

Introduction to Recommender Systems

ESSIR 2019 – Recommender Systems – Paolo Cremonesi



POLITECNICO
MILANO 1863

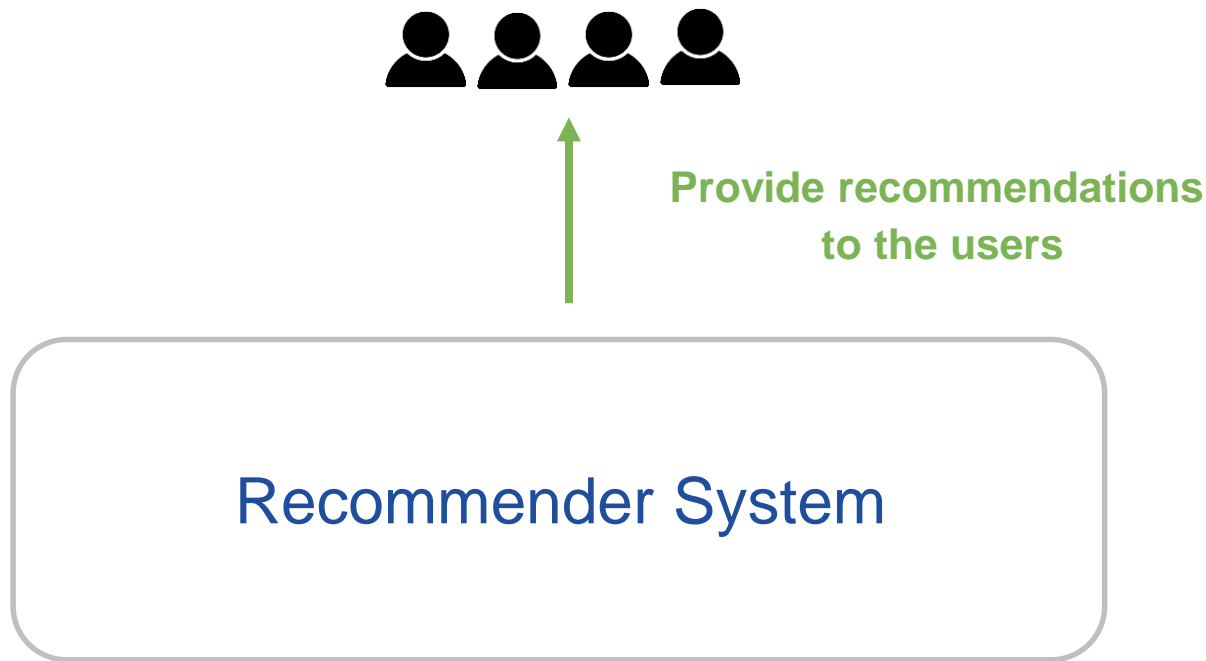
Recommender Systems



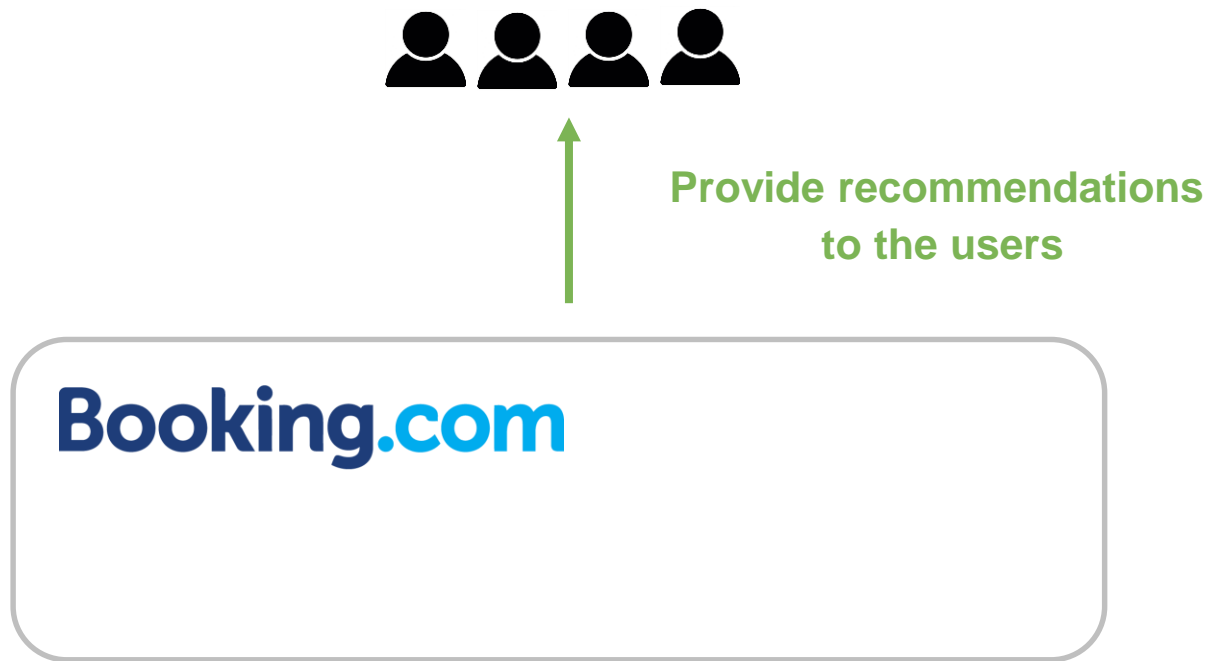
Recommender System



Recommender Systems



Recommender Systems



Recommender Systems



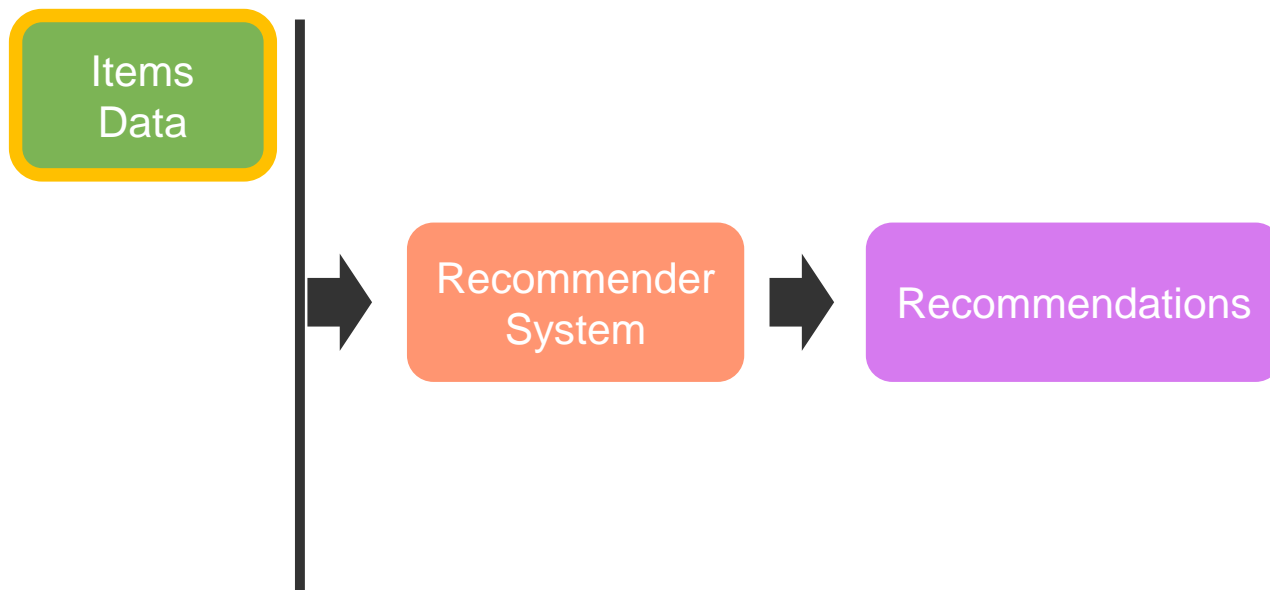
Recommender Systems



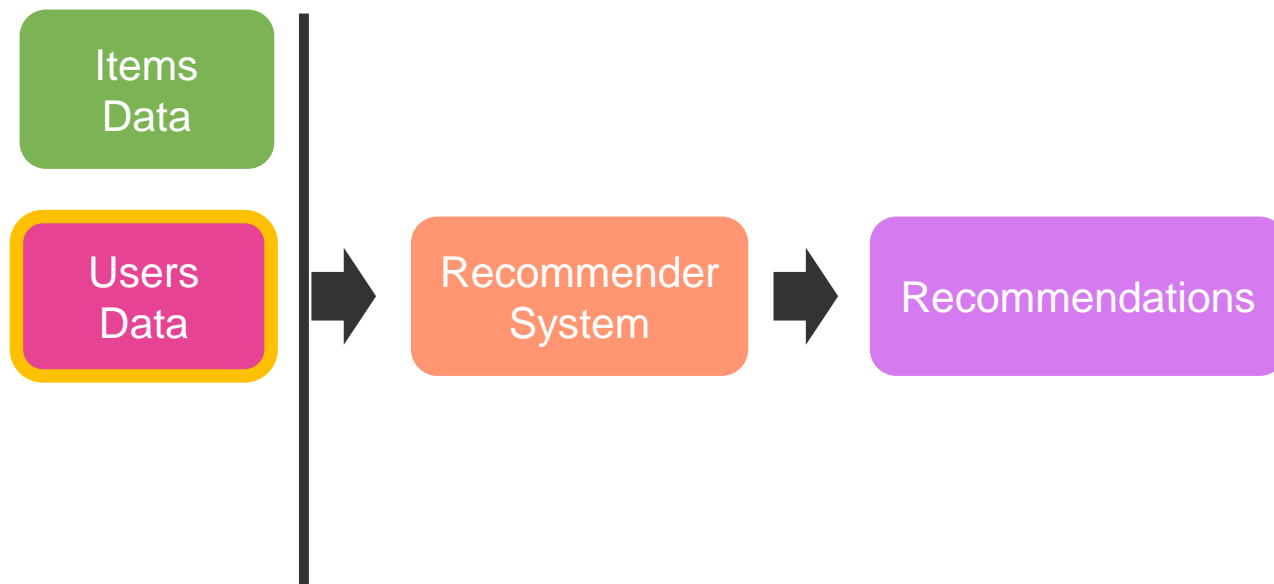
Recommender Systems



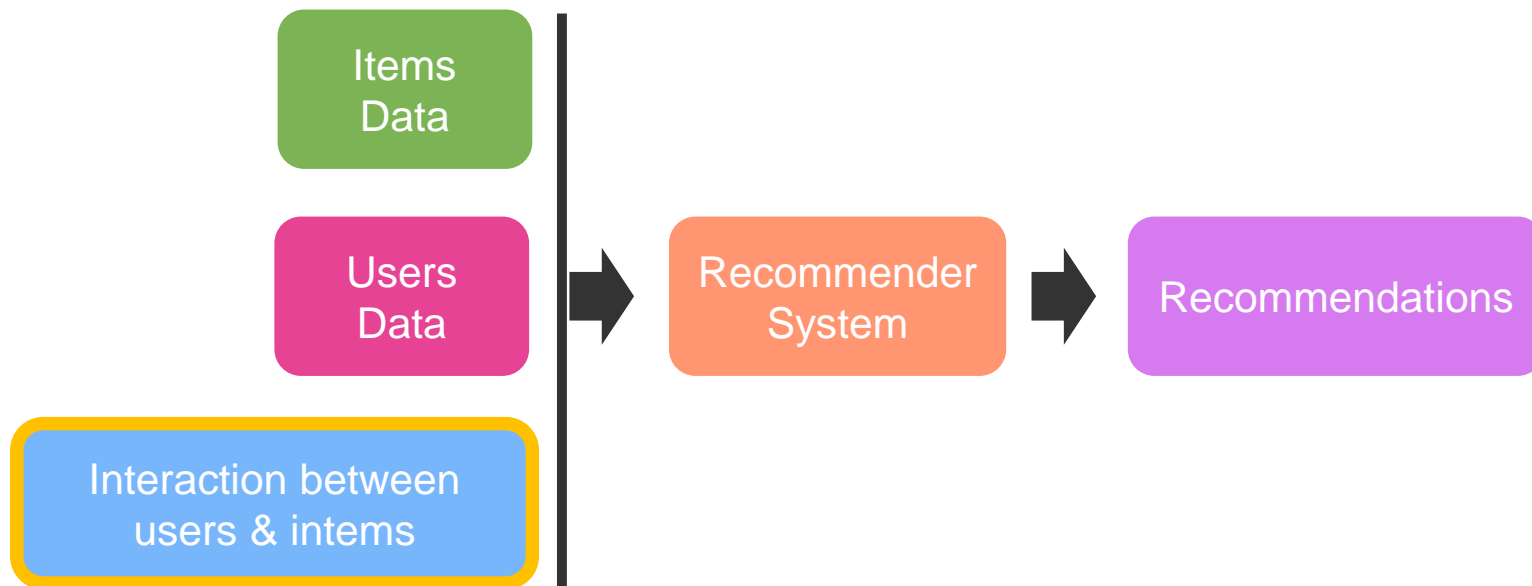
Input Data



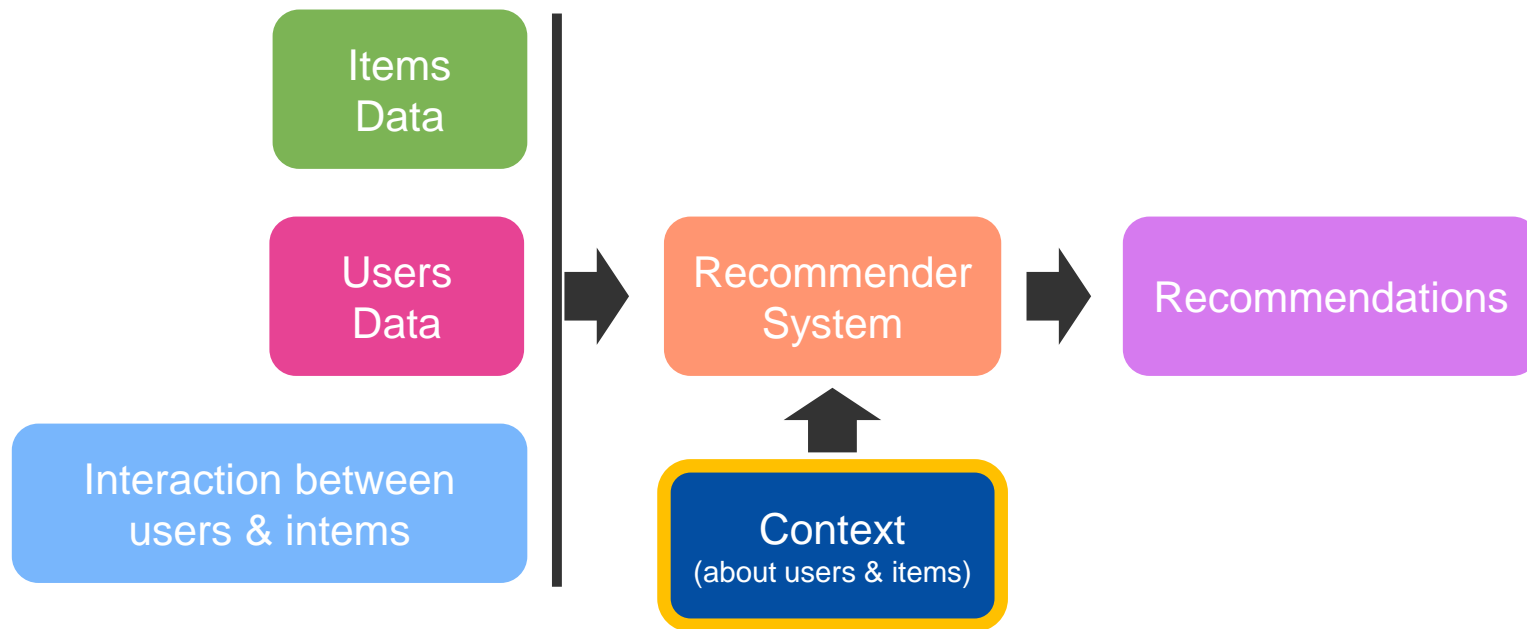
Input Data



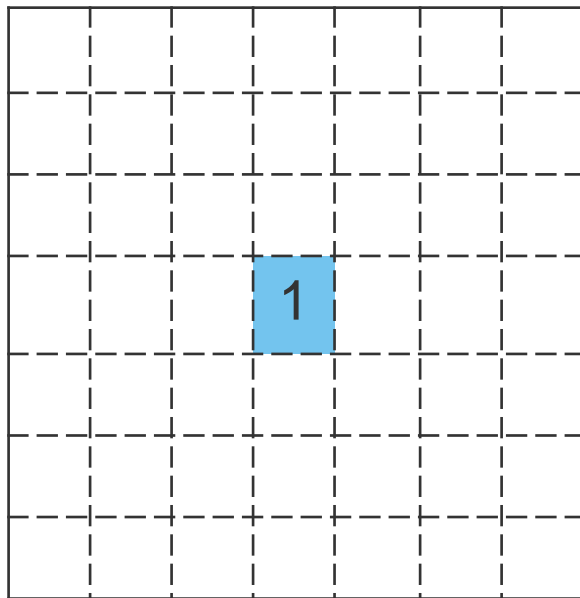
Input Data



Input Data

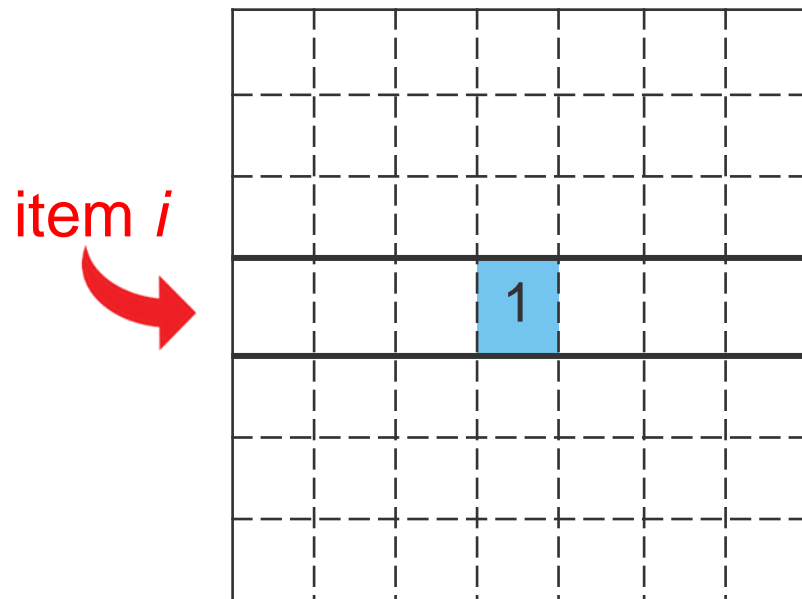


Item Content Matrix (ICM)

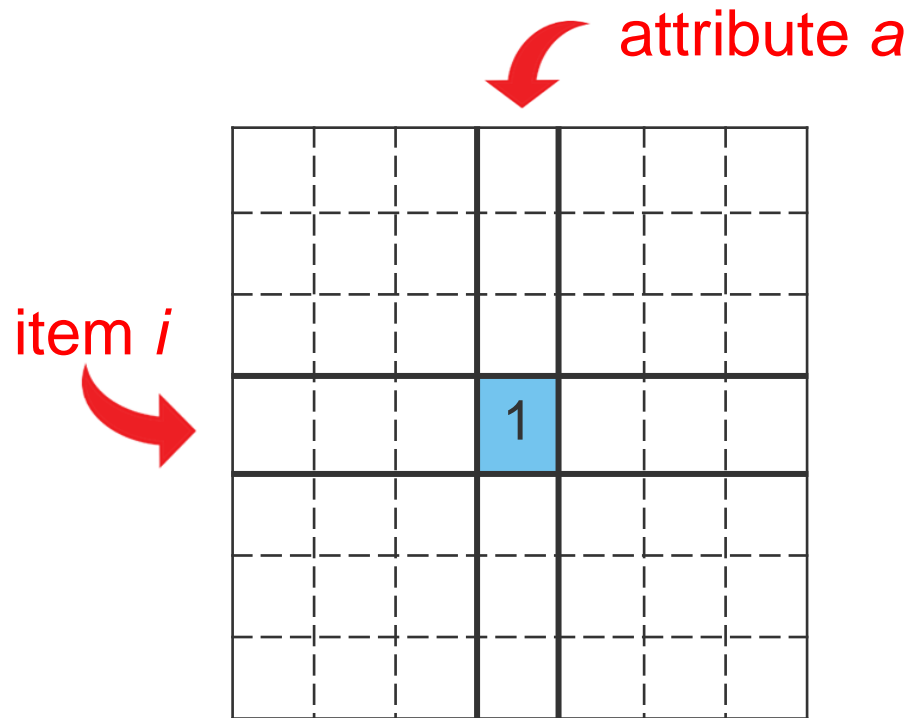


			1			

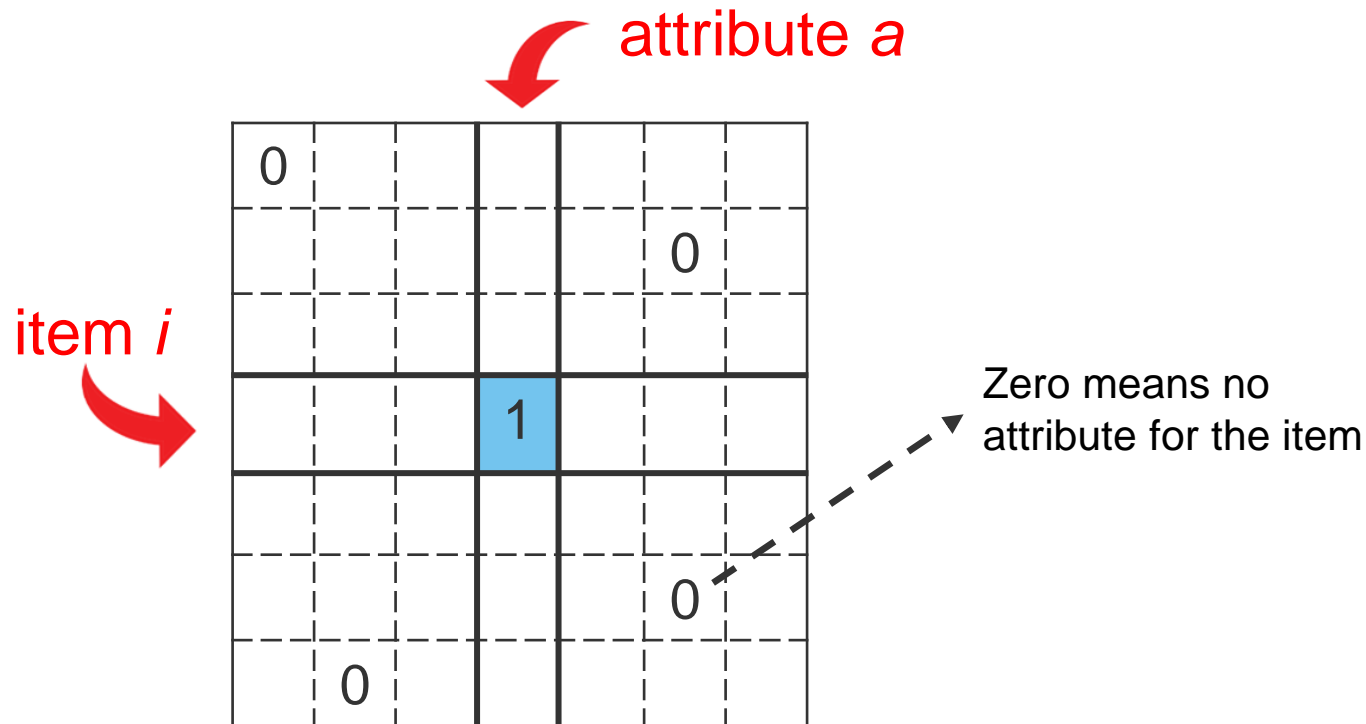
Item Content Matrix (ICM)



