



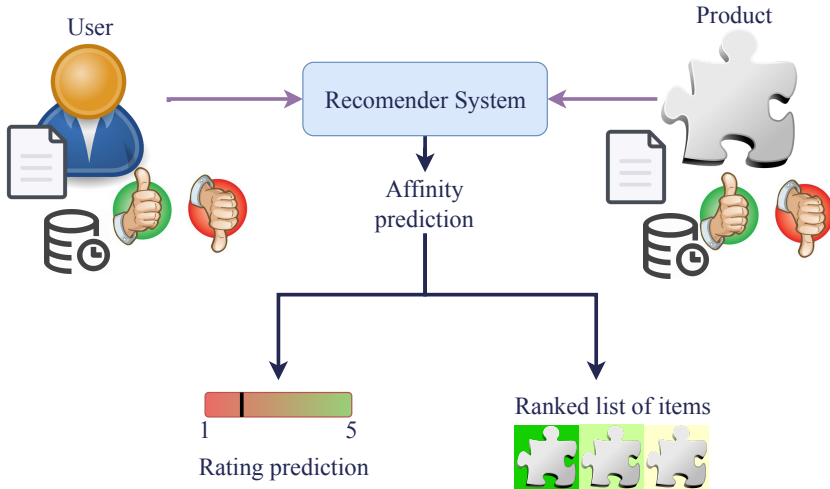
RECOMMENDER SYSTEMS : ADVANCED ARCHITECTURES NEW METRICS

Vincent Guigue



Introduction

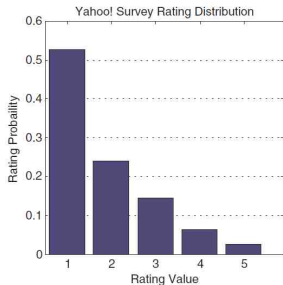
New issues



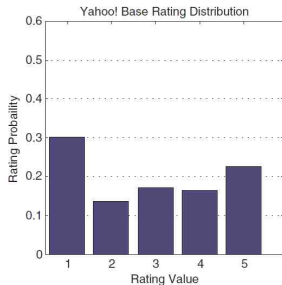
Main issue : the weakness of MCAR hypothesis

Data are not Missing Completely At Random...

Graphs from [Marlin & Zemel '09]:



Survey: ask users to rate a random list of items: approximates **complete** data



Typical Data: users are free to choose which items to rate -> available data are **MNAR** :
instead of giving low ratings, users
tend to not give a rating at all.

Even different in
product/movie domains:

- 60-80% of 4/5 ratings

Main issue : the weakness of MCAR hypothesis

Data are not Missing Completely At Random...

Predicting profile behavior on this kind of data:



H. Steck, KDD, 2010
Training and Testing of Recommender Systems on Data
Missing Not at Random

Table 1: Simplistic Example for ratings missing not at random (MNAR): test data where users rated only what they liked or knew.

		users				
		horror fans		romance lovers		
movie	h	5		5	5	
	o	5	5			
	r		5		5	
	.	5		5	5	
	r					5
eom	o				5	5
	.					5
	m				5	5
	.				5	5

Credit: H. Steck

Main issue : the weakness of MCAR hypothesis

Data are not Missing Completely At Random...

- Changing the **error function**
 - ranking criteria
- Changing **the task**
 - predicting rated item (not the rate)

Several outcomes:

New approaches

Implicit Feedback (SVD++)

For example, a dataset shows that users that rate “Lord of the Rings 3” high also gave high ratings to “Lord of the Rings 1–2”.

⇒ establish high weights from “Lord of the Rings 1–2” to “Lord of the Rings 3”.

Now, if a user did not rate “Lord of the Rings 1–2” at all, his predicted rating for “Lord of the Rings 3” will be penalized.

