

Case Study 4: Collaborative Filtering

Collaborative Filtering Matrix Completion Alternating Least Squares

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

Emily Fox

May 7th, 2015

©Emily Fox 2015

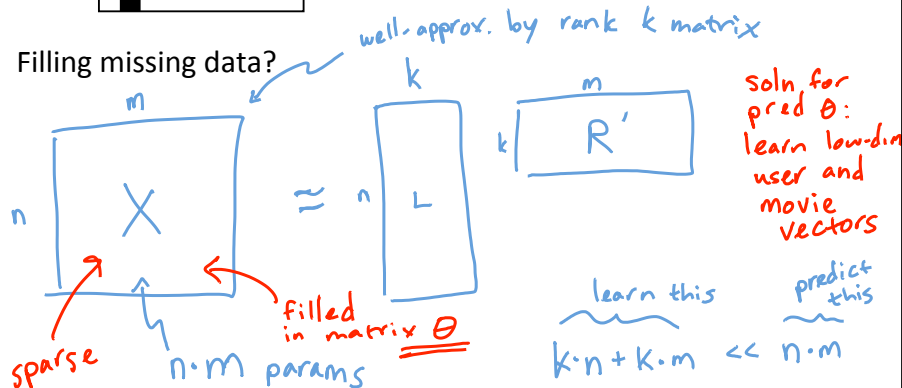
1

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 2015

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)
- Such that: $\Theta_{uv} = r_{uv} \quad \forall r_{uv} \neq ?$
fit $r_{uv} \neq ?$ perfectly match ratings observed in X
- But... want low-rank Θ ★
- Introduce bias:
one possible objective $\left\{ \begin{array}{l} \min \text{rank}(\Theta) \\ \Theta \text{ s.t. } \Theta_{uv} = r_{uv} \quad \forall r_{uv} \neq ? \end{array} \right.$
- Two issues: $\left\{ \begin{array}{l} \text{NP-hard} \\ \text{you can't hope to get exact matching} \end{array} \right.$

©Emily Fox 2015

3

Approximate Matrix Completion

- Minimize squared error:
 – (Other loss functions are possible) relax hard constraints

$$\min_{\Theta} \sum_{(u,v): r_{uv} \neq ?} (\Theta_{uv} - r_{uv})^2$$
allow for some error
- Choose rank k :

$$\hat{\Theta} = \hat{L} \hat{R}' \quad \leftarrow \text{fix rank } k$$
- Optimization problem:

$$\min_{L, R} \sum_{r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2$$
non-convex opt. problem ... local optima only

©Emily Fox 2015

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:

$$\min_{L_1, \dots, L_n} \sum_{(u,v): r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2$$

$V_u \triangleq$ set of movies user u rated

$$= \min_{L_1, \dots, L_n} \sum_u \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2$$

ind. opt. problem for each user

$$= \sum_u \min_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2$$

data parallel problem

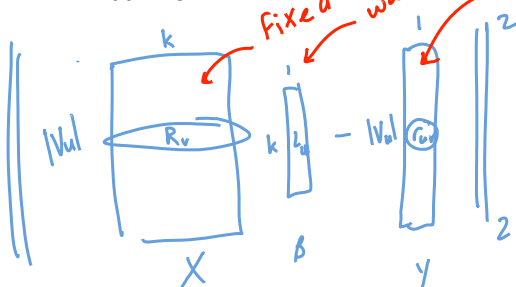
©Emily Fox 2015

5

Minimizing Over User Factors

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

- In matrix form:



Think of as

$$\|XB - y\|_2^2$$

normal LS problem

- Second observation: Solve by

- matrix inversion
- gradient methods

©Emily Fox 2015

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 2015

7

Effect of Regularization

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

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

$$X = L R'$$

If $\|\cdot\| = \|\cdot\|_F^2$
each subproblem
uses $\|L_u\|_2^2 \rightarrow$ ridge
regression

If $\|\cdot\| = \|\cdot\|_1$,
each subproblem $\|L_u\|_1$
 \rightarrow solved via
LASSO methods

©Emily Fox 2015

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 2015

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

May 7th, 2015

©Emily Fox 2015

10

Stochastic Gradient Descent

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

- 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} :

$$\frac{\partial F}{\partial L_u} = \epsilon_t R_v + \lambda_u L_u$$

$$\frac{\partial F}{\partial R_v} = \epsilon_t L_u + \lambda_v R_v$$

$$\nabla F_t = \begin{bmatrix} \epsilon_t R_v^{(t)} + \lambda_u L_u^{(t)} \\ \epsilon_t L_u^{(t)} + \lambda_v R_v^{(t)} \end{bmatrix}$$

- Updates:

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

← step size
← fast + easy to implement

©Emily Fox 2015

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:

$$\min_{\theta} \text{rank}(\theta)$$

$$\theta_{uv} = r_{uv} \quad \forall (u,v,r_{uv}) \in X, r_{uv} \neq ?$$

- Which is NP-hard...

