_{\theta} \text{rank}(\theta) = \|\lambda\|_0 \\ \theta \succeq 0 \quad \theta_{uv} = r_{uv} \end{array} \quad \approx \quad \begin{array}{l} \min_{\theta} \text{trace}(\theta) \\ \theta \succeq 0 \quad \theta_{uv} = r_{uv} \end{array}$$

©Emily Fox 2014

13

Generalizing the Trace Trick

- Non-square matrices ain't got no trace

- For (square) positive semidefinite matrices, matrix factorization: $\lambda_i \geq 0$

$$\theta = P \Lambda P^{-1} \quad \text{diag}(\lambda)$$

- For rectangular matrices, singular value decomposition: (SVD)

$$\theta = U \Sigma V^T$$

diagonal matrix w/ entries $\sigma_i(\theta) \geq 0$
 \uparrow
*i*th singular value

- Nuclear norm:

$$\|\theta\|_* = \sum_{i=1}^m \sigma_i(\theta)$$

$$\begin{array}{l} \min_{\theta} \|\theta\|_* \\ \text{s.t. } \theta_{uv} = r_{uv} \end{array}$$

Convex problem!

©Emily Fox 2014

14

Nuclear Norm Minimization

- Optimization problem: (relaxation)

new

$$\min_{\Theta} \|\Theta\|_*$$

$$\Theta_{uv} = r_{uv}$$

can have no feasible soln

- Possible to relax equality constraints: (relaxation of relaxation)

$$\min_{\Theta, r_{uv}} \sum (\Theta_{uv} - r_{uv})^2 + \lambda \|\Theta\|_*$$

- Both are convex problems!
(solved by semidefinite programming)

Analysis of Nuclear Norm

- Nuclear norm minimization = convex relaxation of rank minimization:

$$\min_{\Theta} \text{rank}(\Theta) \xleftarrow{\text{NP-hard}} \approx \min_{\Theta} \|\Theta\|_* \xleftarrow{\text{convex relaxation}} \xleftarrow{\text{has polytime soln}}$$

$$r_{uv} = \Theta_{uv}, \forall r_{uv} \in X, r_{uv} \neq? \quad r_{uv} = \Theta_{uv}, \forall r_{uv} \in X, r_{uv} \neq?$$

- Theorem [Candes, Recht '08]:

- If there is a true matrix of rank k ,
- And, we observe at least

$$C \cdot k \cdot n^{1.2} \log n$$

random entries of true matrix

Original problem has $n \cdot m$ entries

need $\approx k n^{1.2}$ obs

assuming $n \geq m$

we have $k \cdot n + k \cdot m$ params

- Then true matrix is recovered exactly with high probability via convex nuclear norm minimization!
- Under certain conditions

Nuclear Norm Minimization versus Direct (Bilinear) Low Rank Solutions

- Nuclear norm minimization: $\min_{\Theta} \sum_{r_{uv}} (\Theta_{uv} - r_{uv})^2 + \lambda \|\Theta\|_*$ (*)
 - Annoying because:
 - Θ very large (8B entries in Netflix)
 - SDP solvers are very slow (but polytime)
- Instead: $\min_{L,R} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\|_F^2 + \lambda_v \|R\|_F^2$ (***)
 - Annoying because: - many local optima
 - But $\|\Theta\|_* = \inf \left\{ \frac{1}{2} \|L\|_F^2 + \frac{1}{2} \|R\|_F^2 : \Theta = LR \right\}$
 - So (***) is a non-convex approx to (*)
 - And if we pick rank of $L \cdot R'$ to be slightly higher than $\text{rank}(\Theta^*)$, local optima of (***) are global optima of (*)

Under certain conditions [Burer, Monteiro '04]

©Emily Fox 2014

17

What you need to know...

- Stochastic gradient descent for matrix factorization
- Norm minimization as convex relaxation of rank minimization
 - Trace norm for PSD matrices
 - Nuclear norm in general
- Intuitive relationship between nuclear norm minimization and direct (bilinear) minimization

©Emily Fox 2014

18

Case Study 4: Collaborative Filtering

Nonnegative Matrix Factorization Projected Gradient

Machine Learning for Big Data
CSE547/STAT548, University of Washington

Emily Fox

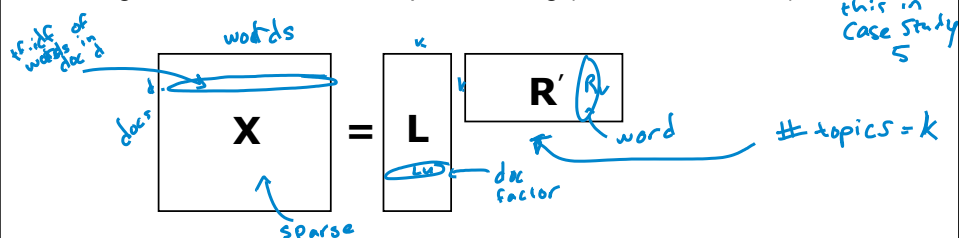
February 13th, 2014

©Emily Fox 2014

19

Matrix factorization solutions can be unintuitive...