- Binary coding : **interesting** / **not interesting**
 - positive rating / simple visit = 1
 - negative rating / missing values = 0
 - Initial predictor f
 - Simple heuristic / expert knowledge...
 - $R(u)$: set of items rated by $u + f$
- ⇒ $N(u)$: set of items implicitly rated by u



Yehuda Koren, KDD 2008

Factorization Meets the Neighborhood: a Multifaceted Collaborative Filtering Model

SVD++

General idea:

$$\hat{r}_{ui} = \underbrace{b + b_u + b_i}_{b_{ui}} + \frac{1}{\sqrt{|R(u)|}} \sum_{j \in R(u)} (r_{uj} - b_{uj}) w_{ij} + \frac{1}{\sqrt{|N(u)|}} \sum_{j \in N(u)} c_{ij}$$

Overweighting deviation for prolific users ($\frac{1}{\sqrt{|R(u)|}}$ instead of $\frac{1}{|R(u)|}$)

w_{ij} Learning deviation meaning wrt b_{uj}

c_{ij} Learning the meaning of j absence wrt i

- ... Too expensive (too many coefficients to learn)

Factorized formulation = SVD++:

$$\hat{r}_{ui} = b_{ui} + \mathbf{i} \cdot \left(\mathbf{u} + \frac{1}{\sqrt{|N(u)|}} \sum_{j \in N(u)} \mathbf{y}_j \right)$$

⇒ Y. Koren winning proposal to the Netflix challenge (2009)

Factorization machine

Back to (huge) linear model !

- Factor = interactions between items
- Easy to encode unseen effect

Feature vector \mathbf{x}																	Target y					
$\mathbf{x}^{(1)}$	1	0	0	...	1	0	0	0	...	0.3	0.3	0.3	0	...	13	0	0	0	0	...	5	$y^{(1)}$
$\mathbf{x}^{(2)}$	1	0	0	...	0	1	0	0	...	0.3	0.3	0.3	0	...	14	1	0	0	0	...	3	$y^{(2)}$
$\mathbf{x}^{(3)}$	1	0	0	...	0	0	1	0	...	0.3	0.3	0.3	0	...	16	0	1	0	0	...	1	$y^{(2)}$
$\mathbf{x}^{(4)}$	0	1	0	...	0	0	1	0	...	0	0	0.5	0.5	...	5	0	0	0	0	...	4	$y^{(3)}$
$\mathbf{x}^{(5)}$	0	1	0	...	0	0	0	1	...	0	0	0.5	0.5	...	8	0	0	1	0	...	5	$y^{(4)}$
$\mathbf{x}^{(6)}$	0	0	1	...	1	0	0	0	...	0.5	0	0.5	0	...	9	0	0	0	0	...	1	$y^{(5)}$
$\mathbf{x}^{(7)}$	0	0	1	...	0	0	1	0	...	0.5	0	0.5	0	...	12	1	0	0	0	...	5	$y^{(6)}$
	A	B	C	...	TI	NH	SW	ST	...	TI	NH	SW	ST	...	Time	TI	NH	SW	ST	...		
	User				Movie					Other Movies rated						Last Movie rated						

$U = \{\text{Alice (A), Bob (B), Charlie (C), ...}\}$

$I = \{\text{Titanic (TI), Notting Hill (NH), Star Wars (SW),
Star Trek (ST), ...}\}$

$S = \{(A, TI, 2010-1, 5), (A, NH, 2010-2, 3), (A, SW, 2010-4, 1),$
 $(B, SW, 2009-5, 4), (B, ST, 2009-8, 5),$
 $(C, TI, 2009-9, 1), (C, SW, 2009-12, 5)\}$



S. Rendle, ICDM 2010
Factorization machines

Bayesian Personalized Ranking

For each user, what are the preferred items?

$$p(i >_u j | \theta) = \frac{1}{1 + \exp(-f_\theta(u, i, j))}$$

For instance (inspired from MF):

$$f_\theta(u, i, j) = \mathbf{u} \cdot \mathbf{i} - \mathbf{u} \cdot \mathbf{j}$$

NB: same nb of parameters than MF

Evaluation =

$$AUC = \frac{1}{n_u} \sum_u \frac{1}{|E(u)|} \sum_{(i,j) \in E(u)} \delta(\mathbf{u} \cdot \mathbf{i} > \mathbf{u} \cdot \mathbf{j})$$