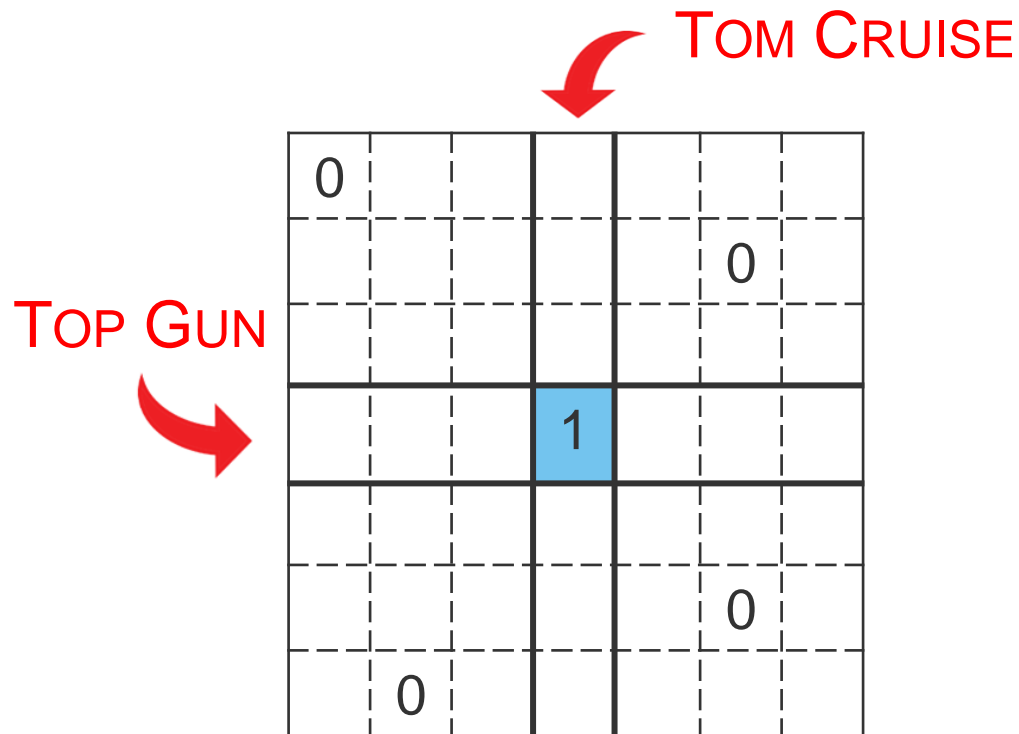
Item Content Matrix (ICM)



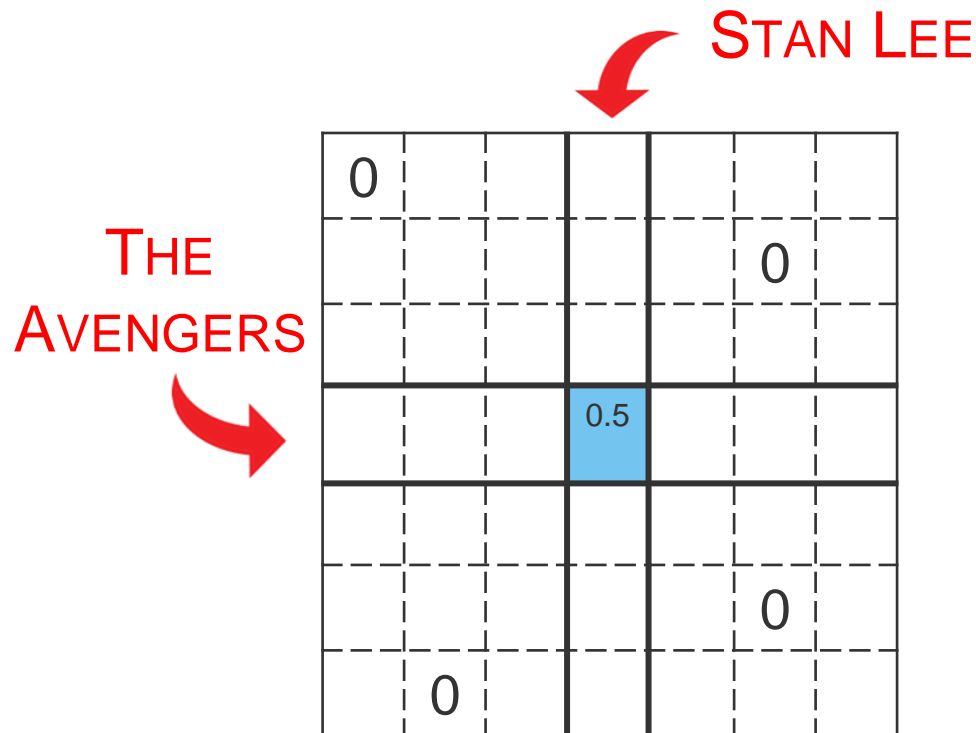
Item Content Matrix (ICM)



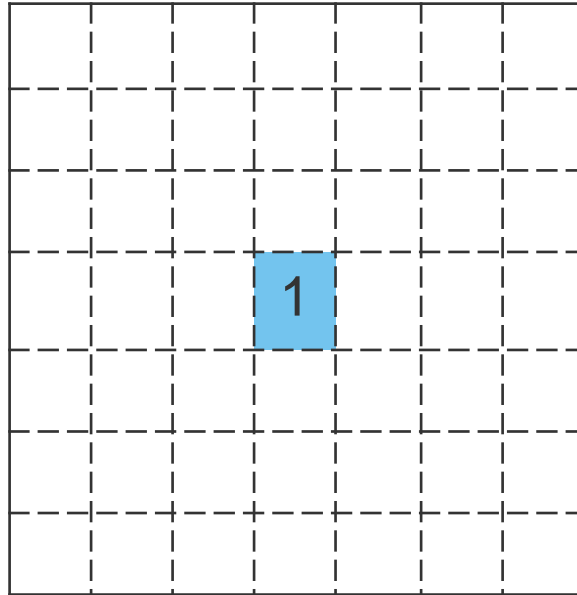
Item Content Matrix (ICM)



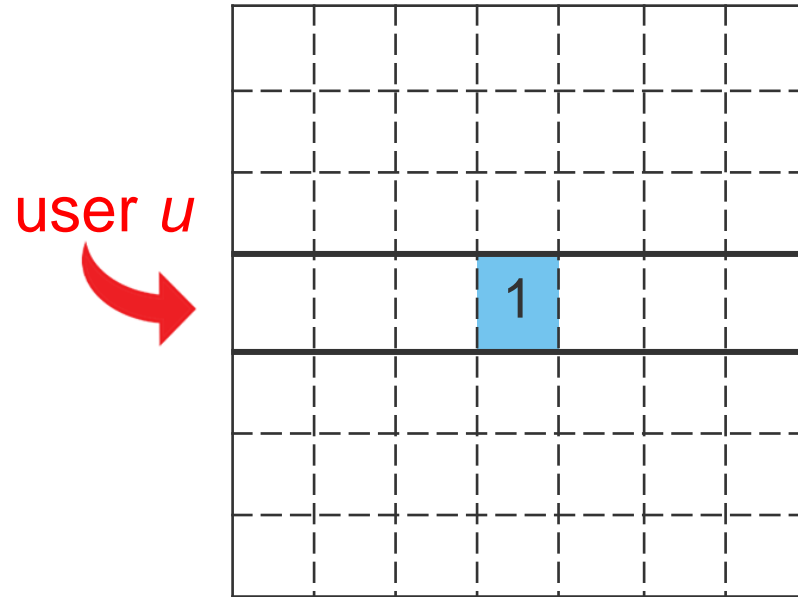
Item Content Matrix (ICM)



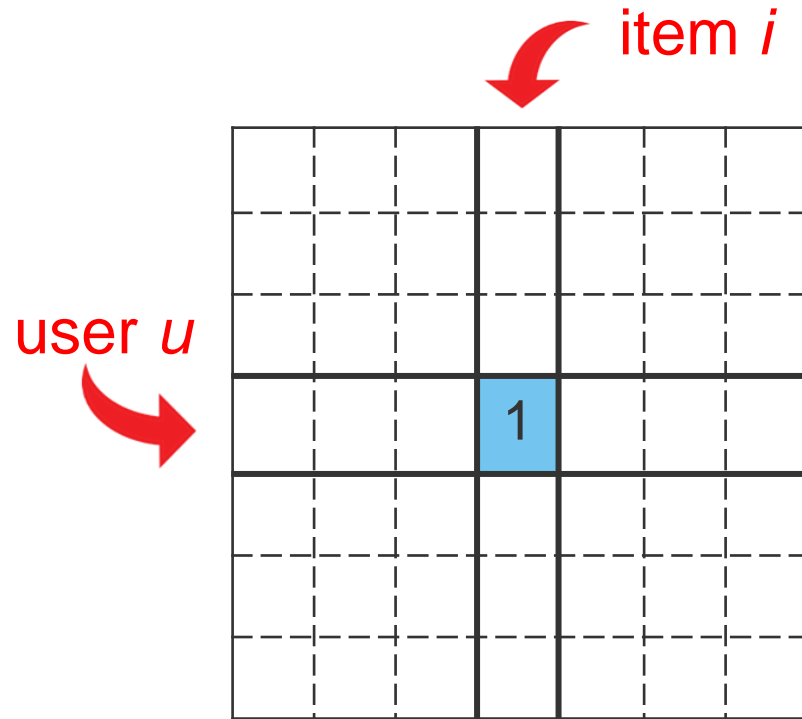
User Rating Matrix (URM)



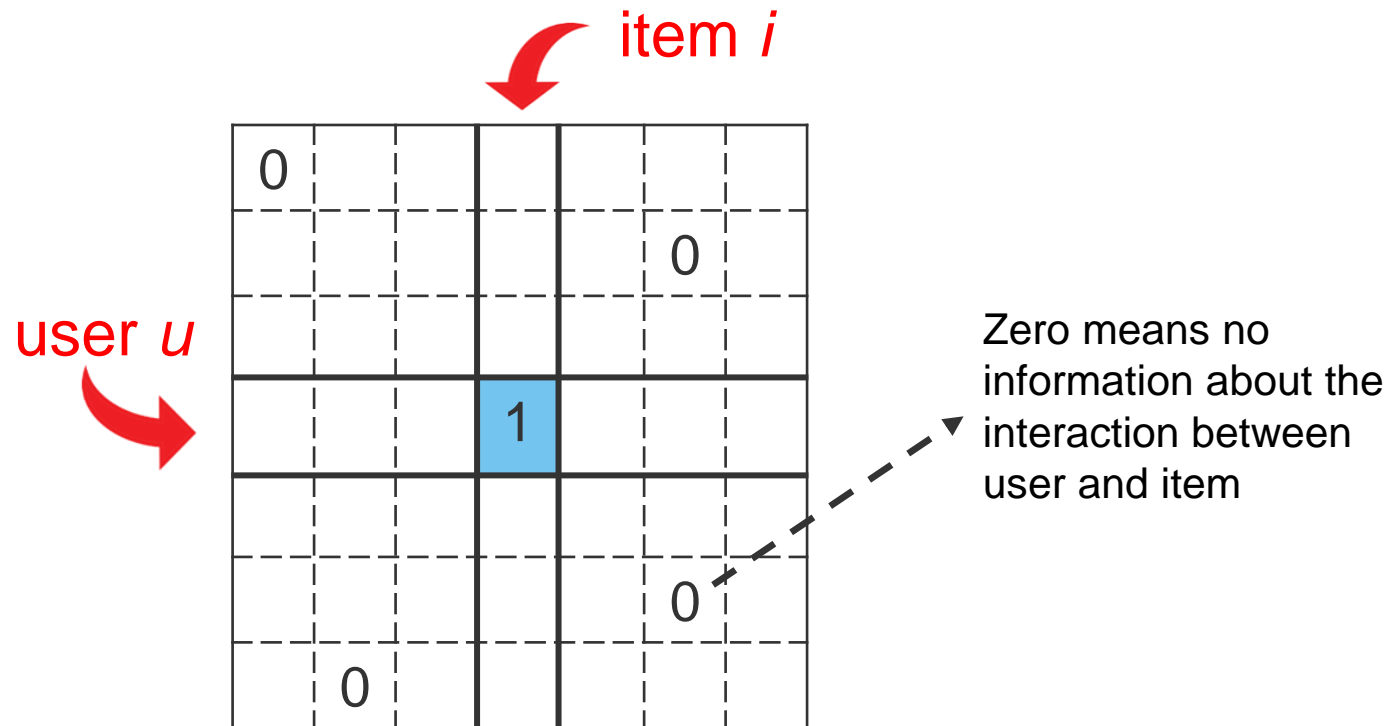
User Rating Matrix (URM)



User Rating Matrix (URM)



