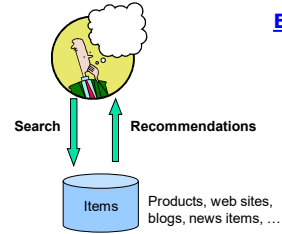


Data Mining Recommendation Models

Prof. Dr. Şule Gündüz Öğüdücü

1

Recommendations



Examples:

- amazon.com
- StumbleUpon
- delicio.us
- NETFLIX
- movielens
- lost.fm
- Google News
- YouTube
- XBOX LIVE

<http://www.mmds.org>

2

Recommender Systems

- Recommend new content

LIVE UK nurses hold biggest strike in NHS history

Nurses in England, Wales and Northern Ireland have begun the first of two day-long walkouts over pay.

Health

- Largest nursing strike in NHS history set to start



- What will nurses' strike mean for patients?



Australia police killers obsessed with guns - father

Ronald Train says his sons took a "dark track" after cutting off their family 20 years ago.

4h • Australia



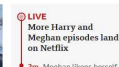
Musk taking legal action over private jet tweets



'My boss made me



Deadly winter storm barrels east across US



LIVE More Harry and Meghan episodes land on Netflix

2m Meghan (Items herself 'to foreign organist')

11m New Harry and

3

Recommender Systems

- Recommend similar items

good or bad??



Obi-Wan Kenobi



Andor



Star Wars: Episode V - The Empire Strikes...



Tales of the Jedi

4

Recommender Systems

- Recommend co-occurred items

Sıklıkla birlikte alınır



Motiv Kadın Polo Yaka T-shirti

189,99 TL



Rod 102215200005 Rodin

189,00 TL



Virgin Kargın Yastır Yemegün Üzeri T-shirti Herkes

899,99 TL



Zincir Kızıl Kızıl Göl

347,19 TL



Lale Baskın Yaka Yemegün

114,99 TL



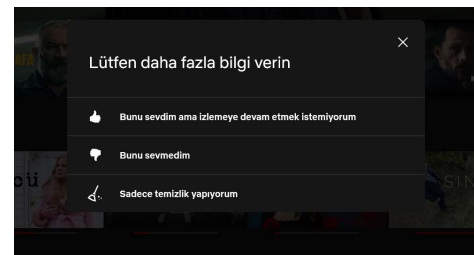
Lale Baskın Yaka Yemegün

289,99 TL

5

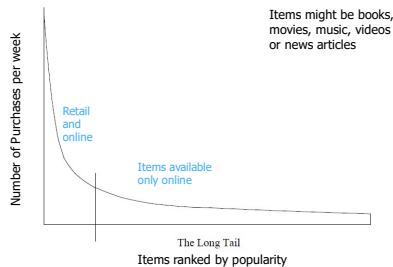
Recommender Systems

- Identify items that we may like



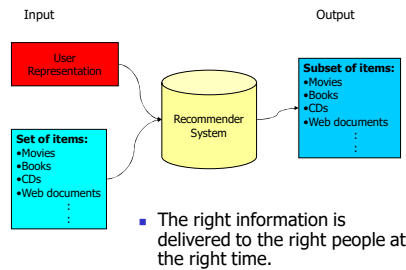
6

The Long Tail



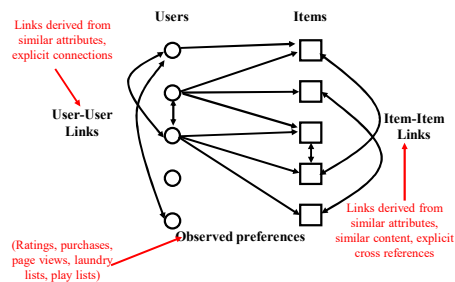
7

Recommendation Process



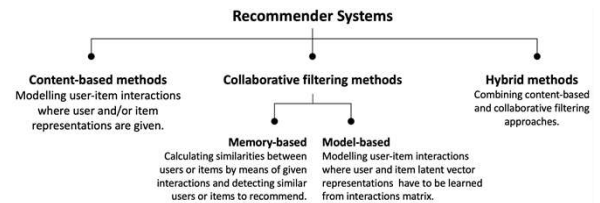
8

The Recommendation System Space



9

Recommender Systems: Methods



How it Works?