- Many, many, many applications of matrix factorization
- E.g., in text data, can do topic modeling (alternative to LDA):



- Would like:
 L_u : how much is doc u about each topic
 R_v : how much a word v contributes to each topic
- But...
Standard matrix factorization: L_u, R_v can be negative

©Emily Fox 2014

20

Nonnegative Matrix Factorization

$$\mathbf{X} = \mathbf{L} \mathbf{R}'$$

- Just like before, but

$$\min_{L \geq 0, R \geq 0} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|\mathbf{L}\|_F^2 + \lambda_v \|\mathbf{R}\|_F^2$$

non-neg.
L, R

- Constrained optimization problem
 - Many, many, many, many solution methods... we'll check out a simple one

©Emily Fox 2014

21

Projected Gradient

- Standard optimization:
 - Want to minimize: $\min_{\Theta} f(\Theta)$
 - Use gradient updates:

$$\Theta^{(t+1)} \leftarrow \Theta^{(t)} - \eta_t \nabla f(\Theta^{(t)})$$

- Constrained optimization:
 - Given convex set \mathcal{C} of feasible solutions
 - Want to find minima within \mathcal{C} : $\min_{\Theta \in \mathcal{C}} f(\Theta)$

- Projected gradient:
 - Take a gradient step (ignoring constraints):

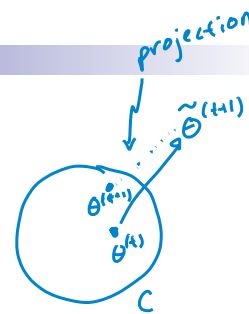
$$\tilde{\Theta}^{(t+1)} \leftarrow \Theta^{(t)} - \eta_t \nabla f(\Theta^{(t)})$$

Projection into feasible set:

$$\Pi_{\mathcal{C}}(\Theta) = \operatorname{argmin}_{\beta \in \mathcal{C}} \|\Theta - \beta\|_2$$

$$\Theta^{(t+1)} = \Pi_{\mathcal{C}}(\tilde{\Theta}^{(t+1)})$$

} often easy to compute (always convex)



©Emily Fox 2014

22

Projected Stochastic Gradient Descent for Nonnegative Matrix Factorization

$$\min_{L \geq 0, R \geq 0} \frac{1}{2} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \frac{\lambda_u}{2} \|L\|_F^2 + \frac{\lambda_v}{2} \|R\|_F^2$$

- Gradient step observing r_{uv} ignoring constraints:

$$\begin{bmatrix} \tilde{L}_u^{(t+1)} \\ \tilde{R}_v^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} (1 - \eta_t \lambda_u) L_u^{(t)} - \eta_t \epsilon_t R_v^{(t)} \\ (1 - \eta_t \lambda_v) R_v^{(t)} - \eta_t \epsilon_t L_u^{(t)} \end{bmatrix}$$

- Convex set: $L_u \geq 0 \quad R_v \geq 0 \quad \forall u, v$
- Projection step: $\Pi_C(\theta) = \underset{\beta \in C}{\operatorname{argmin}} \|\theta - \beta\|_2^2 \leftarrow$ totally ind. problems per dimension

Single dimension:

$$= \underset{\beta \geq 0}{\operatorname{argmin}} (\theta - \beta)^2$$

$$= \begin{cases} \theta & \text{if } \theta \geq 0 \\ 0 & \text{if } \theta < 0 \end{cases} = (\theta)_+$$

$$\begin{bmatrix} L_u^{(t+1)} \\ R_v^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} \tilde{L}_u^{(t+1)} \\ \tilde{R}_v^{(t+1)} \end{bmatrix}$$

set all neg. coords to zero
easy!!!
(phev)

©Emily Fox 2014

23

What you need to know...

- In many applications, want factors to be nonnegative
- Corresponds to constrained optimization problem
- Many possible approaches to solve, e.g., projected gradient