– How do these things relate???

©Emily Fox 2015

12

Eigenvalue Decompositions for PSD Matrices

- Given a (square) symmetric positive semidefinite matrix:

– 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:

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

- Property of trace:

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

- Thus, approximate rank minimization by:

$$\begin{array}{ll} \min_{\theta} & \text{rank}(\theta) = \|\lambda\|_0 \\ \text{s.t.} & \theta_{uv} = r_{uv} \\ & \theta \succeq 0 \end{array}$$

$$\approx \begin{array}{ll} \min_{\theta} & \text{trace}(\theta) = \|\lambda\|_1 \\ \text{s.t.} & \theta_{uv} = r_{uv} \\ & \theta \succeq 0 \end{array}$$

©Emily Fox 2015

13

Generalizing the Trace Trick

- Non-square matrices ain't got no trace

- For (square) positive semidefinite matrices, eigendecomposition:

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

- For rectangular matrices, singular value decomposition:

$$\begin{array}{c} m \\ \theta \end{array} = \begin{array}{c} n \\ U \end{array} \begin{array}{c} n \\ \Sigma \end{array} \begin{array}{c} m \\ V' \end{array}$$

diagonal matrix
w/ entries
 $\sigma_i(\theta) \geq 0$
singular values

- Nuclear norm:

$$\|\theta\|_* = \sum_i \sigma_i(\theta)$$

↑
nuclear norm

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

convex problem!

©Emily Fox 2015

14

Nuclear Norm Minimization

- Optimization problem:

new

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

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

no feasible soln

- Possible to relax equality constraints:

(relaxation of relaxation)

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

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

©Emily Fox 2015

15

Analysis of Nuclear Norm

- Nuclear norm minimization = convex relaxation of rank minimization:

$$\min_{\Theta} \text{rank}(\Theta)$$

NP-hard

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

convex relaxation

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

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

- Theorem [Candes, Recht '08]:

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

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

random entries of true matrix

original problem has $n \cdot m$ entries

assuming $n \geq m$

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

©Emily Fox 2015

16

Nuclear Norm Minimization vs. Direct (Bilinear) Low Rank Solutions

- Nuclear norm minimization: $\min_{\Theta} \sum_{r_{uv}} (\Theta_{uv} - r_{uv})^2 + \lambda \|\Theta\|_*$ (*)
convex, global opt. close to truth
 - Annoying because:
 - Θ very large (8B entries in Netflix)
 - SDP solvers are very slow (but polytime)
 - Instead: $\min_{L,R} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\|_F^2 + \lambda_v \|R\|_F^2$ (**)
 $\sim 10-100$ M params, but very fast solvers
 - Annoying because: *many local optima*
 - But $\|\Theta\|_* = \inf \left\{ \min_{L,R} \frac{1}{2} \|L\|_F^2 + \frac{1}{2} \|R\|_F^2 : \Theta = LR' \right\}$
 - So (**) is a non-convex approx to (*)
 - And if we pick rank of $L \cdot R$ to be slightly larger than $\text{rank}(\Theta^*)$, local optima of (**) are global optima of (*)
- Under certain conditions [Burer, Monteiro '04] ©Emily Fox 2015 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

Case Study 4: Collaborative Filtering

Nonnegative Matrix Factorization Projected Gradient

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

Emily Fox

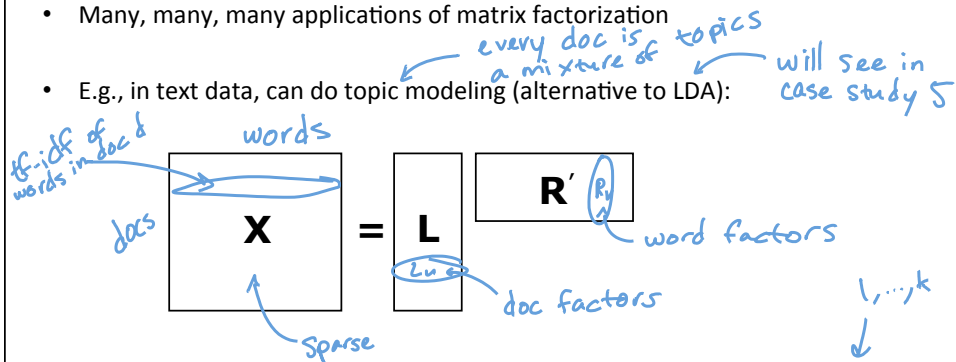
May 6th, 2015

©Emily Fox 2015

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 a doc is about each topic
- R_v : how much a word contributes to a topic