- Each user has a **profile**
- Users **rate** items
 - Explicitly: score from 1..5
 - Implicitly: web usage mining
 - Time spent in viewing the item
 - Navigation path
 - Etc...
- System does the rest, How?
 - This is what we will show today

Formal Model

- X = set of **Customers**
- S = set of **Items**
- Utility function** $u: X \times S \rightarrow R$
 - R = set of ratings
 - R is a totally ordered set
 - e.g., 0-5 stars, real number in $[0,1]$

http://www.mmds.org

Utility Matrix

	Avatar	LOTR	Matrix	Pirates
Alice	1		0.2	
Bob		0.5		0.3
Carol	0.2		1	
David				0.4

<http://www.mmds.org>

Key Problems

- (1) Gathering "known" ratings for matrix
 - How to collect the data in the utility matrix
- (2) Extrapolate unknown ratings from the known ones
 - Mainly interested in high unknown ratings
 - We are not interested in knowing what you don't like but what you like
- (3) Evaluating extrapolation methods
 - How to measure success/performance of recommendation methods

(1) Gathering Ratings

- Explicit**
 - Ask people to rate items
 - Doesn't work well in practice – people can't be bothered
- Implicit**
 - Learn ratings from user actions
 - E.g., purchase implies high rating
 - What about low ratings?

<http://www.mmds.org>

(2) Extrapolating Utilities

- Key problem:** Utility matrix U is **sparse**
 - Most people have not rated most items
- Cold start:**
 - New items have no ratings
 - New users have no history
- Three approaches to recommender systems:**
 - 1) Content-based
 - 2) Collaborative
 - 3) Hybrid Models

Collaborative Filtering

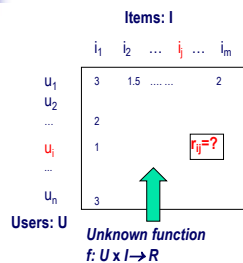
- Method**
 - Correlation between user's interests (such as votes and trails)
 - Results are captured in a generally sparse matrix (users x items)
 - Similarities between items
- Problems**
 - Sparsity
 - Cold-start
 - Diversity

Data Structures

- Each user has a profile**
- Users rate items**
 - Explicitly: a user consciously express his or her preference for a title
 - score from 1..5
 - Implicitly: interpret user behavior or selections to impute a vote or preference using web usage mining
 - Time spent in viewing the item
 - Navigation path
 - purchase history
 - Etc...

		Items: I					
		i_1	i_2	...	i_j	...	i_m
Users: U	U_1	3	1.5	...			2
	U_2						
	...	2					
	U_i	1			r_{ij}		
	...						
	U_n	3					

Collaborative Filtering: A Framework



The task:
Q1: Find Unknown ratings?
Q2: Which items should we recommend to this user?

Collaborative Filtering

- **User-Centric**
 - The preferences of a large group of people are registered. These preferences could be items bought by the user on a e-commerce Web site, or Web pages visited by the user
 - According to a similarity metric a subgroup of people are selected whose preferences are similar to the current user's preferences.
 - An average of the preferences of that subgroup is calculated.
 - Options on which the current user has no experience yet are selected to generate a recommendation set.
- **Item Centric**
 - Build an item-item matrix
 - Determine the relationship between pairs of items

Algorithms for CF

- r_{ij} = vote of user i on item j
- I_i = items for which user i has voted
- Mean vote for user i is

$$\bar{r}_i = \frac{1}{|I_i|} \sum_{j \in I_i} r_{ij}$$

- Predicted vote for "active user" a is weighted sum

$$r_{aj} = \bar{r}_a + \kappa \sum_{i=1}^n \underbrace{w(a,i)}_{\text{similarity between user } i \text{ and active user } a} (r_{ij} - \bar{r}_i)$$

normalizer

How to Measure Similarity?