©Emily Fox 2014

24

Case Study 4: Collaborative Filtering

Cold Start Problem

Machine Learning for Big Data
CSE547/STAT548, University of Washington

Emily Fox

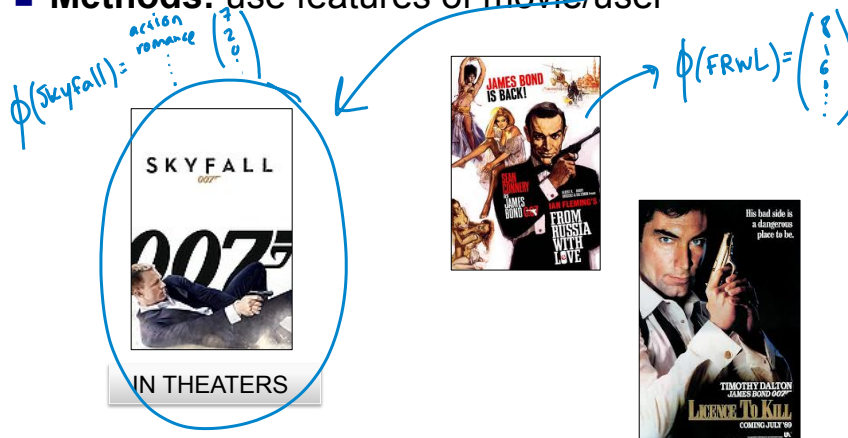
February 13th, 2014

©Emily Fox 2014

25

Cold-Start Problem

- **Challenge:** Cold-start problem (new movie or user)
- **Methods:** use features of movie/user



©Emily Fox 2014

26

Cold-Start Problem More Formally

- Consider a new user u' and predicting that user's ratings

- No previous observations

$r_{u'v} = ? \quad \forall v$ want $L_{u'}$ to predict ratings for this user

- Objective considered so far:

$$\min_{L,R} \frac{1}{2} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \frac{\lambda_u}{2} \|L\|_F^2 + \frac{\lambda_v}{2} \|R\|_F^2$$

all $r_{uv} \neq ?$
(only obs. values)

doesn't depend on $L_{u'}$ (no obs. $r_{u'v}$)

- Optimal user factor:

$L_{u'} = 0$ only penalty term in opt. $L_{u'}$ (min $\|L_{u'}\|_2^2$)

only term appearing in ALS step for u'

- Predicted user ratings:

always predict $r_{u'v} = 0 \quad \forall v$... problem

©Emily Fox 2014

27

An Alternative Formulation

- A simpler model for collaborative filtering

- We would not have this issue if we assumed all users were identical

- i.e. all users shared a feature vector w

- w informed by all ratings + can use it for new user u'

- What about for new movies? What if we had side information?

Create a movie feature vector: \rightarrow synopsis/info

$\phi(v) = (\text{'action', 1994, Tarantino, ...})$
genre year director \leftarrow dim k

- What dimension should w be?

- Fit linear model:

\leftarrow same length as movie feature vector

For all users, u , $r_{uv} \approx w \cdot \phi(v)$

\leftarrow want to learn this

- Minimize:

$$\min_w \sum_{r_{uv}} (w \cdot \phi(v) - r_{uv})^2 + \lambda_w \|w\|$$

\leftarrow LS Lasso

©Emily Fox 2014

28

Personalization

- If we don't have any observations about a user, use wisdom of the crowd
 - Address cold-start problem

For user u' , predict $r_{u'v} \approx w \cdot \phi(v)$

- Clearly, not all users are the same ... shared w is strong assumption
- Just as in personalized click prediction, consider model with global and user-specific parameters

Consider user-specific deviations w_u from the crowd w

$r_{uv} \approx (w + w_u) \cdot \phi(v)$
init. to 0

- As we gain more information about the user, forget the crowd
 w_u more informed

User Features...

- In addition to movie features, may have information about the user:

$\phi(u) = (25, F, MS, A^+, \dots)$
age, gender, education, grade in Big Data

- Combine with features of movie:

$\phi(u, v) = (\dots, \phi(u), \dots, \dots, \phi(v), \dots, \dots)$
cross features...

- Unified linear model:

$r_{uv} \approx (w + w_u) \cdot \phi(u, v)$

Feature-based Approach versus Matrix Factorization

- Feature-based approach:
 - Feature representation of user and movies fixed
 - Can address cold-start problem

- Matrix factorization approach:
 - Suffers from cold-start problem
 - User & movie features are learned from data

- A unified model:

©Emily Fox 2014

31

Unified Collaborative Filtering via SGD

$$\min_{L,R,w,\{w_u\}_u} \frac{1}{2} \sum_{r_{uv}} (L_u \cdot R_v + (w + w_u) \cdot \phi(u, v) - r_{uv})^2 + \frac{\lambda_u}{2} \|L\|_F^2 + \frac{\lambda_v}{2} \|R\|_F^2 + \frac{\lambda_w}{2} \|w\|_2^2 + \frac{\lambda_{wu}}{2} \sum_u \|w_u\|_2^2$$

- Gradient step observing r_{uv}
 - For L,R
$$\begin{bmatrix} L_u^{(t+1)} \\ R_v^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} (1 - \eta_t \lambda_u) L_u^{(t)} - \eta_t \epsilon_t R_v^{(t)} \\ (1 - \eta_t \lambda_v) R_v^{(t)} - \eta_t \epsilon_t L_u^{(t)} \end{bmatrix}$$
 - For w and w_u :

©Emily Fox 2014

32

What you need to know...

- Cold-start problem
- Feature-based methods for collaborative filtering
 - Help address cold-start problem
- Unified approach