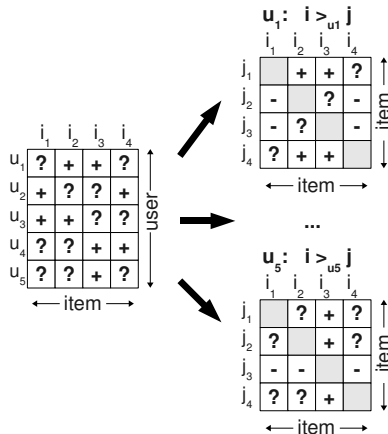


S. Rendle et al., UAI 2009

BPR: Bayesian Personalized Ranking from Implicit Feedback

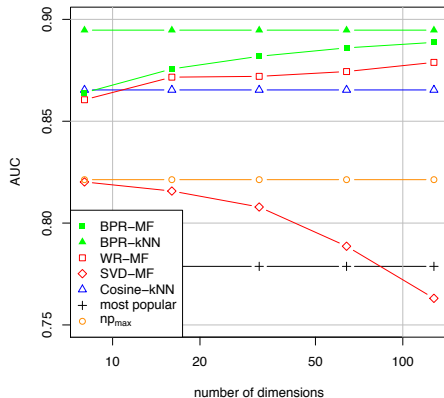
u_1 prefers $(i_2 \text{ to } i_1)$ & $(i_3 \text{ to } i_1)$

$\Rightarrow +$

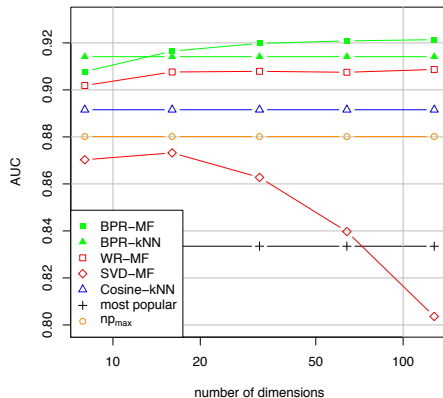


Bayesian Personalized Ranking

Online shopping: Rossmann



Video Rental: Netflix



The ranking criterion enables us to exploit high dimensional user/item representation



S. Rendle et al., UAI 2009

BPR: Bayesian Personalized Ranking from Implicit Feedback

AllRank

- $r_{ui}^{(o\&i)} \Rightarrow$ All binary rating

- 0 for missing
- 1 relevant

- $W_{ui} = \begin{cases} 1 & \text{if observed} \\ w_m & \text{if missing} \end{cases}$

- $b_{ui}^{(o\&i)} = \begin{cases} b_{ui} & \text{if observed} \\ r_m & \text{if missing} \end{cases}$

- w_m, r_m (& λ) have to be tuned

Cost function:

$$\mathcal{L} = \sum_u \sum_i W_{ui} \left[\left(r_{ui}^{(o\&i)} - (b_{ui}^{(o\&i)} + \mathbf{u} \cdot \mathbf{i}) \right)^2 + \lambda (\|\mathbf{u}\|^2 + \|\mathbf{i}\|^2) \right]$$



H. Steck, RecSys 2010

Training and Testing of Recommender Systems on Data Missing Not at Random

Learning to rank

- Pointwise :
 - Ranking score based on regression (or classification)
 - RS seen example: SVD approaches
- Pairwise :
 - Loss function is defined on pair-wise preferences
 - RankSVM, RankBoost, RankNet, FRank...
 - RS seen example: BPR
- Listwise :
 - Gradient descent on smoothed version of objective function (e.g. CLiMF presented at Recsys 2012 or TFMAP at SIGIR 2012)
 - SVM-MAP relaxes the MAP metric by adding it to the SVM constraints
 - AdaRank (modified version of Adaboost)

Evaluation:
evaluation metrics
vs
learning metrics

How to evaluate RS performance?

Warning

We should not confuse **evaluation metrics** & **learning metrics**

⇒ MSE is a convenient learning metrics

(easily differentiable + convex ...)