- **k-nearest neighbour:**
 - Compute distance between all other users and active user
- aggregate ratings from K nearest neighbors to predict active user's rating

$$d_{ai} = \sum_{j=1}^M (r_{aj} - r_{ij})^2$$

$$w_k(a,i) = \begin{cases} 1 & \text{if } i \in \text{neighbors}(a) \\ 0 & \text{else} \end{cases}$$

Finding "Similar" Users

- Let r_a be the vector of user a 's ratings
- **Jaccard similarity measure**
 - **Problem:** Ignores the value of the rating
- **Cosine similarity measure**
 - $w_c(a,i) = \frac{r_a \cdot r_i}{||r_a|| ||r_i||}$
 - **Problem:** Treats missing ratings as "negative"
- **Pearson correlation coefficient**
 - $S_{ai} = I_a \cap I_i$ = items rated by both users a and i

$$r_a = [*, -, *, *, **] \\ r_i = [*, -, **, *, -]$$

$$r_a, r_i \text{ as sets:} \\ r_a = \{1, 4, 5\} \\ r_i = \{1, 3, 4\}$$

$$r_a, r_i \text{ as points:} \\ r_a = \{1, 0, 0, 1, 3\} \\ r_i = \{1, 0, 2, 2, 0\}$$

$$w_p(a,i) = \frac{\sum_{j \in S_{ai}} (r_{aj} - \bar{r}_a)(r_{ij} - \bar{r}_i)}{\sqrt{\sum_{j \in S_{ai}} (r_{aj} - \bar{r}_a)^2} \sqrt{\sum_{j \in S_{ai}} (r_{ij} - \bar{r}_i)^2}}$$



<http://www.mmds.org>

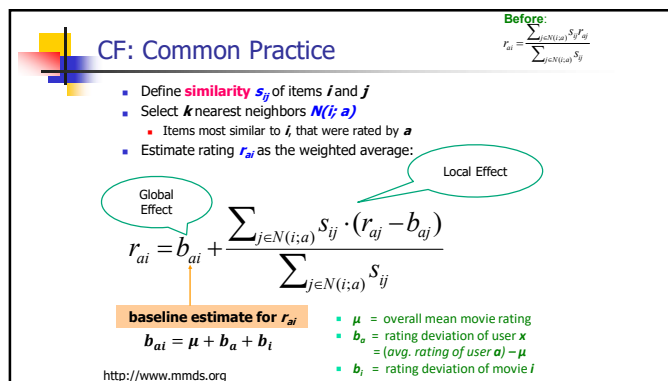
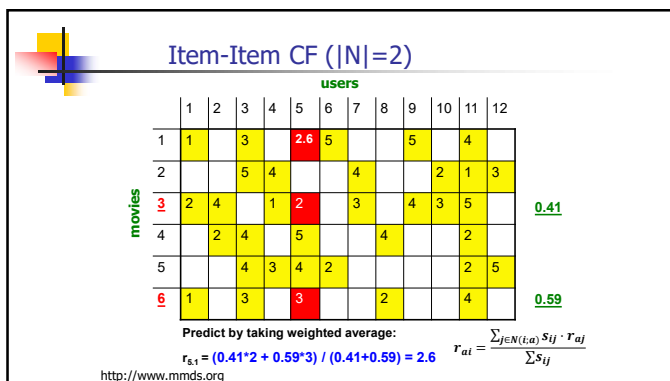
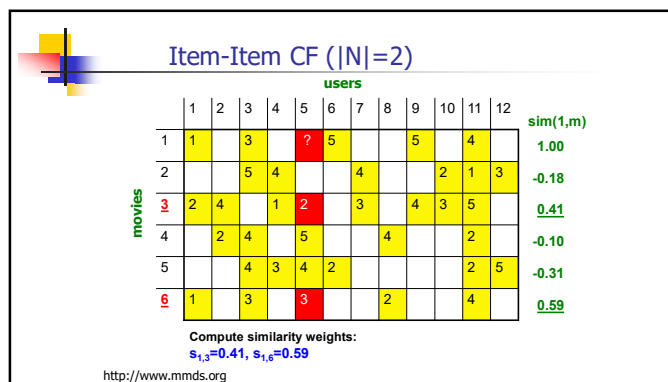
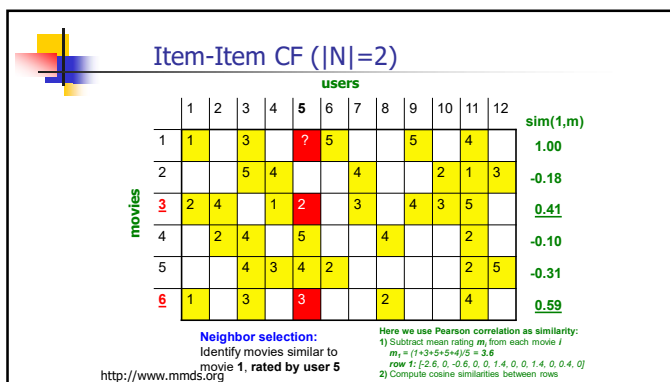
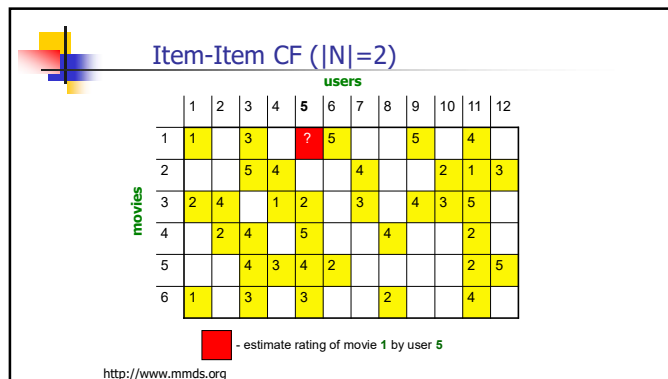
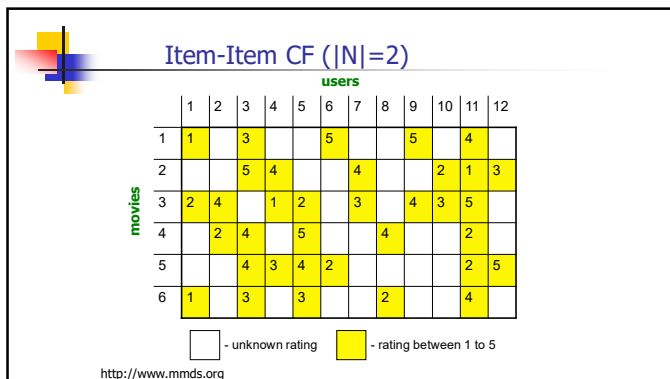
Item-Item Collaborative Filtering