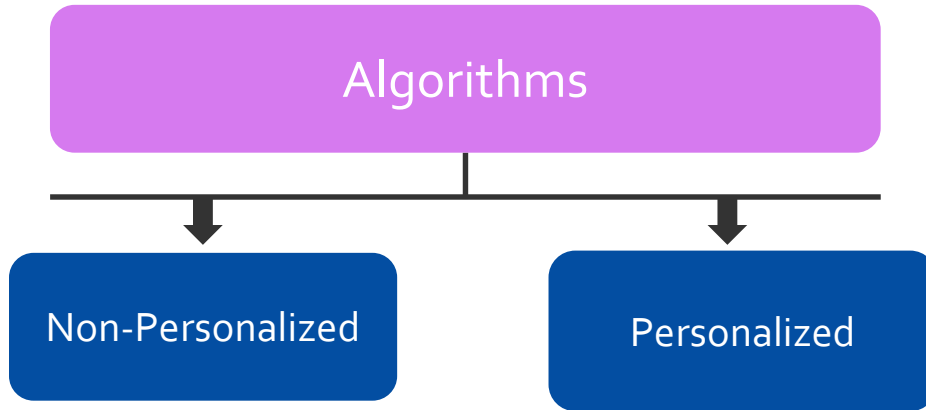
User Rating Matrix (URM)