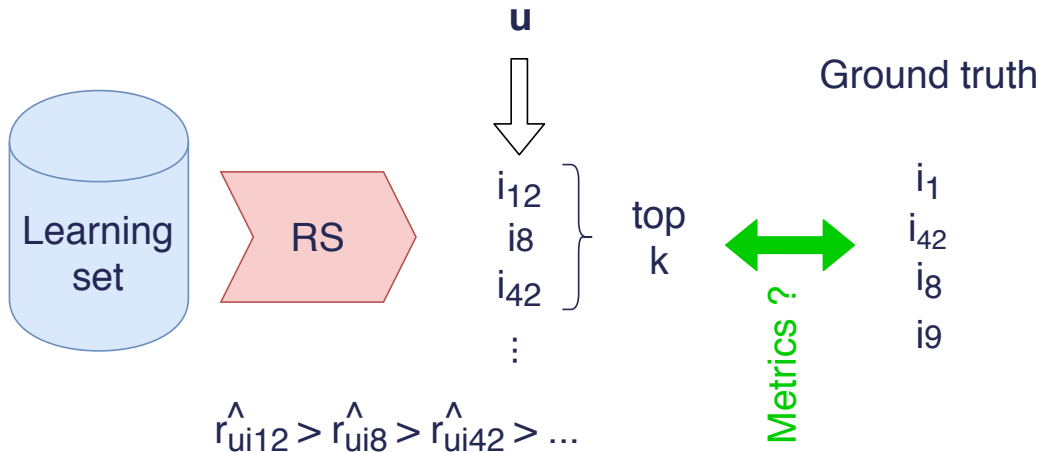
... but it is a poor evaluation metrics

... cf Netflix Challenge feedbacks

It do not tell us if we provide relevant suggestions

- What are the other available metrics?
 - Have a look towards the IR community
- Can we use those metrics during the learning step?

Precision / Recall



1/0 labeling, AUC metrics

- Rendle popularize both 1/0 prediction & AUC metrics
- AUC = tradeoff between precision & recall
 - Percentage of correct binary ranking for ONE user
 - Aggregation over n_u users

$$AUC = \frac{1}{n_u} \sum_u \frac{1}{|E(u)|} \sum_{(i,j) \in E(u)} \delta(\mathbf{u} \cdot \mathbf{i} > \mathbf{u} \cdot \mathbf{j})$$

- + k not required
- **top** of the list = same impact as **bottom** of the list

Mean Average Precision (from the IR domain)

RS aim at proposing an ordered list of suggestion...

Which **head** is far more important than the rest.

For a user u with 4 liked items to discover:

$$query = u \Rightarrow RS_1 \Rightarrow \left[\begin{array}{c} i_{12} \\ i_8 \\ i_{42} \\ i_1 \end{array} \right] \Leftrightarrow \left[\begin{array}{c} i_1 \\ i_{42} \\ i_8 \\ i_9 \end{array} \right] = GT$$

- Average precision (one query/user) :

$$\frac{1}{K} \sum_{k=1}^K precision@K = \frac{1}{4} \left(0 + \frac{1}{2} + \frac{2}{3} + \frac{3}{4} \right) = 0.478$$

- Mean Average Precision =

Averaging over the whole population 13/20

Mean Average Precision (from the IR domain)

RS aim at proposing an ordered list of suggestion...

Which **head** is far more important than the rest.

For a user u with 4 liked items to discover:

$$query = u \Rightarrow RS_2 \Rightarrow \left[\begin{array}{c} i_1 \\ i_8 \\ i_{42} \\ \downarrow \\ i_{12} \end{array} \right] \Leftrightarrow \left[\begin{array}{c} i_1 \\ i_{42} \\ i_8 \\ i_9 \end{array} \right] = GT$$

■ Average precision :

$$\frac{1}{4} \sum_{k=1}^4 precision@K = \frac{1}{4} (1 + 1 + 1 + \frac{3}{4}) = 0.9375$$

■ Mean Average Precision =

Averaging over the whole population 13/20

Mean Reciprocal Rank

At which rank is the first relevant item?

$$query = \mathbf{u} \Rightarrow RS \Rightarrow \left[\begin{array}{c} i_{12} \\ i_8 \\ i_{42} \\ i_1 \end{array} \right] \Leftrightarrow \left[\begin{array}{c} i_1 \\ i_{42} \\ i_8 \\ i_9 \end{array} \right] = GT$$

$$RR = \frac{1}{rank_i} = \frac{1}{2} \text{ on previous example}$$

Mean Reciprocal Rank =

Averaging over the whole population

$\Rightarrow \approx$ How many iterations to obtain a relevant item?

nDCG : Normalized Discounted Cumulative Gain

We assume that we have a relevance score for each item...

$$query = \mathbf{u} \Rightarrow RS \Rightarrow \left[\begin{array}{c|c} i_{12} & ind = 1 \\ i_8 & ind = 2 \\ i_{42} & ind = 3 \\ \vdots & \\ i_1 & ind = 4 \end{array} \right] \Leftrightarrow \left[\begin{array}{c} 0 \\ 2 \\ 3 \\ 3 \end{array} \right] = relevance$$

$$DCG_p = \sum_{ind=1}^p \frac{relev_{ind}}{\log_2(ind+1)} = 0 + 1.26 + 1.5 + 1.29 = 4.05$$

$$nDCG = \frac{DCG}{IdealDCG} = \frac{4.05}{3 + 1.89 + 1 + 0/0.86} = 0.69/0.6$$

Relative ideal (among suggestions) vs Absolute ideal (among all items)