- **So far: User-user collaborative filtering**
- **Another view: Item-item**
 - For item i , find other similar items
 - Estimate rating for item i based on ratings for similar items
 - Can use same similarity metrics and prediction functions as in user-user model

$$r_{ai} = \frac{\sum_{j \in N(i;a)} S_{ij} \cdot r_{aj}}{\sum_{j \in N(i;a)} S_{ij}}$$

S_{ij} ... similarity of items i and j
 r_{aj} ... rating of user a on item j
 $N(i;a)$... set items rated by a similar to i

<http://www.mmds.org>



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*

⇒ **Final estimate:**

Joe will rate *The Sixth Sense* 3.8 stars



$$b_{ai} = \mu + b_a + b_i = 3.7 + (-0.2) + 0.5 = 4$$

- μ = overall mean movie rating
- b_a = rating deviation of user x
= (avg. rating of user a) - μ
- b_i = rating deviation of movie i



<http://www.mmds.org>

Pros/Cons of Collaborative Filtering

+ Works for any kind of item

- No feature selection needed

- Cold Start:

- Need enough users in the system to find a match

- Sparsity:

- The user/ratings matrix is sparse
- Hard to find users that have rated the same items

- First rater:

- Cannot recommend an item that has not been previously rated
- New items, Esoteric items

- Popularity bias:

- Cannot recommend items to someone with unique taste
- Tends to recommend popular items

- Similarity:

- Similarity measures are "arbitrary": Pairwise similarities neglect interdependencies among users

$$r_{ai} = b_{ai} + \frac{\sum_{j \in N(i,a)} s_{ij} \cdot (r_{aj} - b_{aj})}{\sum_{j \in N(i,a)} s_{ij}}$$

<http://www.mmds.org>

Idea: Interpolation Weights w_{ij}

- Use a **weighted sum** rather than **weighted avg.**:

$$\hat{r}_{xi} = b_{xi} + \sum_{j \in N(i,x)} w_{ij} (r_{xj} - b_{xi})$$

A few notes:

- $N(i; x)$... set of movies rated by user x that are similar to movie i
- 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)

<http://www.mmds.org>

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_{ij} that minimize **SSE on training data!**