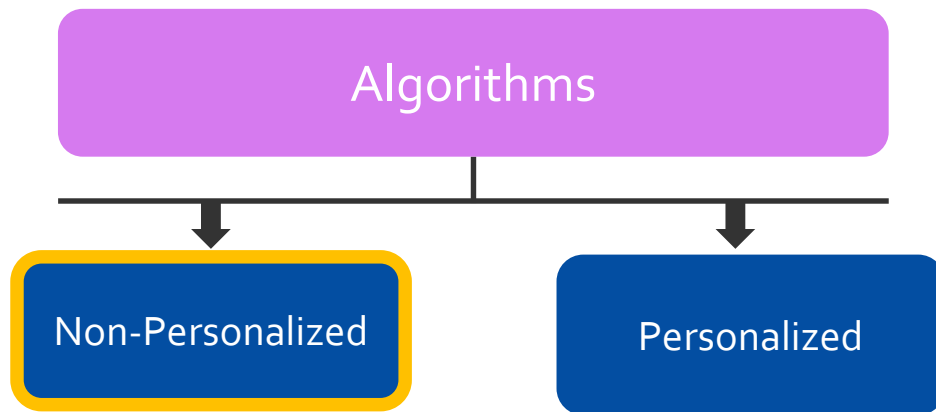
Taxonomy of Recommender Systems



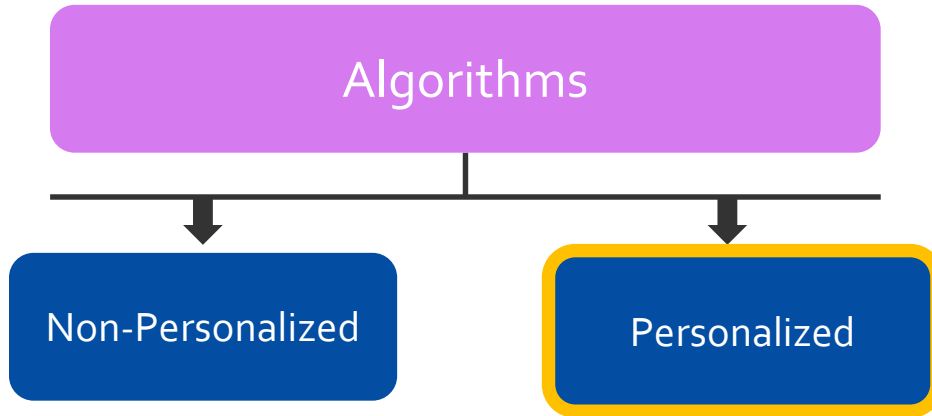
Taxonomy of Rec. Systems



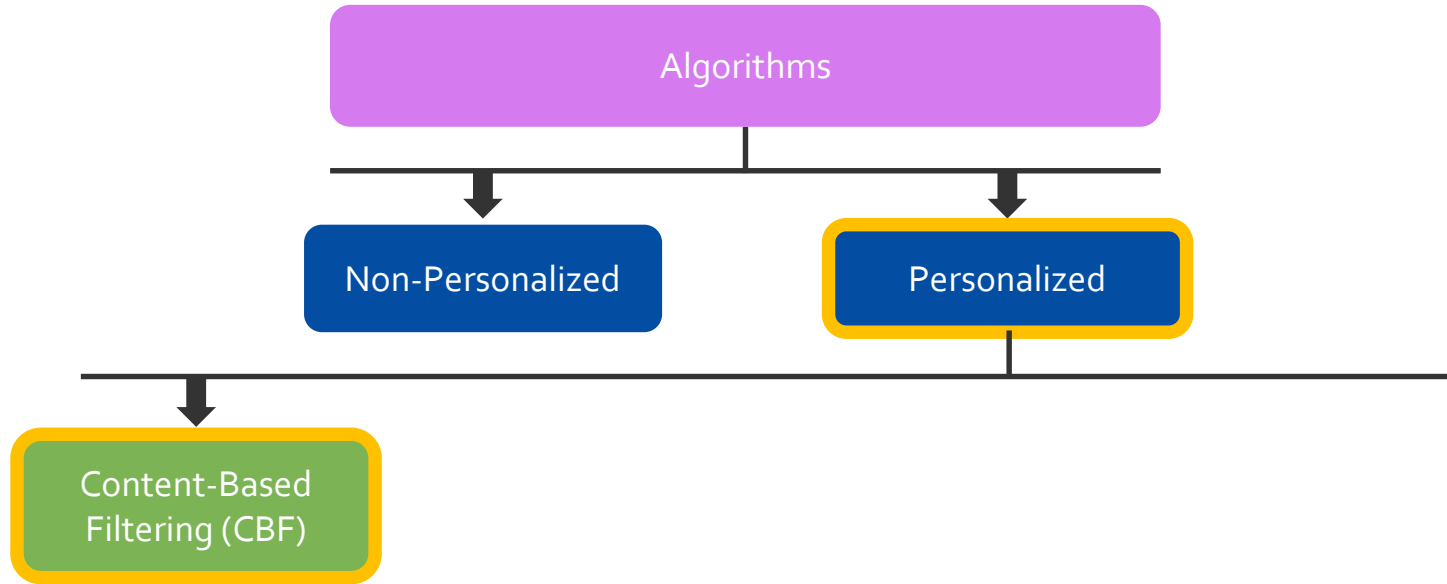
Taxonomy of Rec. Systems



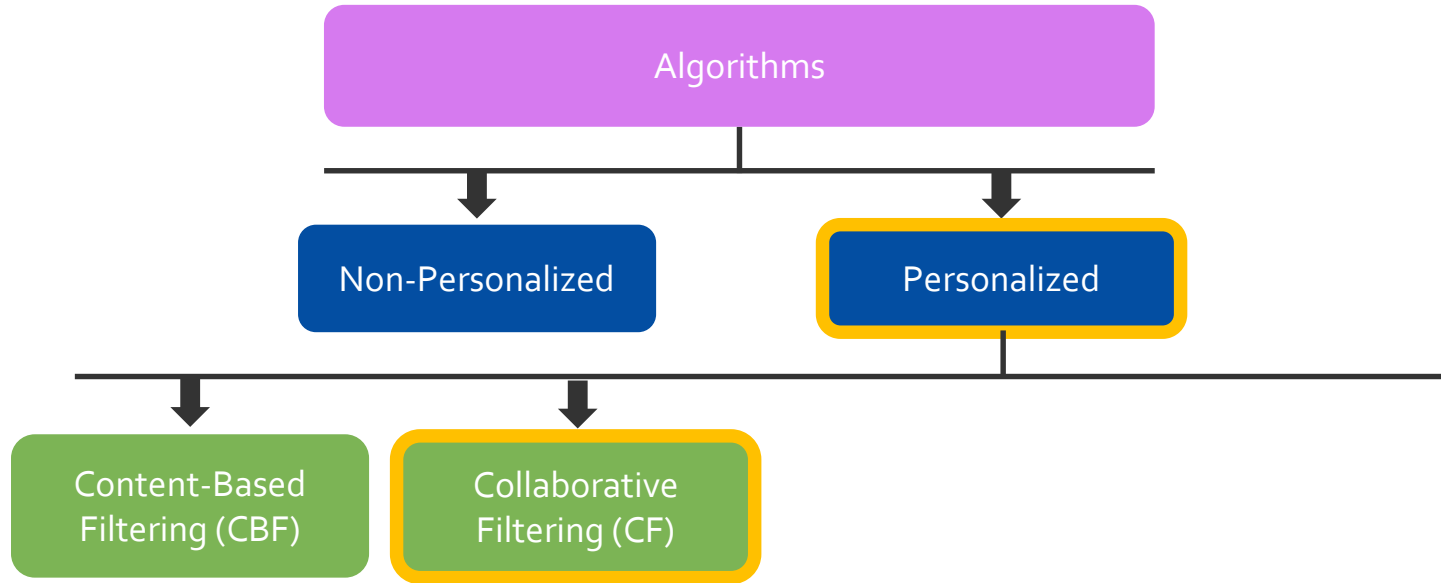
Taxonomy of Rec. Systems



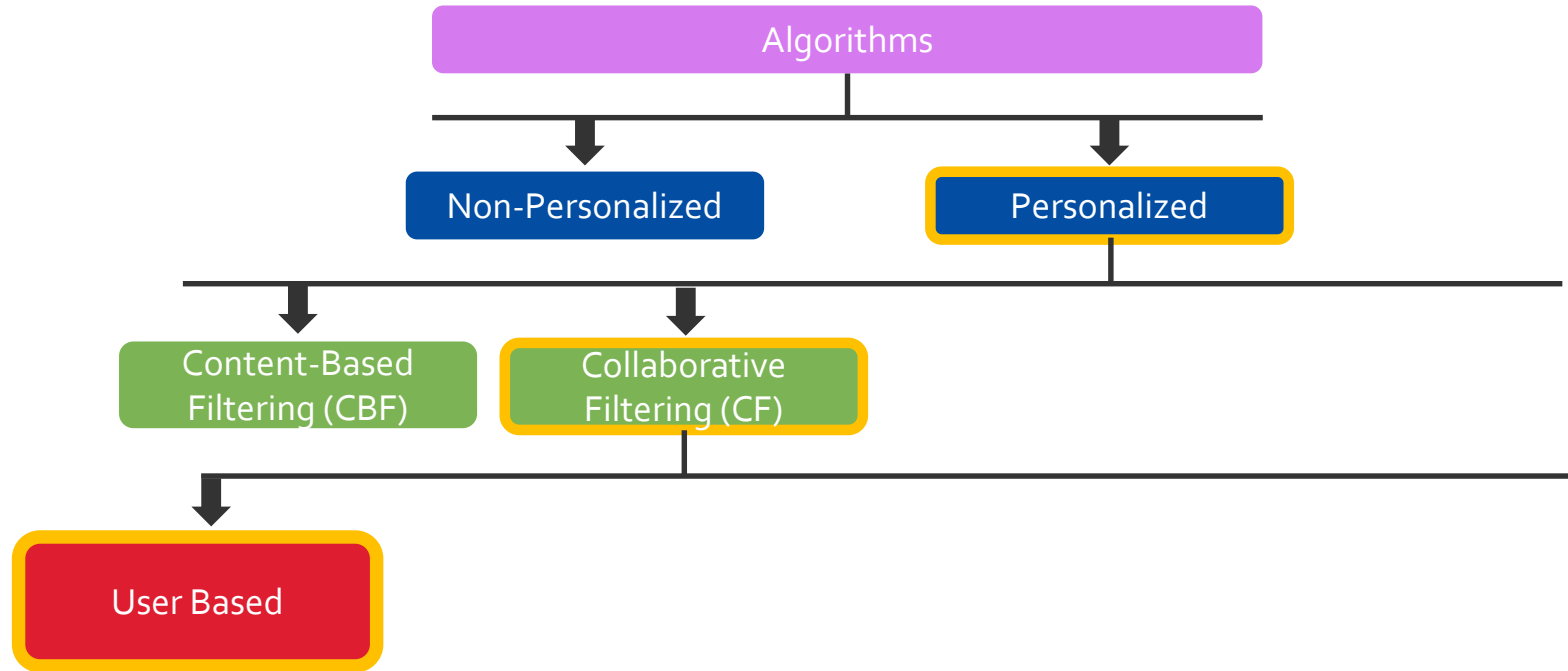
Taxonomy of Rec. Systems



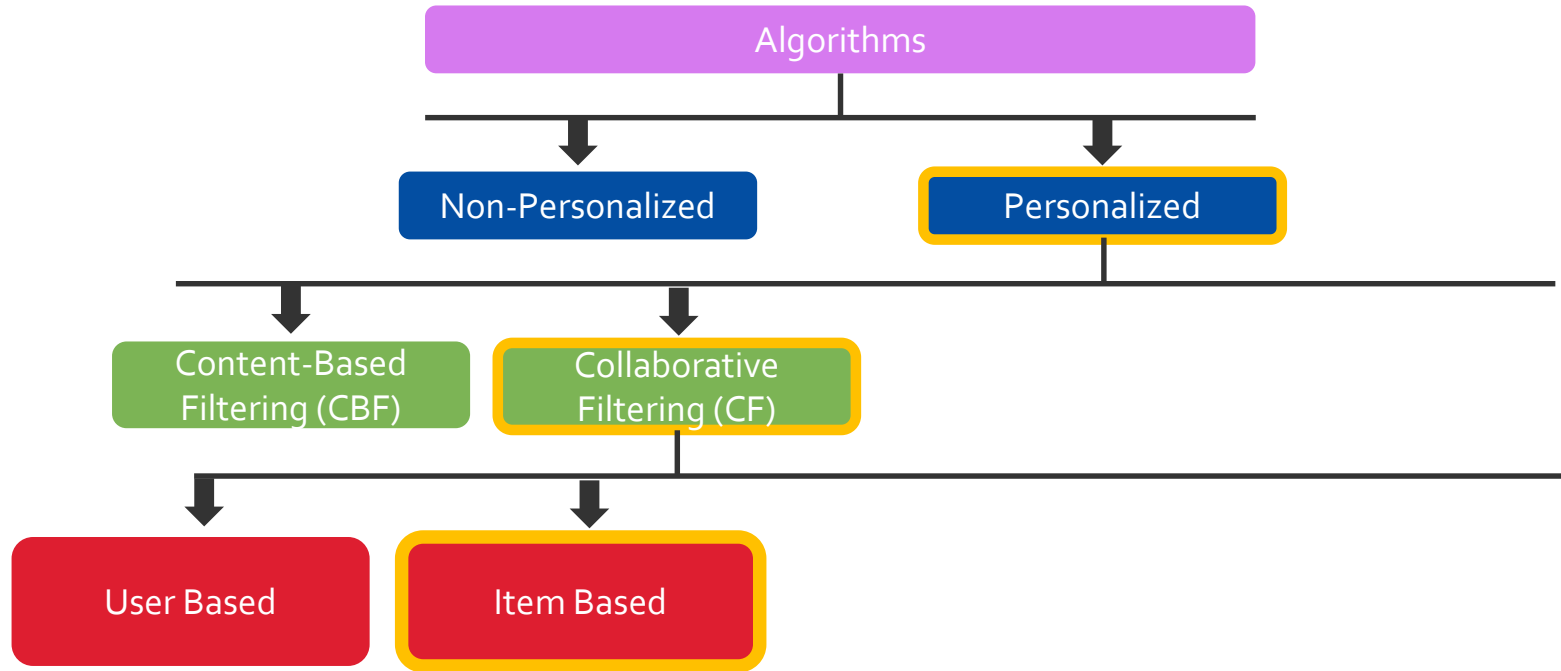
Taxonomy of Rec. Systems



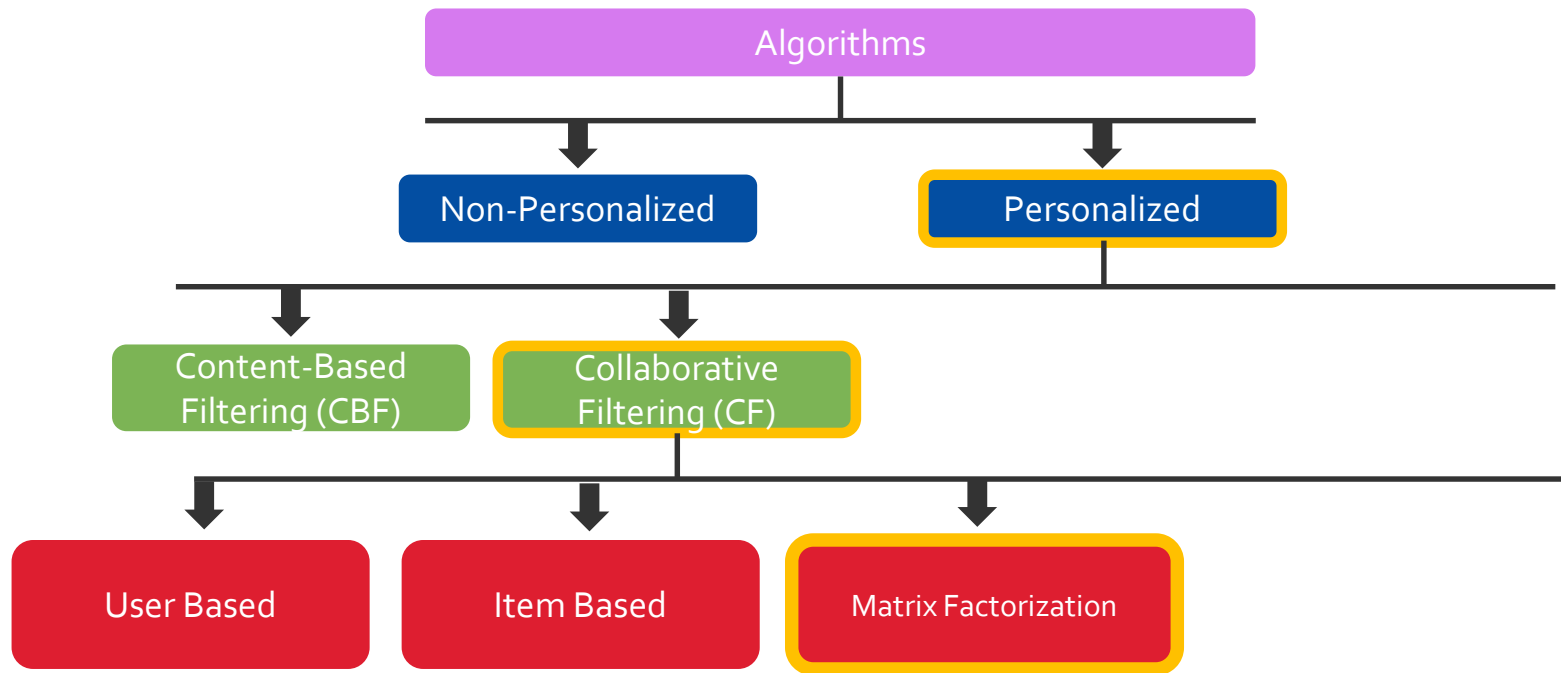
Taxonomy of Rec. Systems



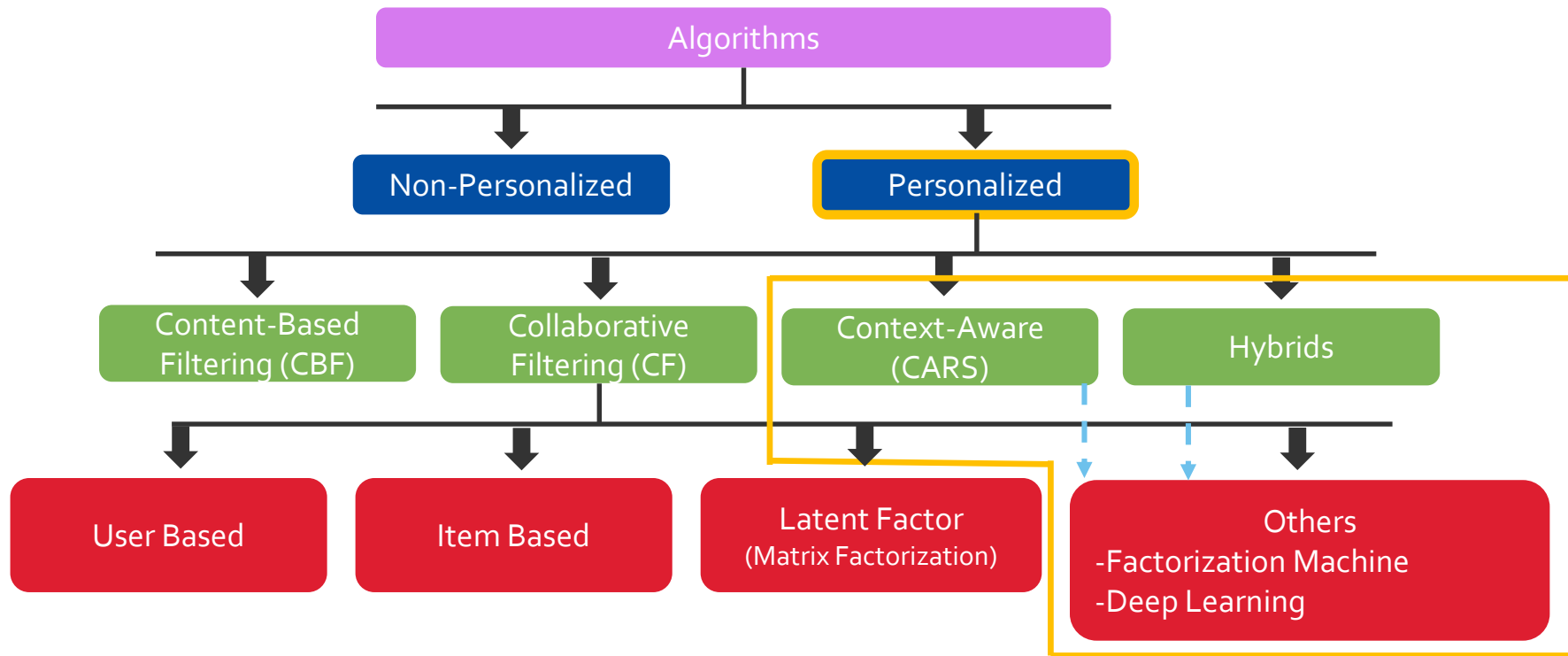
Taxonomy of Rec. Systems



Taxonomy of Rec. Systems



