

Introduction to Data Mining

Lecture #18: Recommendation 2 -Latent Factor Model

U Kang Seoul National University

U Kang



In This Lecture

- Learn the weight learning approach for collaborative filtering
- Understand the main idea of latent factor model
- Learn the advanced techniques for latent factor model, including regularization and bias extension



Outline

Netflix Prize; Weight Learning in CF

- Latent Factor Model
- □ Regularization for LF
- Bias Extension for LF
- □ Netflix Challenge



The Netflix Prize



Training data

- 100 million ratings, 480,000 users, 17,770 movies
- G years of data: 2000-2005

Test data

- Last few ratings of each user (2.8 million)
- Evaluation criterion: Root Mean Square Error (RMSE) =

$$\frac{1}{|R|} \sqrt{\sum_{(i,x)\in R} (\hat{r}_{xi} - r_{xi})^2}$$

Netflix's system RMSE: 0.9514

Competition

- 2,700+ teams
- \$1 million prize for 10% improvement on Netflix



The Netflix Utility Matrix R

480,000 users

Matrix R

17,700 movies

						-
1	1	3	4			
		3	5			5
			4	5		5
			3			
S			3			
	2			2		2
					5	
		2	1			1
		3			3	
	1					



Utility Matrix *R***: Evaluation**





BellKor Recommender System

- The winner of the Netflix Challenge!
- Multi-scale modeling of the data:

Combine top level, "regional" modeling of the data, with a refined, local view:

- **Global:**
 - Overall deviations of users/movies
- Factorization:
 - Addressing "regional" effects
- Collaborative filtering:
 - Extract local patterns



Modeling Local & Global Effects

Global:

- Mean movie rating: 3.7 stars
- □ The Sixth Sense is **0.5** stars above avg.
- □ Joe rates 0.2 stars below avg.
 ⇒ Baseline estimation: Joe will rate The Sixth Sense 4 stars
- Local neighborhood (CF/NN):
 - Joe didn't like related movie Signs
 - $\Box \Rightarrow Final estimate:$
 - Joe will rate The Sixth Sense 3.8 stars







- Earliest and most popular collaborative filtering method
 - Infer unknown ratings from those of "similar" movies (it em-item variant)
 - Define similarity measure s_{ii} of items i and j
 - Select k-nearest neighbors, compute the rating
 - N(i; x): items most similar to i that were rated by x

$$\hat{r}_{xi} = \frac{\sum_{j \in N(i;x)} S_{ij} \cdot r_{xj}}{\sum_{j \in N(i;x)} S_{ij}}$$

s_{ij}... similarity of items *i* and *j* r_{xj}...rating of user *x* on item *j N(i;x)*... set of items similar to item *i* that were rated by *x*



Modeling Local & Global Effects

■ In practice we get better estimates if we model deviations: $\hat{r}_{xi} = b_{xi} + \frac{\sum_{j \in N(i;x)} s_{ij} \cdot (r_{xj} - b_{xj})}{\sum_{j \in N(i;x)} s_{ij}}$

baseline estimate for r_{xi}

$$\boldsymbol{b}_{xi} = \boldsymbol{\mu} + \boldsymbol{b}_x + \boldsymbol{b}_i$$

- μ = overall mean rating
- $\boldsymbol{b}_{\boldsymbol{x}}$ = rating deviation of user \boldsymbol{x}
 - = (avg. rating of user \mathbf{x}) $\boldsymbol{\mu}$

$$\boldsymbol{b}_i = (avg. rating of movie \boldsymbol{i}) - \boldsymbol{\mu}$$

Problems/Issues:

Similarity measures are "arbitrary"
 Taking a weighted average can be restricting
 Solution: Instead of s_{ij} use w_{ij} that we learn from data



Idea: Interpolation Weights w_{ij}

Use a weighted sum rather than weighted avg.:

$$\widehat{r_{xi}} = b_{xi} + \sum_{j \in N(i;x)} w_{ij} (r_{xj} - b_{xj})$$

A few notes:

- N(i; x) ... set of movies rated by user x that are similar to movie i
- \square w_{ij} is the interpolation weight (some real number)
 - We allow: $\sum_{j \in N(i,x)} w_{ij} \neq 1$
- *w_{ij}* models interaction between pairs of movies (it does not depend on user *x*)



Idea: Interpolation Weights w_{ij}

- $\widehat{r_{xi}} = b_{xi} + \sum_{j \in N(i,x)} w_{ij} (r_{xj} b_{xj})$
- How to set w_{ij}?