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

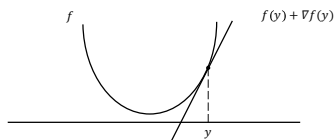
- Think of w as a vector of numbers

<http://www.mmds.org>

Detour: Minimizing a function

A simple way to minimize a function $f(x)$:

- Compute the take a derivative ∇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**



<http://www.mmds.org>

Interpolation Weights

- We have the optimization problem, now what?**

$$J(w) = \sum_x \left(\left(b_{xi} + \sum_{j \in N(i,x)} w_{ij} (r_{xj} - b_{xi}) \right) - r_{xi} \right)^2$$

Gradient decent:

- Iterate until convergence: $w \leftarrow w - \eta \nabla_w J$**

- where $\nabla_w J$ is the gradient (derivative evaluated on data):**

η ... learning rate

$$\nabla_w J = \left[\frac{\partial J(w)}{\partial w_{ij}} \right]$$

$$= 2 \sum_{x,i} \left(\left(b_{xi} + \sum_{k \in N(i,x)} w_{ik} (r_{xk} - b_{xi}) \right) - r_{xi} \right) (r_{xj} - b_{xi})$$

for $j \in \{N(i; x), \forall i, \forall x\}$

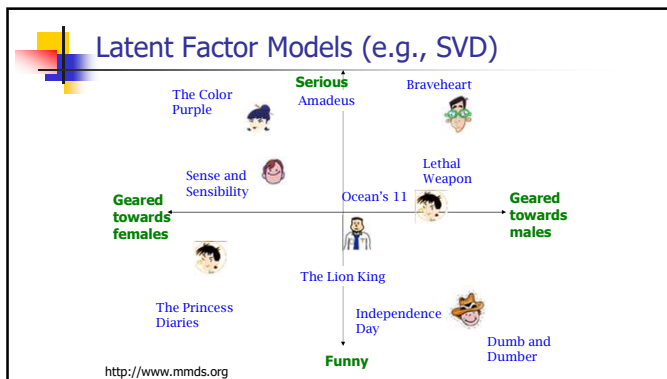
else $\frac{\partial J(w)}{\partial w_{ij}} = 0$

- 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 $|w_{new} - w_{old}| > \epsilon$:

$w_{old} = w_{new}$

$w_{new} = w_{old} - \eta \cdot \nabla_w J$



Recap: SVD

Remember SVD:

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

So in our case:
"SVD" on Netflix data: $R \approx Q \cdot P^T$

$A = R, Q = U, P^T = \Sigma V^T$

We already know that SVD gives minimum reconstruction error (Sum of Squared Errors):

$$\min_{U, V} \sum_{i,j \in A} (A_{ij} - [UV^T]_{ij})^2$$

Note two things:

- SSE and RMSE are monotonically related:
 - $RMSE \propto \sqrt{SSE}$ **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!

<http://www.mmms.org>

Latent Factor Models

"SVD" on Netflix data: $R \approx Q \cdot P^T$ SVD: $A = U \Sigma V^T$

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

- 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

<http://www.mmms.org> 39

Ratings as Products of Factors

How to estimate the missing rating of user x for item i ?

$$\hat{r}_{xi} = q_i \cdot p_x$$

$$= \sum_f q_{if} \cdot p_{xf}$$

q_i = row i of Q
 p_x = column x of P^T

<http://www.mmms.org> 40

Ratings as Products of Factors

How to estimate the missing rating of user x for item i ?

$$\hat{r}_{xi} = q_i \cdot p_x$$

$$= \sum_f q_{if} \cdot p_{xf}$$

q_i = row i of Q
 p_x = column x of P^T

<http://www.mmms.org> 41

Ratings as Products of Factors

How to estimate the missing rating of user x for item i ?

$$\hat{r}_{xi} = q_i \cdot p_x$$

$$= \sum_f q_{if} \cdot p_{xf}$$

q_i = row i of Q
 p_x = column x of P^T

<http://www.mmms.org> 42

Latent Factor Models

users					factors				
1	3	5	4	2	1	4	2		
2	5	4	2	1	1.5	6	5		
3	1	2	3	5	1.1	2.1	3		
4	3	4	2	1	1.1	2.1	3		
5	3	4	2	1	1.1	2.1	3		

- SVD isn't defined when entries are missing!
- Use specialized methods to find P, Q