Taxonomy of Rec. Systems



Ratings, predictions and recommendations



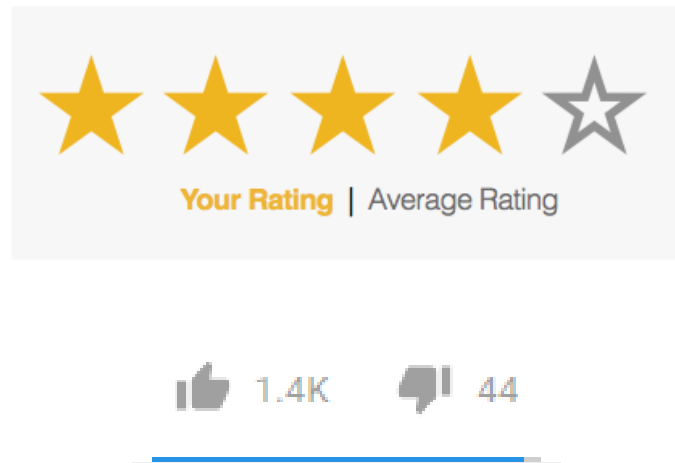
Rating Systems

Explicit ratings



Rating Systems

Explicit ratings



Ratings Distribution



Rating Systems

Explicit ratings

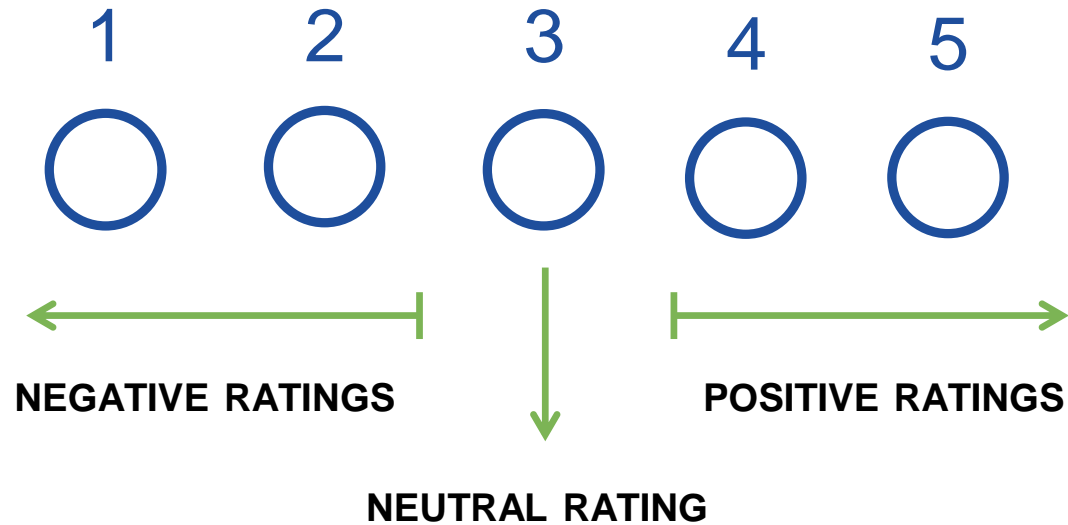
Implicit ratings



Even Ratings Scale



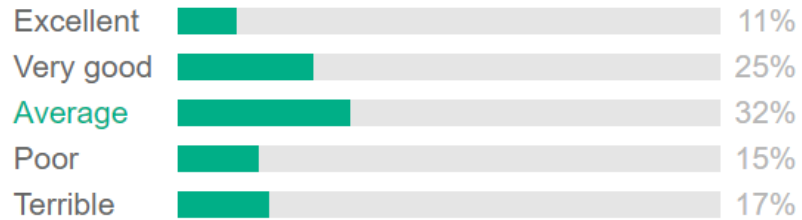
Odd Ratings Scale



Ratings Distribution

3.0 

51 reviews



Rating Systems

Explicit ratings

Implicit ratings

- Viewing time of a movie
- Song streaming count
- Purchases
- Other user activity (e.g. opened links)