• Remember, error metric is: $\frac{1}{|R|} \sqrt{\sum_{(i,x) \in R} (\hat{r}_{xi} - r_{xi})^2}$ or e quivalently SSE: $\sum_{(i,x) \in R} (\hat{r}_{xi} - r_{xi})^2$

- Find w_{ij} that minimize SSE on training data!
 - Models relationships between item *i* and its neighbors *j*
- w_{ij} can be learned/estimated based on x and all other users that rated i

Why is this a good idea?

Recommendations via Optimization

• Goal: Make good recommendations

- Quantify goodness using RMSE:
 Lower RMSE ⇒ better recommendations
- Want to make good recommendations on items that user has not yet seen. Very difficult task!
- Let's build a system such that it works well on known (user, item) ratings
 And hope the system will also predict well the unknown ratings



Recommendations via Optimization

- Idea: Let's set values w such that they work well on known (user, item) ratings
- How to find such values w?
- Idea: Define an objective function and solve the optimization problem

Find w_{ii} that minimize SSE on training data!

$$J(w) = \sum_{x,i} \left(\left[b_{xi} + \sum_{\substack{j \in N(i;x) \\ \text{Predicted rating}}} w_{ij} (r_{xj} - b_{xj}) \right] - r_{xi} \right)^2$$

True rating
Think of *w* as a vector of numbers



Detour: Minimizing a function

- A simple way to minimize a function *f*(*x*):
 - Take a gradient ∇f
 - Start at some point y and evaluate $\nabla f(y)$
 - Make a step in the reverse direction of the gradient: $y = y - \nabla f(y)$
 - Repeat until converged



15



Interpolation Weights

We have the optimization problem, now what?

$$J(w) = \sum_{x,i} \left(\left[b_{xi} + \sum_{k \in N(i;x)} w_{ik} (r_{xk} - b_{xk}) \right] - r_{xi} \right)^2$$

- Gradient decent:
 - □ Iterate until convergence: $w \leftarrow w \eta \nabla_w J$ $\eta \dots$ learning rate
 - where $\nabla_w J$ is the gradient:

$$\nabla_{w}J = \left[\frac{\partial J(w)}{\partial w_{ij}}\right] = 2\sum_{x,i} \left(\left[b_{xi} + \sum_{k \in N(i;x)} w_{ik}(r_{xk} - b_{xk}) \right] - r_{xi} \right) (r_{xj} - b_{xj})$$

for $j \in \{N(i; x), \forall i, \forall x\}$
else $\frac{\partial J(w)}{\partial w_{ij}} = \mathbf{0}$

■ Note: We fix movie *i*, go over all r_{xi} , for every movie $j \in N(i; x)$, we compute $\frac{\partial J(w)}{\partial w_{ij}}$ while $|w_{new} - w_{old}| > \varepsilon$: UKang $w_{new} = w_{old} - \eta \cdot \nabla w_{old}$ 16



Interpolation Weights

• So far:
$$\widehat{r_{xi}} = b_{xi} + \sum_{j \in N(i;x)} w_{ij} (r_{xj} - b_{xj})$$

- Weights *w_{ij}* learned based on their role; no use of an arbitrary similarity measure (*w_{ij}* ≠ *s_{ij}*)
- Explicitly account for interrelationships among the neighboring movies

Next: Latent factor model

Extract "regional" correlations



Performance of Various Methods

Global average: 1.1296

User average: 1.0651 Movie average: 1.0533 Netflix: 0.9514 **Basic Collaborative filtering: 0.94 CF+Biases+learned weights: 0.91** (Collaborative filtering ++) Grand Prize: 0.8563



Outline

- Metflix Prize; Weight Learning in CF
- 🔷 🔲 Latent Factor Model
 - □ Regularization for LF
 - Bias Extension for LF
 - □ Netflix Challenge





Latent Factor Models

SVD: $A = U \Sigma V^T$

• "SVD" on Netflix data: $\mathbf{R} \approx \mathbf{Q} \cdot \mathbf{P}^{T}$



For now let's assume we can approximate the rating matrix *R* as a product of "thin" *Q* · *P*^T

R has missing entries but let's ignore that for now!