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

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

$$\hat{r}_{xi} = q_i \cdot p_x$$

<http://www.mmms.org>

Latent Factor Models

- Our goal is to find P and Q such that:

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

users					factors				
1	3	5	4	2	1	4	2		
2	5	4	2	1	1.5	6	5		
3	1	2	3	5	1.1	2.1	3		
4	3	4	2	1	1.1	2.1	3		
5	3	4	2	1	1.1	2.1	3		

<http://www.mmms.org>

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 test 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



<http://www.mmms.org>

Dealing with Missing Entries

- To solve overfitting we introduce **regularization**:

- Allow rich model where there are sufficient data
- Shrink aggressively where data are scarce



$$\min_{P, Q} \sum_{\text{training}} (r_{xi} - q_i p_x)^2 + \left[\lambda_1 \sum_x \|p_x\|^2 + \lambda_2 \sum_i \|q_i\|^2 \right]$$

"error" "length"

$\lambda_1, \lambda_2 \dots$ user set regularization parameters

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

<http://www.mmms.org>

Stochastic Gradient Descent

- Want to find matrices P and Q

$$\min_{P, Q} \sum_{\text{training}} (r_{xi} - q_i p_x)^2 + \left[\lambda_1 \sum_x \|p_x\|^2 + \lambda_2 \sum_i \|q_i\|^2 \right]$$

- Gradient descent:**

- Initialize P and Q (using SVD, pretend missing ratings are 0)

- Do gradient descent:

$$P \leftarrow P - \eta \nabla P$$

$$Q \leftarrow Q - \eta \nabla Q$$

where ∇Q is gradient/derivative of matrix Q

$$\nabla Q = [\nabla q_{ij}] \text{ and } \nabla q_{ij} = \sum_x (r_{xi} - q_i p_x) p_{xj} + 2\lambda_2 q_{ij}$$

Here q_{ij} is entry f of row q_i of matrix Q

How to compute gradient of a matrix?
Compute gradient of every element independently!

<http://www.mmms.org>

Stochastic Gradient Descent

- Gradient Descent (GD) vs. Stochastic GD**

- Observation:** $\nabla Q = [\nabla q_{ij}]$ where

$$\nabla q_{ij} = \sum_x -2(r_{xi} - q_i p_x) p_{xj} + 2\lambda_2 q_{ij} = \sum_x \nabla Q(r_{xi})$$

- Here q_{ij} is entry f of row q_i of matrix Q

$$Q = Q - \eta \nabla Q = Q - \eta [\sum_x \nabla Q(r_{xi})]$$

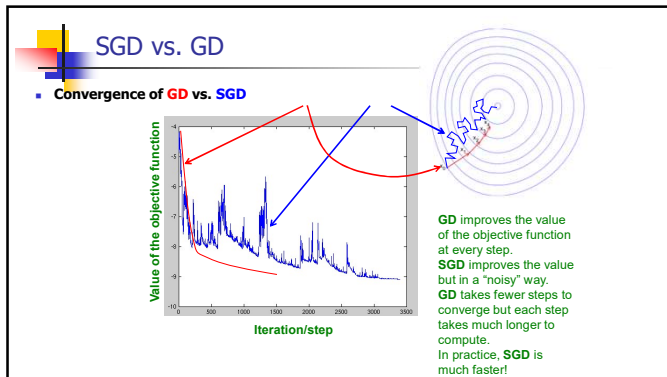
- Idea:** Instead of evaluating gradient over all ratings evaluate it for each individual rating and make a step

$$\text{GD: } Q \leftarrow Q - \eta [\sum_x \nabla Q(r_{xi})]$$

$$\text{SGD: } Q \leftarrow Q - \eta \nabla Q(r_{xi})$$

- Faster convergence!**

- Need more steps but each step is computed much faster



Stochastic Gradient Descent

- Stochastic gradient descent:**
 - Initialize \mathbf{P} and \mathbf{Q} (using SVD, pretend missing ratings are 0)
 - Then iterate over the ratings (multiple times if necessary) and update factors:
 - For each r_{xi} :
 - $\epsilon_{xi} = 2(r_{xi} - q_i \cdot p_x)$ (derivative of the "error")
 - $q_i \leftarrow q_i + \eta_1 (\epsilon_{xi} p_x - \lambda_2 q_i)$ (update equation)
 - $p_x \leftarrow p_x + \eta_2 (\epsilon_{xi} q_i - \lambda_3 p_x)$ (update equation)
 - 2 for loops:
 - For until convergence:
 - For each r_{xi}
 - Compute gradient, do a "step" $\eta \dots$ learning rate