Inferring Preferences

$URM = u$

i

0						
					0	
			r_{ui}			
					0	
	0					

r_{ui} = Rating that user u gave to item i

Inferring Preferences

Diagram illustrating a 2D lattice structure. The horizontal axis is labeled i and the vertical axis is labeled j . The lattice is divided into four quadrants by a central vertical line and a central horizontal line. The central cell, at the intersection of the two lines, is highlighted in blue and labeled r_{ui} . The four quadrants are labeled 0 in the top-left, top-right, bottom-left, and bottom-right corners.

 $r_{ui} \in \{0,1\}$ (Implicit)
$$r_{ui} = \text{Rating that user } u \text{ gave to item } i$$

Inferring Preferences

$URM = u$

i					
0					
				0	
			r_{ui}		
				0	
	0				

$r_{ui} \in \{0,1\} \leftarrow$ implicit

$r_{ui} \in \{1,2,3,4,5\} \leftarrow$ explicit

r_{ui} = Rating that user u gave to item i

URM Density

typical URM density < 0.01 %



URM Density

typical URM density < 0.01 %

Netflix URM density $\approx 0.002\%$

MovieLens URM density $\approx 0.005\%$

Non-Personalized Recommenders



Non-Personalized Recommenders

Top Popular



Non-Personalized Recommenders

Top Popular



	i				
u	3	0	0	1	2
	4	2	0	2	3
	0	0	0	1	2
	0	0	5	0	0
	3	0	1	0	5

Non-Personalized Recommenders

Top Popular



	i				
u	3	0	0	1	2
	4	2	0	2	3
	0	0	0	1	2
	0	0	5	0	0
	3	0	1	0	5

3	1	1	3	4
---	---	---	---	---

ratings per item

Non-Personalized Recommenders

Top Popular



	<i>i</i>				<i>k</i>
<i>u</i>	3	0	0	1	2
	4	2	0	2	3
	0	0	0	1	2
	0	0	5	0	0
	3	0	1	0	5
	3	1	1	3	4

item *k* is the
top popular
↓
largest number
of ratings

Non-Personalized Recommenders

Top Popular

Best Rated



Non-Personalized Recommenders

Top Popular

Best Rated

	i				
u	3	0	0	1	2
	4	2	0	2	3
	0	0	0	1	2
	0	0	5	0	0
	3	0	1	0	5

Non-Personalized Recommenders

Top Popular

Best Rated

	i				
u	3	0	0	1	2
	4	2	0	2	3
	0	0	0	1	2
	0	0	5	0	0
	3	0	1	0	5

3.3	2	3	1.3	3
-----	---	---	-----	---

Avg ratings per item

Non-Personalized Recommenders

Top Popular

Best Rated

	j	i			
u	3	0	0	1	2
	4	2	0	2	3
	0	0	0	1	2
	0	0	5	0	0
	3	0	1	0	5
	3.3	2	3	1.3	3

item j is the
best rated

↓
largest
average
rating

Item bias

Average rating of item i

$$b_i = \frac{\sum_u r_{ui}}{N_i}$$



Item bias

Average rating of item i

$$b_i = \frac{\sum_u r_{ui}}{N_i}$$

r_{ui} : rating given by user u to item i (non zero ratings)

N_i : number of users who rated item i

Item bias: support

Shrunked avg. rating of item i

$$b_i = \frac{\sum_u r_{ui}}{N_i + C}$$

Item bias: support

Shrunked avg. rating of item i

$$b_i = \frac{\sum_u r_{ui}}{N_i + C}$$

r_{ui} : rating given by user u to item i (non zero ratings)

N_i : number of users who have rated item i

C : shrink term (constant value)

Global Effects



Global Effects: Step 1

Avg. ratings for all items and users

$$\mu = \frac{\sum_i \sum_u r_{ui}^+}{N^+}$$



Global Effects: Step 1

Avg. ratings for all items and users

$$\mu = \frac{\sum_i \sum_u r_{ui}^+}{N^+}$$

μ : overall average of ratings, for all users and all items

r_{ui}^+ : explicit rating given by user u to item i

N^+ : total number of non zero ratings

*Note: the + symbol denotes that we are not computing this average on the full URM, but **only** on the **non zero elements***

Global Effects: Step 2

Normalized rating

$$r'_{ui} = r_{ui}^+ - \mu$$

to be computed for each user u and item i , only on non-zero ratings

Global Effects: Step 3

Item bias

$$b_i = \frac{\sum_u r'_{ui}}{N_i + C}$$



Global Effects: Step 3

Item bias

$$b_i = \frac{\sum_u r'_{ui}}{N_i + C}$$

N_i : number of users who have rated item i
to be computed for each item i

Global Effects: Step 4

Recompute rating

$$r''_{ui} = r'_{ui} - b_i$$

to be computed for each user u and item i , only on non-zero ratings

Global Effects: Step 5

User bias

$$b_u = \frac{\sum_i r_{ui}''}{N_u + C}$$

Global Effects: Step 5

User bias

$$b_u = \frac{\sum_i r_{ui}''}{N_u + C}$$

N_u : number of items i rated by user u

to be computed for each user u

Global Effects: Step 6

Global effects final formula estimated rating

$$\tilde{r}_{ui} = \mu + b_i + b_u$$

to be computed for each user u and item i , only on non-zero ratings

Global Effects: Recap

- Step 1: compute the average of all ratings (μ)

$$\tilde{r}_{ui} = \mu + b_i + b_u$$



Global Effects: Recap

- Step 1: compute the average of all ratings (μ)
- Step 2: remove this quantity from the URM

$$\tilde{r}_{ui} = \mu + b_i + b_u$$



Global Effects: Recap

- Step 1: compute the average of all ratings (μ)
- Step 2: remove this quantity from the URM
- Step 3: compute the bias for each item (b_i)

$$\tilde{r}_{ui} = \mu + b_i + b_u$$

Global Effects: Recap

- Step 1: compute the average of all ratings (μ)
- Step 2: remove this quantity from the URM
- Step 3: compute the bias for each item (b_i)
- Step 4: remove this quantity from the URM

$$\tilde{r}_{ui} = \mu + b_i + b_u$$



Global Effects: Recap

- Step 1: compute the average of all ratings (μ)
- Step 2: remove this quantity from the URM
- Step 3: compute the bias for each item (b_i)
- Step 4: remove this quantity from the URM
- Step 5: compute the bias for each user (b_u)

$$\tilde{r}_{ui} = \mu + b_i + b_u$$

Global Effects: Recap

- Step 1: compute the average of all ratings (μ)
- Step 2: remove this quantity from the URM
- Step 3: compute the bias for each item (b_i)
- Step 4: remove this quantity from the URM
- Step 5: compute the bias for each user (b_u)
- Step 6: final formula creating a new URM

$$\tilde{r}_{ui} = \mu + b_i + b_u$$

Evaluation of Recommender Systems

ESSIR 2019 – Recommender Systems – Paolo Cremonesi



POLITECNICO
MILANO 1863

FUNCTIONAL REQUIREMENTS

Requirements

FUNCTIONAL REQUIREMENTS

What the software does



Requirements

FUNCTIONAL REQUIREMENTS

What the software does

NON-FUNCTIONAL REQUIREMENTS

Requirements

FUNCTIONAL REQUIREMENTS

What the software does

NON-FUNCTIONAL REQUIREMENTS

How the software does its job

Non-Functional Requirements

RESPONSE TIME



Non-Functional Requirements

RESPONSE TIME

- How long does it take for the system to generate one recommendation?

Non-Functional Requirements

RESPONSE TIME

SCALABILITY



Non-Functional Requirements

RESPONSE TIME

SCALABILITY

- How many recommendations per second the system is able to generate?

Non-Functional Requirements

RESPONSE TIME

SCALABILITY

PRIVACY AND SECURITY

Non-Functional Requirements

RESPONSE TIME

SCALABILITY

PRIVACY AND SECURITY

- Protect against reverse engineering
- Protect against intrusions from outside

Non-Functional Requirements

RESPONSE TIME

SCALABILITY

PRIVACY AND SECURITY

USER INTERFACE

Non-Functional Requirements

RESPONSE TIME

SCALABILITY

PRIVACY AND SECURITY

USER INTERFACE

- Which is the best place to show recommendations?
- How many items should be recommended?

Quality indicators for Recommender Systems



Quality Indicators

RELEVANCE



Quality Indicators

RELEVANCE

- Recommend items that users like

Quality Indicators

RELEVANCE

COVERAGE



Quality Indicators

RELEVANCE

COVERAGE

- Ability to recommend most of the items in a catalogue



Quality Indicators

RELEVANCE

COVERAGE

NOVELTY



Quality Indicators

RELEVANCE

COVERAGE

NOVELTY

- Recommend items unknown to the user

Quality Indicators

RELEVANCE

COVERAGE

NOVELTY

DIVERSITY



Quality Indicators

RELEVANCE

COVERAGE

NOVELTY

DIVERSITY

- Diversify the items recommended

Quality Indicators

RELEVANCE

COVERAGE

NOVELTY

DIVERSITY

CONSISTENCY

Quality Indicators

RELEVANCE

COVERAGE

NOVELTY

DIVERSITY

CONSISTENCY

- Recommendations should not change to often

Quality Indicators

RELEVANCE

COVERAGE

NOVELTY

DIVERSITY

CONSISTENCY

CONFIDENCE



Quality Indicators

RELEVANCE

COVERAGE

NOVELTY

DIVERSITY

CONSISTENCY

CONFIDENCE

- How much a system is sure about a recommendation

Quality Indicators

RELEVANCE

COVERAGE

NOVELTY

DIVERSITY

CONSISTENCY

CONFIDENCE

SERENDIPITY



Quality Indicators

RELEVANCE

COVERAGE

NOVELTY

DIVERSITY

CONSISTENCY

CONFIDENCE

SERENDIPITY

- The ability of *surprising* the user
- The ability to recommend items that users would have never been able to discover by themselves