ATOP

recall@k =

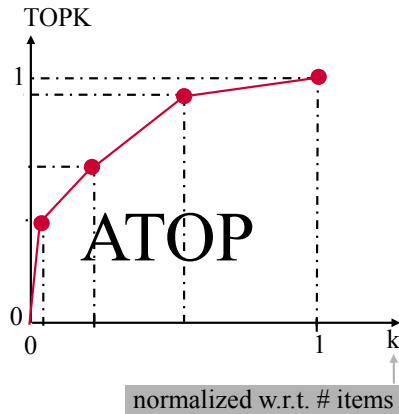
$$\frac{\text{\#relevant items in top } k}{\text{\#relevant items}}$$

Compute all recall@k...

until k match $R(u)$

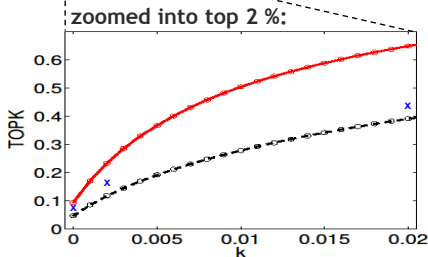
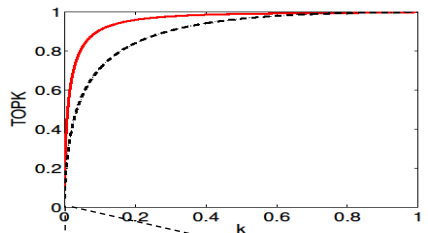
Compute the area under the curve

- focus on rated items
- numerical indicator + graphical details



H. Steck, KDD, 2010

Training and Testing of Recommender Systems on Data Missing Not at Random



39 % 50 % larger Top-k Hit-Rate: AllRank vs. integrated model

Comparison of Approaches:

- AllRank (RMSE = 1.106)
- ignore missing ratings (RMSE = 0.921)
- x integrated model [Koren '08] (RMSE = 0.887)
(trained to minimize RMSE)

Large increase in Top-k Hit-Rate when accounting also for missing ratings when training on MNAR data.

A/B testing & production launch

In a real situation

Designing an online Recommender System offers new performance indicators

- Online click, purchase, etc

A/B testing:

- 1 Defining some performance indicator with expert
- 2 Re-direct a small part of the customers to the new system *B*
 - make sure that the redirection is random (not biased)
- 3 Compare indicators from *A* and *B*

⇒ **Best evaluation...**

But only available **online** & with **access to the backoffice**

Serendipity : another important factor to evaluate...

... But very difficult to quantify

- Exploration / exploitation dilemma
- Clustering / categorization exploitation
 - propose items from different region
- Post processing / HMI issue

CF can offer serendipity

- increase neighborhood,
- increase implicit feedback weight
- ...



M. Ge et al., RecSys, 2010

Beyond Accuracy: Evaluating Recommender Systems by Coverage and Serendipity

Idea to design a metric

- 1 Learn a strong baseline (SVD)
- 2 New system RS
- 3 Unexpectedness = $RS \setminus SVD$
- 4 Serendipity =
usefulness(Unexpectedness)

CB is not well adapted

- Clustering heuristics
- bad performance

Conclusion

Conclusions

- For many applications such as Recommender Systems (but also Search, Advertising, and even Networks) understanding data and users is vital
- Algorithms can only be as good as the data they use as input
- But the inverse is also true: you need a good algorithm to leverage your data
- Importance of User/Data Mining is going to be a growing trend in many areas in the coming years
- RS have the potential to become as important as Search is now

⇒ there are still many open questions and a lot of interesting research to do!

Performance = industrial global deployment

Navigation traces & user generated contents = behavior sensors

- Collecting the data
 - Storing many trace for further exploitation
 - Large amount of data = cost
 - ... But trying to reduce the **noise**
 - missing/corrupted data, inadvertently operations...
 - ... and extract implicit **feedback** (and/or specific features)
 - e.g. video watching statistics
- UI integration
 - 50% \Rightarrow 90% of the job (Netflix !)
- ROI: return on invest
- Other expected benefits:
 - Fighting against **adversarial noise**
 - spam, web spam, review spam