 Basically, we will want the reconstruction error to be small on known ratings and we don't care about the values on the missing ones



Ratings as Products of Factors

How to estimate the missing rating of











fa	2.1	4	.6	1.7	2.4	.9	3	.4	.8	.7	6	.1
• actc	8	.7	.5	1.4	.3	-1	1.4	2.9	7	1.2	1	1.3
rs	1.1	2	.3	.5	-2	5	.8	4	.3	1.4	2.4	9



Ratings as Products of Factors

How to estimate the missing rating of



items	.1	4	.2				
	5	.6	.5				
	2	.3	.5				
	1.1	2.1	.3				
	7	2.1	-2				
	-1	.7	.3				
factors							



 $\hat{r}_{xi} = q_i \cdot p_x$ $\mathbf{q}_{if} \cdot \mathbf{p}_{xf}$ $\boldsymbol{q}_i = \text{row } \boldsymbol{i} \text{ of } \boldsymbol{Q}$ $p_x = \text{column } x \text{ of } P^T$

users

	_	_							_			
) S	1.1	2	.3	.5	-2	5	.8	4	.3	1.4	2.4	9
icto	8	.7	.5	1.4	.3	-1	1.4	2.9	7	1.2	1	1.3
fa	2.1	4	.6	1.7	2.4	.9	3	.4	.8	.7	6	.1
fa	2.1	4	.6	1.7	2.4	.9	3	.4	.8	.7	6	.1

PI



Ratings as Products of Factors

How to estimate the missing rating of





ors	1.1	2	.3	.5	-2	5	.8	4	.3	1.4	2.4	9
act	8	.7	.5	1.4	.3	-1	1.4	2.9	7	1.2	1	1.3
ff	2.1	4	.6	1.7	2.4	.9	3	.4	.8	.7	6	.1

DI

IISARS

 $\hat{\boldsymbol{r}}_{xi} = \boldsymbol{q}_i \cdot \boldsymbol{p}_x$

 $\boldsymbol{q}_i = \text{row } \boldsymbol{i} \text{ of } \boldsymbol{Q}$

 $p_x = \text{column } x \text{ of } P^T$

ρ_λ q_{if} · p_{xf}

U Kang



Latent Factor Models



Latent Factor Models



Singular Value Decomposition(SVD)

SVD:

- A: Input data matrix
- U: Left singular vecs
- V: Right singular vecs
- Σ : Singular values



So in our case: "SVD" on Netflix data: $R \approx Q \cdot P^T$ $A = R, Q = U, P^T = \sum V^T$



SVD: More good stuff

SVD gives minimum reconstruction error (Sum of S quared Errors):

$$\min_{U,V,\Sigma} \sum_{ij\in A} \left(A_{ij} - [U\Sigma V^{\mathrm{T}}]_{ij} \right)^2$$

Note two things:

- **SSE** and **RMSE** are monotonically related:
 - $RMSE = \frac{1}{c}\sqrt{SSE}$ Great news: SVD is minimizing RMSE
- Complication: The sum in SVD error term is over all entries (no-rating is interpreted as zero-rating).
 But our *R* has missing entries!



Latent Factor Models



SVD isn't defined when entries are missing!

Use specialized methods to find P, Q

$$= \min_{P,Q} \sum_{(i,x)\in\mathbb{R}} (r_{xi} - q_i \cdot p_x)^2 \qquad \hat{r}_{xi} = q_i \cdot p_x$$

- Note:
 - We don't require cols of P, Q to be orthogonal/unit length
 - P, Q map users/movies to a latent space
 - The most popular model among Netflix contestants U Kang



Outline

- Metflix Prize; Weight Learning in CF
- Latent Factor Model
- Þ 🔲 Regularization for LF
 - □ Bias Extension for LF
 - □ Netflix Challenge



Latent Factor Models

Our goal is to find P and Q such tat:

$$\min_{P,Q}\sum_{(i,x)\in R} (r_{xi}-q_i\cdot p_x)^2$$





Back to Our Problem

- Want to minimize SSE for unseen test data
- Idea: Minimize SSE on training data
 - Want large k (# of factors) to capture all the signals
 - But, **SSE** on <u>test</u> data begins to rise for k > 2
- This is a classical example of overfitting:
 - With too much freedom (too many free parameters) the model starts fitting noise
 - That is it fits too well the training data and thus not generalizing well to unseen test data





Dealing with Missing Entries

To solve overfitting we introduce regularization:



Allow rich model where there are sufficient data

Shrink aggressively where data are scarce



 $\lambda_1, \lambda_2 \dots$ user set regularization parameters (≥ 0)

Note: We do not care about the "raw" value of the objective function, but we care in P,Q that achieve the minimum of the objective

U Kang





















Outline

- Metflix Prize; Weight Learning in CF
- Latent Factor Model
- Regularization for LF
- Þ 🔲 Bias Extension for LF
 - Netflix Challenge

Modeling Biases and Interactions





Baseline predictor

- Separates users and movies
- Benefits from insights into user's behavior
- Among the main practical contributions of the competition

$$\mu$$
 = overall mean rating

b_x = bias of user
$$\mathbf{x}$$

b $_{i}$ = bias of movie i



User-Movie interaction

- Characterizes the matching between users and movies
- Attracts most research in the field
- Benefits from algorithmic and mathematical innovations



Baseline Predictor

We have expectations on the rating by user x of movie i, even without estimating x's attitude towards movies like i



- Rating scale of user x
- Values of other ratings user gave recently



(Recent) popularity of movie *i*



Putting It All Together



Example:

- Mean rating: $\mu = 3.7$
- You are a critical reviewer: your ratings are 1 star lower than the mean: b_x = -1
- Star Wars gets a mean rating of 0.5 higher than average movie: b_i = + 0.5
- Predicted rating for you on Star Wars (w/o interaction):
 = 3.7 1 + 0.5 = 3.2



Fitting the New Model

Solve:

$$\begin{split} \min_{Q,P} \sum_{(x,i)\in R} & \left(r_{xi} - (\mu + b_x + b_i + q_i p_x)\right)^2 \\ & \text{goodness of fit} \\ & + \left(\lambda_1 \sum_i \left\|q_i\right\|^2 + \lambda_2 \sum_x \left\|p_x\right\|^2 + \lambda_3 \sum_x \left\|b_x\right\|^2 + \lambda_4 \sum_i \left\|b_i\right\|^2\right) \\ & \wedge \text{ is selected via cross-} \end{split}$$

validation

Stochastic gradient decent to find parameters

Note: Both biases b_x, b_i as well as interactions q_i, p_x are t reated as parameters (we estimate them)

Performance of Various Methods

Global average: 1.1296 User average: 1.0651 Movie average: 1.0533 Netflix: 0.9514 Basic Collaborative filtering: 0.94 Collaborative filtering++: 0.91 Latent factors: 0.90 Latent factors+Biases: 0.89 Grand Prize: 0.8563



Temporal Biases Of Users

U Kang

- Sudden rise in the average movie rating (early 2004)
 - Improvements in Netflix
 - GUI improvements
 - Meaning of rating changed

Movie age

 Older movies receive higher ratings than newer ones

Y. Koren, Collaborative filtering with temporal dynamics, KDD '09





Temporal Biases & Factors

Original model:

$$r_{xi} = \mu + b_x + b_i + q_i \cdot p_x$$

Add time dependence to biases:

$$r_{xi} = \mu + b_x(t) + b_i(t) + q_i \cdot p_x$$

- Make parameters $\boldsymbol{b}_{\boldsymbol{x}}$ and $\boldsymbol{b}_{\boldsymbol{i}}$ to depend on time
- (1) Parameterize time-dependence by linear trends
 (2) Each bin corresponds to 10 consecutive weeks

$$b_i(t) = b_i + b_{i,\operatorname{Bin}(t)}$$

Add temporal dependence to factors

$\square p_x(t)$... user preference vector on day t

Y. Koren, Collaborative filtering with temporal dynamics, KDD '09 U Kang 46

Performance of Various Methods

Basic Collaborative filtering: 0.94 Collaborative filtering++: 0.91 Latent factors: 0.90

Latent factors+Biases: 0.89

Latent factors+Biases+Time: 0.876

Global average: 1.1296

<u>User</u> average: 1.0651 Movie average: 1.0533

Netflix: 0.9514

Still no prize! 🛞 Getting desperate.

Grand Prize: 0.8563



Outline

- Metflix Prize; Weight Learning in CF
- Latent Factor Model
- Regularization for LF
- Bias Extension for LF
- Þ 🔲 Netflix Challenge



Final Solution

- Many solutions proposed
 - Baseline
 - Basic collaborative filtering
 - Basic collaborative filtering w/ weight learning
 - Latent factor model
 - Latent factor w/ time bias

• ...