Evaluation Techniques



Evaluation Techniques

ONLINE



Evaluation Techniques

ONLINE

OFF-LINE



Online Evaluation

DIRECT USER FEEDBACK



Online Evaluation

DIRECT USER FEEDBACK



Online Evaluation

DIRECT USER FEEDBACK

A/B TESTING



Online Evaluation

DIRECT USER FEEDBACK

A/B TESTING



Online Evaluation

DIRECT USER FEEDBACK

A/B TESTING

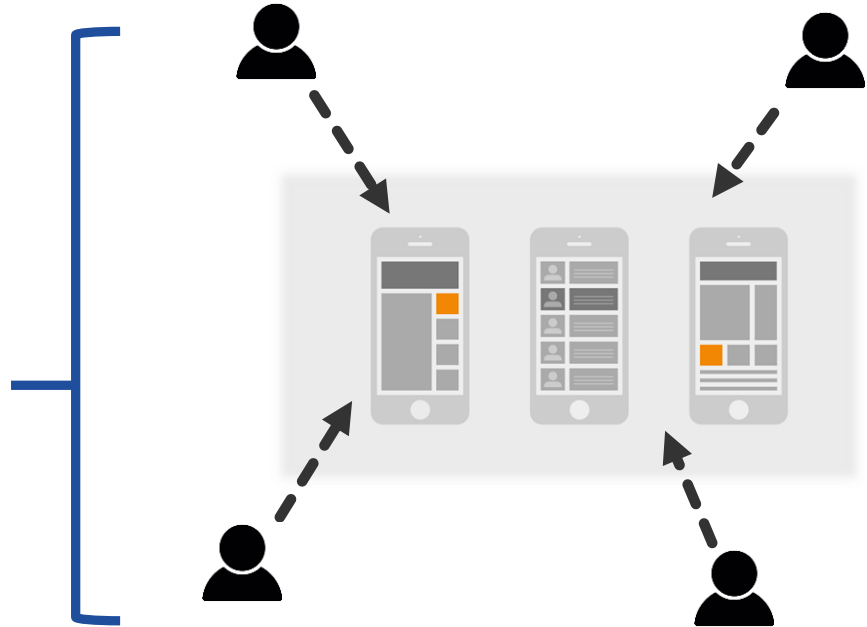
CONTROLLED EXPERIMENTS

Online Evaluation

DIRECT USER FEEDBACK

A/B TESTING

CONTROLLED EXPERIMENTS



Online Evaluation

DIRECT USER FEEDBACK

A/B TESTING

CONTROLLED EXPERIMENTS

CROWDSOURCING

Online Evaluation

DIRECT USER FEEDBACK

A/B TESTING

CONTROLLED EXPERIMENTS

CROWDSOURCING



Off-line Evaluation

TASK



Off-line Evaluation

TASK

DATASET

Off-line Evaluation

TASK

DATASET

PARTITIONING



Off-line Evaluation

TASK

DATASET

PARTITIONING

METRICS



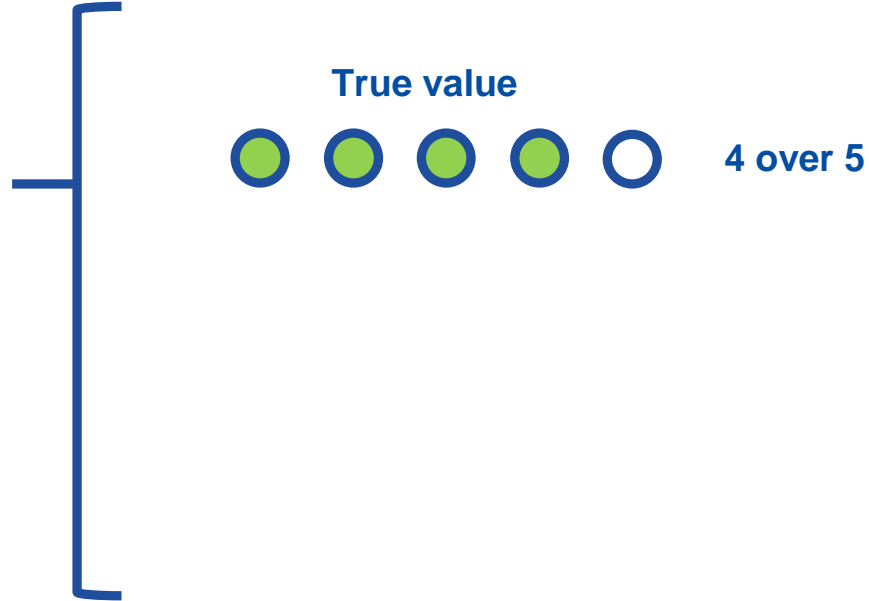
Off-line Evaluation: Task

RATING PREDICTION



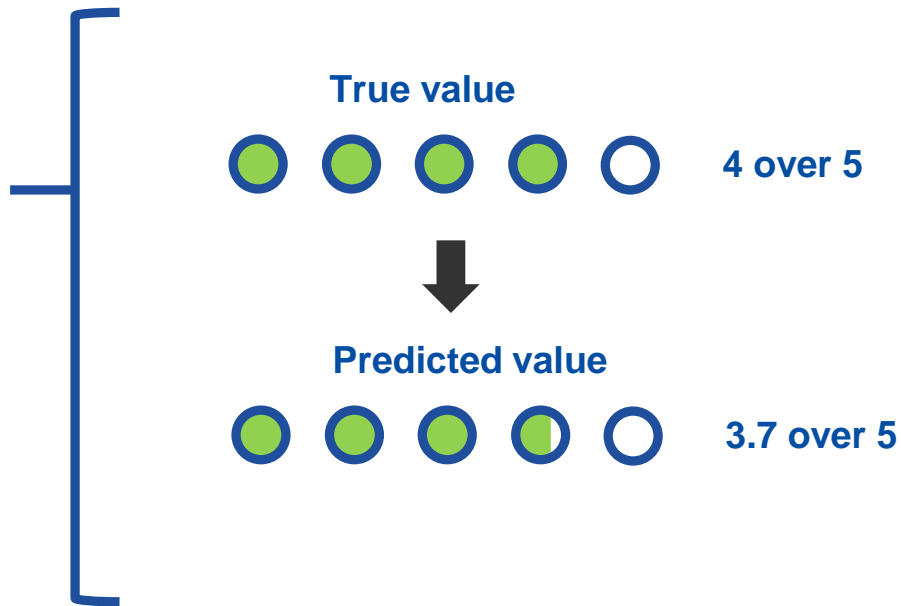
Off-line Evaluation: Task

RATING PREDICTION



Off-line Evaluation: Task

RATING PREDICTION



Off-line Evaluation: Task

RATING PREDICTION

TOP-N RECOMMENDATION



Off-line Evaluation: Task

RATING PREDICTION

TOP-N RECOMMENDATION



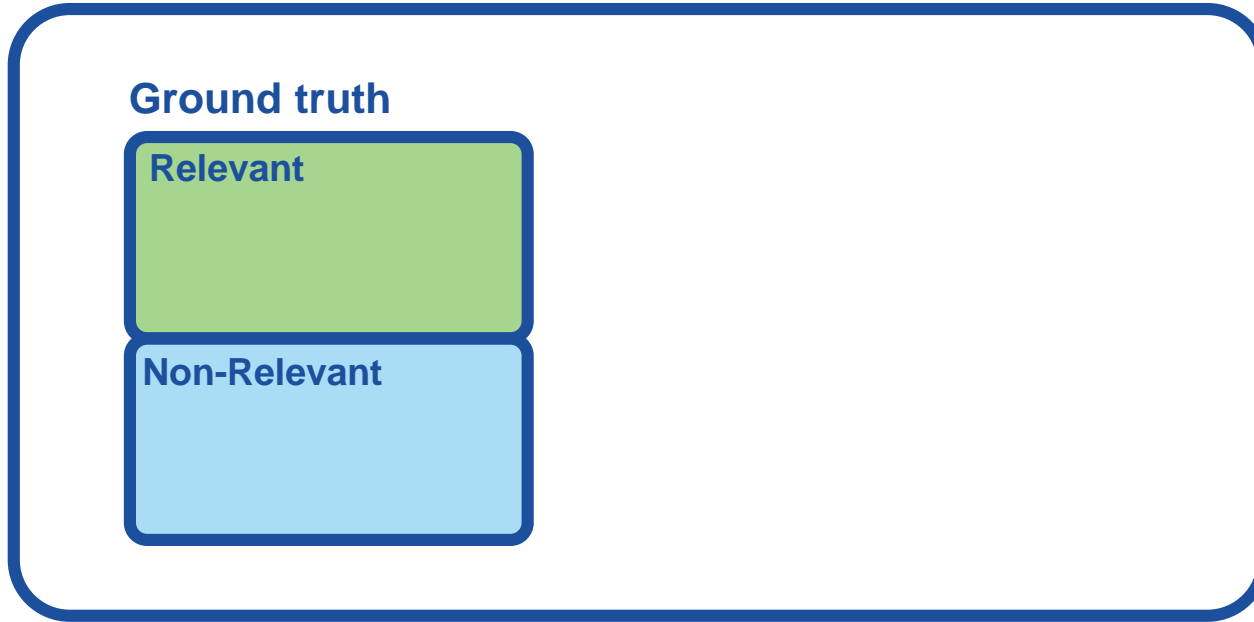
Off-line Evaluation: Dataset

Dataset



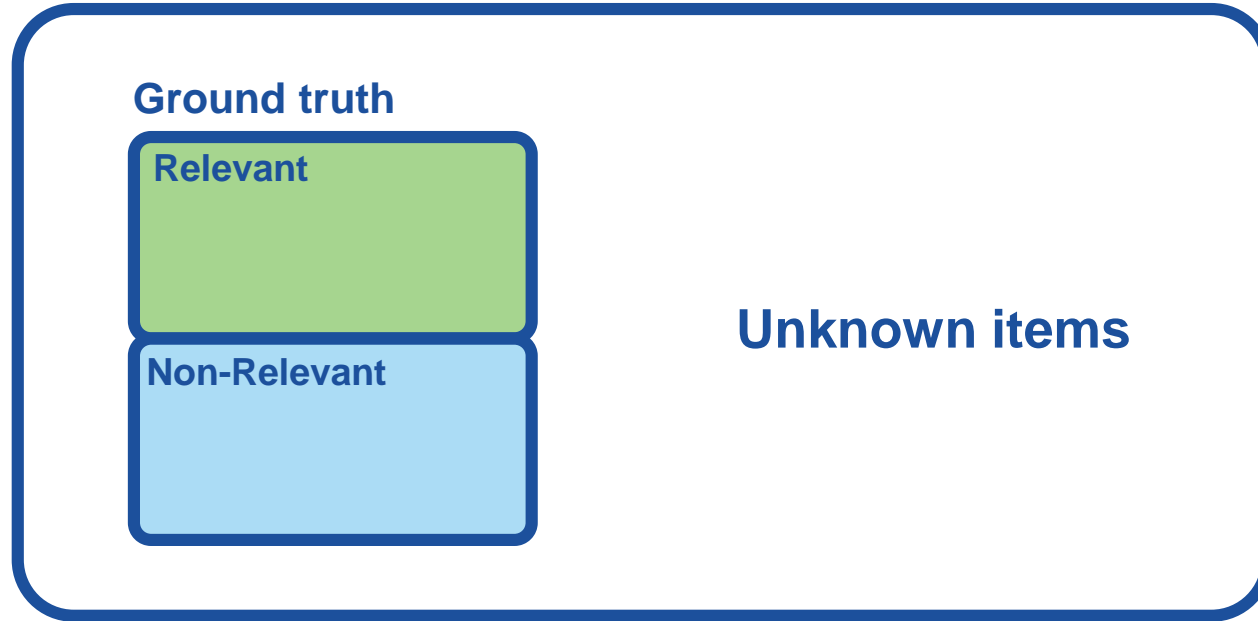
Off-line Evaluation: Dataset

Dataset



Off-line Evaluation: Dataset

Dataset



Off-line Evaluation: Partitioning

- $\text{Model} = f(\text{URM})$
- $\text{Estimated ratings} = g(\text{model}, \text{user profile})$

Off-line Evaluation: Partitioning

- $\text{Model} = f(\text{URM})$
- $\text{Estimated ratings} = g(\text{model}, \text{user profile})$

E.G. $\left\{ \begin{array}{l} \text{Model} = \text{Star Wars is similar to Avatar} \\ \text{User profile} = \text{Paolo Cremonesi likes Star Wars} \end{array} \right.$

Off-line Evaluation: Partitioning

- $\text{Model} = f(\text{URM})$
- $\text{Estimated ratings} = g(\text{model}, \text{user profile})$

E.G. $\left\{ \begin{array}{l} \text{Model} = \text{Star Wars is similar to Avatar} \\ \text{User profile} = \text{Paolo Cremonesi likes Star Wars} \end{array} \right.$

- $\text{Estimated ratings} \leftrightarrow \text{True recommendation}$

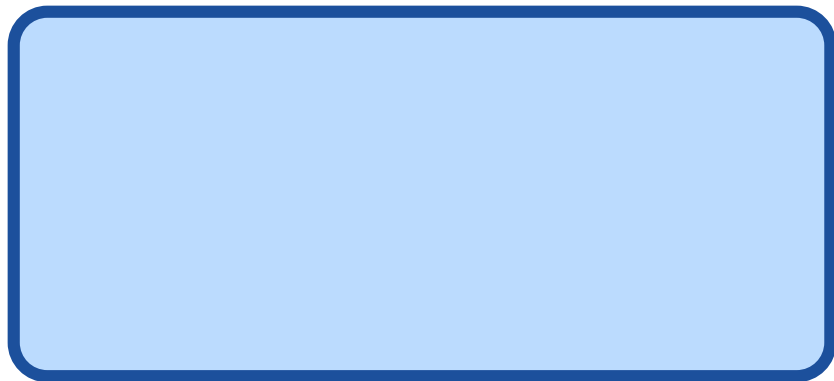
Partitioning: Hold Out of Ratings

- Model = $f(X)$
- Estimated ratings = $g(\text{model}, Y)$
- Estimated ratings $\leftrightarrow Z$