Putting It All Together

$$r_{xi} = \underbrace{\mu}_{\text{Overall mean rating}} + \underbrace{b_x}_{\text{Bias for user } x} + \underbrace{b_i}_{\text{Bias for movie } i} + \underbrace{q_i \cdot p_x}_{\text{User-Movie interaction}}$$

- 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: $= 3.7 - 1 + 0.5 = 3.2$

Fitting the New Model

- Solve:**

$$\min_{\mathbf{Q}, \mathbf{P}} \sum_{(x,i) \in R} (r_{xi} - (\mu + b_x + b_i + q_i \cdot p_x))^2 + \left(\lambda_1 \sum_i \|q_i\|^2 + \lambda_2 \sum_x \|p_x\|^2 + \lambda_3 \sum_x \|b_x\|^2 + \lambda_4 \sum_i \|b_i\|^2 \right)$$

λ is selected via grid-search on a validation set

- Stochastic gradient descent to find parameters**
 - Note:** Both biases b_x, b_i as well as interactions q_i, p_x are treated as parameters (we estimate them)

Evaluation Metrics for Recommendation Systems

- Recall@K** = $\frac{\text{Number of Relevant Items in Top } K}{\text{Total Number of Relevant Items}}$
- Precision@K** = $\frac{\text{Number of Relevant Items in Top } K}{K}$
- F1@K** = $\frac{2 \times \text{Precision@K} \times \text{Recall@K}}{\text{Precision@K} + \text{Recall@K}}$
- Mean Average Precision**
 - A movie recommender system: Recommend K movies, number of relevant items Q
 - Recommendation List: {Star Wars-Return of Jedi, Back To the Future, The Matrix}
 - Ground Truth: {Terminator 2, Back To the Future, The Matrix}
 - $\text{Precision@K} = [0, 1/2, 2/3]$
 - $AP = (1/3) \times [(1/2) + (2/3)] = 0.38$

A movie recommender system: Recommend 10 movies for every user. A user has seen 5 movies, the recommendation list has 3 of them.

$$\text{Recall@K} = \frac{3}{5} = 0.6$$


$$\text{Precision@K} = \frac{3}{10} = 0.3$$

Evaluation Metrics for Recommendation Systems

- Mean Reciprocal Rank:** the position of the first relevant item in the recommendation list
- $$\text{MRR} = \frac{1}{N} \sum_{i=1}^N \frac{1}{\text{rank}_i}$$
 - A movie recommender system: Recommend K movies, number of relevant items N
 - Recommendation List: {Star Wars-Return of Jedi, Back To the Future, The Matrix}
 - Ground Truth: {Back To the Future, The Matrix}
 - $\text{MRR} = 1/2 \times [(1/2) + (1/3)] = 0.41$
- Normalized Cumulative Discounted Gain (NDCG):** a measure of how good a ranked list is
- $$\text{NDCG@k} = \frac{\text{DCG@k}}{\text{IDCG@k}}$$
 - $\text{rel}(i)$: relevancy score of item i
 - IDCG@k : DCG@k of the "ideal" recommendation algorithm
 - Ground Truth: {Terminator 2, Back To the Future, The Matrix}
 - $\text{rel}(\text{Terminator 2}) = 1$, $\text{rel}(\text{Back To the Future}) = 1$, $\text{rel}(\text{The Matrix}) = 1$
 - List of Recommendations: {Star Wars-Return of Jedi, Back To the Future, The Matrix}
- $$\text{DCG@3} = \frac{0}{\log_2(1+1)} + \frac{1}{\log_2(2+1)} + \frac{1}{\log_2(3+1)}$$

$$\text{IDCG@3} = \frac{1}{\log_2(1+1)} + \frac{1}{\log_2(2+1)} + \frac{1}{\log_2(3+1)}$$

$$\text{NDCG@3} = \frac{\text{DCG@3}}{\text{IDCG@3}}$$



References

- Lecture slides of Mining Massive Datasets

55