- But... Standard matrix factorization: L_u, R_v can be negative

©Emily Fox 2015

20

Nonnegative Matrix Factorization

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

- Just like before, but

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

non-negative
L, R

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

©Emily Fox 2015

21

Recall: Projected Gradient

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

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

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

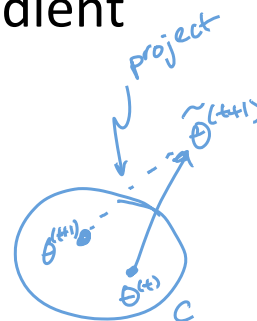
- Projected gradient:
 - Take a gradient step (ignoring constraints):

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

- Projection into feasible set:

$$\Pi_C(\Theta) \equiv \arg \min_{B \in C} \|\Theta - B\|_2^2 \leftarrow \text{often easy to compute (always convex)}$$

$$\Theta^{(t+1)} = \Pi_C(\tilde{\Theta}^{(t+1)})$$



©Emily Fox 2015

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:

don't necessarily satisfy pos. 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:

$\Pi_C(\theta) = \arg \min_{\beta \in C} \|\theta - \beta\|_2^2 \leftarrow$ totally ind. problems per dimension

Single dim: $\arg \min_{\beta \geq 0} (\theta - \beta)^2 = \begin{cases} \theta & \text{if } \theta \geq 0 \\ 0 & \text{if } \theta < 0 \end{cases} = (\theta)_+$

threshold

set all neg. coord. to 0

easy

©Emily Fox 2015

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

Case Study 4: Collaborative Filtering

Cold Start Problem

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

Emily Fox

May 6th, 2015

©Emily Fox 2015

25

Cold-Start Problem

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



$\phi_{\text{RWL}} = \begin{pmatrix} 8 \\ 1 \\ 6 \\ 0 \\ \vdots \end{pmatrix}$ action
romance



©Emily Fox 2015

26

Cold-Start Problem More Formally

- Consider a new user u' and predicting that user's ratings
 - No previous observations

$$r_{u'v} = ? \quad \forall v$$

- Objective considered so far:

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

does not depend on $L_{u'}$ (no obs. ratings)

only term that appears in ALS step

$L_{u'}$

- Optimal user factor:

$$L_{u'} = 0$$

only penalty term is present

- Predicted user ratings:

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

©Emily Fox 2015

27

An Alternative Formulation

- A simpler model for collaborative filtering

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

- If all users shared a feature vector w , with w informed by all ratings, then can use w for new user

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

create movie feature vector

$$\phi(v) = (\text{'action' , 1994 , Tarantino ...})$$

genre , year , director ...

- What dimension should w be?

same length as movie feature vector

- Fit linear model:

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

- Minimize:

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

only 1 param

fixed

LS, Lasso, ridge

©Emily Fox 2015

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
- 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

©Emily Fox 2015

29

User Features...

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

$$\phi(u) = (25, F, MSc, A^+, \dots)$$

age gender education big data grade

- Combine with features of movie:

$$\phi(u, v) = (\phi(u) \dots \phi(v) \dots \text{cross features})$$

- Unified linear model:

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

©Emily Fox 2015

30

Feature-Based Approach vs. Matrix Factorization

- Feature-based approach:
 - Feature representation of user and movies fixed ← *important side info* $\phi(u,v)$
 - Can address cold-start problem
- Matrix factorization approach:
 - Suffers from cold-start problem
 - User & movie features are learned from data L_u, R_v
- A unified model: *combine both ideas*

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

solve via ALS
SGD
⋮

©Emily Fox 2015

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 $r_{uv}^{(t)}$ $\epsilon_t = L_u^{(t)} \cdot R_v^{(t)} + (w^{(t)} + w_u^{(t)}) \cdot \phi(u, v) - r_{uv}^{(t)}$
 - 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 : $\nabla F^{(t)} = \epsilon_t \phi(u, v) + \lambda_w w^{(t)}$

$$\begin{bmatrix} w^{(t+1)} \\ w_u^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} (1 - \eta_t \lambda_w) w^{(t)} - \eta_t \epsilon_t \phi(u, v) \\ (1 - \eta_t \lambda_{w_u}) w_u^{(t)} - \eta_t \epsilon_t \phi(u, v) \end{bmatrix}$$

↖ only update w_u for user u in $r_{uv}^{(t)}$

©Emily Fox 2015

32

What you need to know...

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