- 'Blending' the solutions leads to the best performance
 - □ Linear combination of N (≥ 500) predictors

$$\Box \ \widehat{r_{xi}} = \sum_{k=1}^{N} w_i \cdot pred_k(x, i) \qquad \text{pred}_k: \text{ k th predictor}$$

The big picture Solution of BellKor's Pragmatic Chaos



Michael Jahrer / Andreas Töscher – Team BigChaos – September 21, 2009



Standing on June 26th 2009

TFL	I X			
Ne	etflix Priz	ze		
Rul	es Leaderboard Register	Update Su	bmit Download	
Lea	aderboard		Display top	20 leaders.
Rank	Team Name	Best Score	% Improvement	Last Submit Time
1 ;	BellKor's Pragmatic Chaos	0.8558	10.05	2009-06-26 18:42:37
Grand	<u>d Prize</u> - RMSE <= 0.8563			
2	PragmaticTheory	0.8582	9.80	2009-06-25 22:15:51
3	BellKor in BigChaos	0.8590	9.71	2009-05-13 08:14:09
4 }	Grand Prize Team	0.8593	9.68	2009-06-12 08:20:24
5	Dace	0.8604	9.56	2009-04-22 05:57:03
5	BigChaos	0.8613	9.47	2009-06-23 23:06:52
Progr	ress Prize 2008 - RMSE = 0.8	616 - Winning Ti	eam: BellKor in BigC	haos
7	BellKor	0.8620	9.40	2009-06-24 07:16:02
B	Gravity	0.8634	9.25	2009-04-22 18:31:32
9	Opera Solutions	0.8638	9.21	2009-06-26 23:18:13
10	BruceDengDaoCiYiYou	0.8638	9.21	2009-06-27 00:55:55
11	pengpengzhou	0.8638	9.21	2009-06-27 01:06:43
12	xivector	0.8639	9.20	2009-06-26 13:49:04
13	xiangliang	0.8639	9.20	2009-06-26 07:47:34
201	and the second	No. Contraction of the second		

June 26th submission triggers 30-day "last call"

NETFLIX

Rules

Netflix Prize

Home

Leaderboard

Download

Update

Leaderboard

Showing Test Score. Click here to show quiz score

COMPLETED

Display top 20 \$ leaders.

Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
Gran	<u>d Prize</u> - RMSE = 0.8567 - Winning Te	am: BeliKor's Press	natic Chaos	_
1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:18:28
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22
3	Grand Prize Team	0.058		0210
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31
5	Vandelay Industries !	0.8591	9.81	2009-07-10 00:32:20
6	PragmaticTheory	0.8594	9.77	2009-06-24 12:06:56
7	BellKor in BigChaos	0.8601	9.70	2009-05-13 08:14:09
8	Dace	0.8612	9.59	2009-07-24 17:18:43
9	Feeds2	0.8622	9.48	2009-07-12 13:11:51
10	BigChaos	0.8623	9.47	2009-04-07 12:33:59
11	Opera Solutions	0.8623	9.47	2009-07-24 00:34:07
12	BellKor	0.8624	9.46	2009-07-26 17:19:11
Prog	<u>ress Prize 2008</u> - RMSE = 0.8627 - Wi	inning Team: BellKo	r in BigChaos	
13	xiangliang	0.8642	9.27	2009-07-15 14:53:22
14	Gravity	0.8643	9.26	2009-04-22 18:31:32
15	Ces	0.8651	9.18	2009-06-21 19:24:53
16	Invisible Ideas	0.8653	9.15	2009-07-15 15:53:04
17	Just a guy in a garage	0.8662	9.06	2009-05-24 10:02:54
18	<u>J Dennis Su</u>	0.8666	9.02	2009-03-07 17:16:17
19	Craig Carmichael	0.8666	9.02	2009-07-25 16:00:54
20	acmehill	0.8668	9.00	2009-03-21 16:20:50

Progress Prize 2007 - RMSE = 0.8723 - Winning Team: KorBell



Million \$ Awarded Sept 21st 2009

	dhan alla	2009	
	MELETIX	DATE 09.21.09	
7	CROSER OF BellKor's Pragmatic Chaos	s 1,000,000 ≌	
		00/100	
	FOR The Netflix Prize Reed F	Hastings	



Questions?