Partitioning: Hold Out

- Model = $f(X)$
- Estimated ratings = $g(\text{model}, Y)$
- Estimated ratings $\leftrightarrow Z$

X = training



Partitioning: Hold Out of Ratings

- Model = $f(X)$
- Estimated ratings = $g(\text{model}, Y)$
- Estimated ratings $\leftrightarrow Z$

X = training

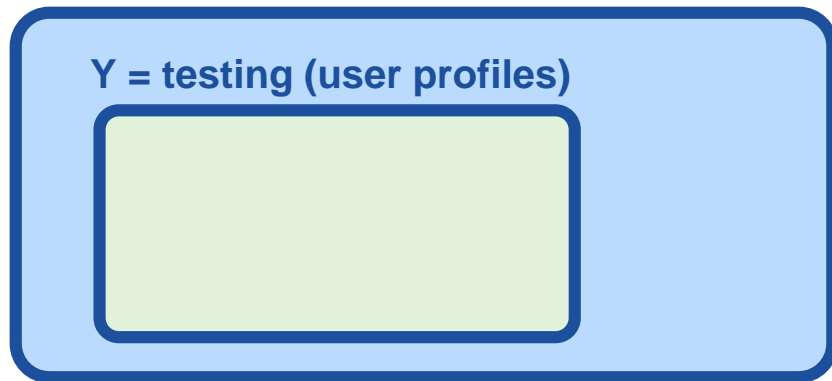
Y = testing (user profiles)



Partitioning: Hold Out of Ratings

- Model = $f(X)$
- Estimated ratings = $g(\text{model}, Y)$
- Estimated ratings $\leftrightarrow Z$

X = training



Z = testing (hidden ratings)



Partitioning: Hold Out of Users

X = training

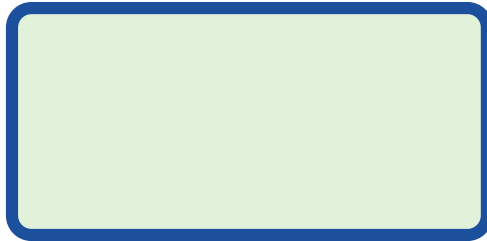


Partitioning: Hold Out of Users

X = training



Y = testing (user profiles)

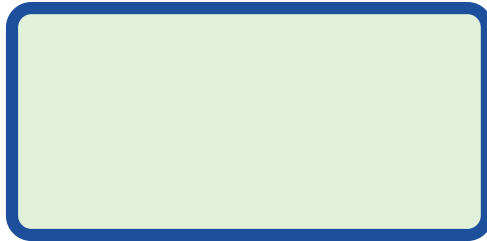


Partitioning: Hold Out of Users

X = training



Y = testing (user profiles)



Z = testing (hidden ratings)



Off-line Evaluation: Error Metrics



Off-line Evaluation: Error Metrics



Off-line Evaluation: Error Metrics



$$\text{Error: } e_{ui} = r_{ui} - \hat{r}_{ui}$$

Off-line Evaluation: Error Metrics



$$\text{Error: } e_{ui} = r_{ui} - \hat{r}_{ui} = 4 - 3.7 = 0.3$$

Off-line Evaluation: Error Metrics



$$\text{Error: } e_{ui} = r_{ui} - \hat{r}_{ui} = 4 - 3.7 = 0.3$$

\hat{r}_{ui} : rating estimated by the recommender system

r_{ui} : true rating in the test set

Off-line Evaluation: Error Metrics

Mean absolute error:

$$MAE = \frac{\sum_{u,i \in T} |r_{ui} - \hat{r}_{ui}|}{|T|}$$

Off-line Evaluation: Error Metrics

Mean absolute error:

$$MAE = \frac{\sum_{u,i \in T} |r_{ui} - \hat{r}_{ui}|}{|T|}$$

Mean square error:

$$MSE = \frac{\sum_{u,i \in T} (r_{ui} - \hat{r}_{ui})^2}{|T|}$$

Off-line Evaluation: Error Metrics

Mean absolute error:

$$MAE = \frac{\sum_{u,i \in T} |r_{ui} - \hat{r}_{ui}|}{|T|}$$

Mean square error:

$$MSE = \frac{\sum_{u,i \in T} (r_{ui} - \hat{r}_{ui})^2}{|T|}$$

T: test set

\hat{r}_{ui} : rating estimated by recommender system

r_{ui} : true rating in the test set

Off-line Evaluation: Error Metrics

Dataset

Ground truth

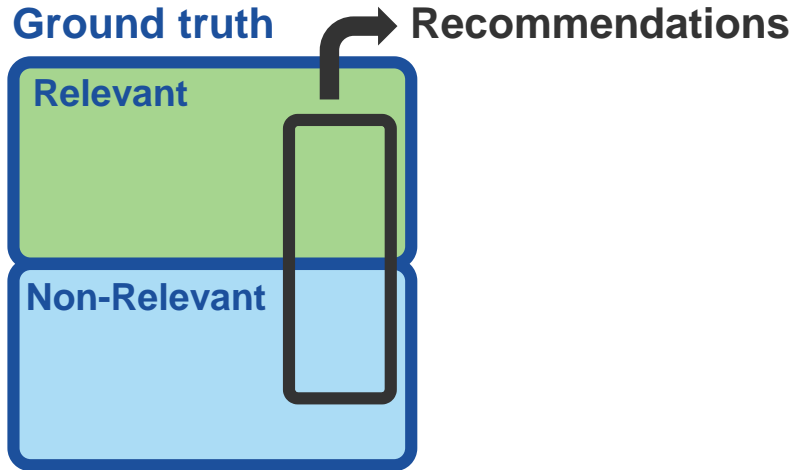
Relevant

Non-Relevant

Unknown ratings

Off-line Evaluation: Error Metrics

Dataset

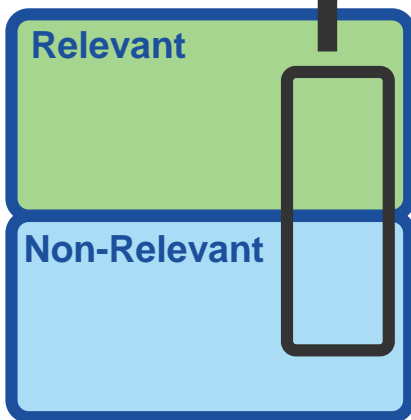


Unknown ratings

Off-line Evaluation: Error Metrics

Dataset

Ground truth



Recommendations

MISSING AS RANDOM ASSUMPTION (MAR):

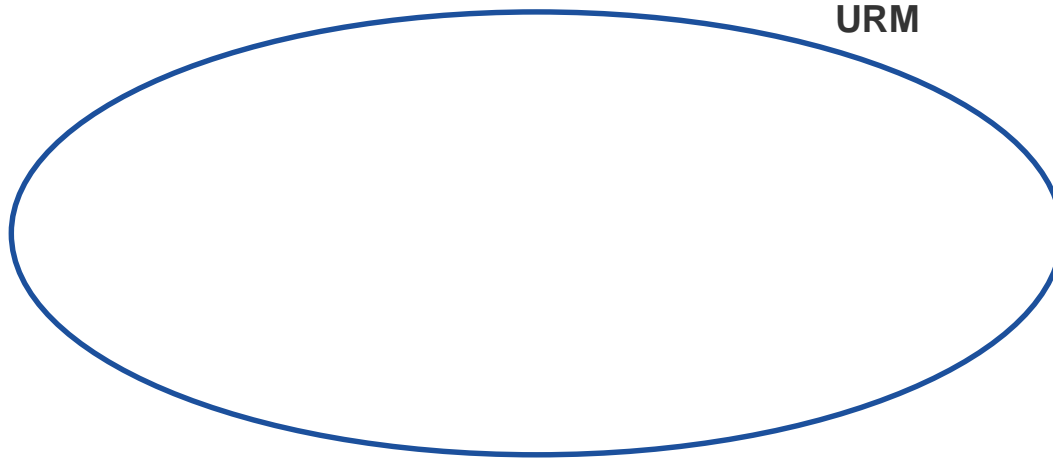
The distribution of missing ratings is equal to the distribution of available ratings

Why?

Because we measure errors only on available ratings, and we assume the error will be the same for unknown ratings

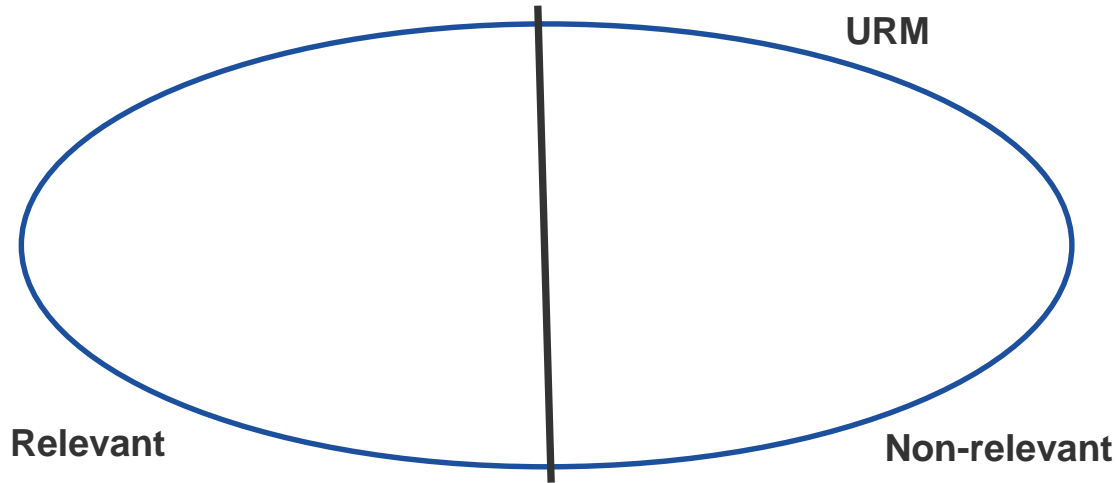
Unknown ratings

Off-line Evaluation: Classification Metrics

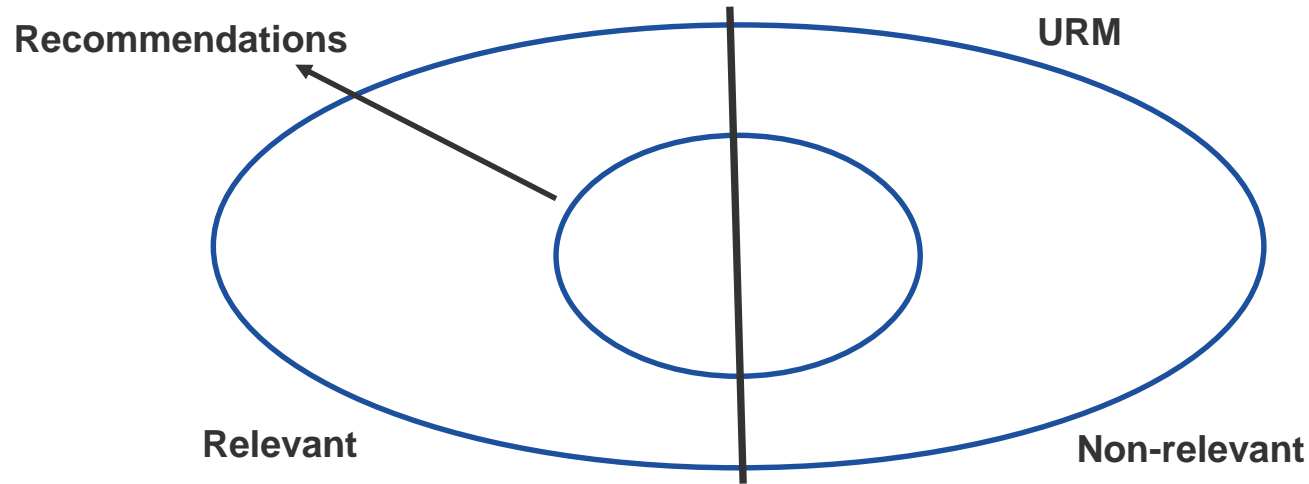


URM

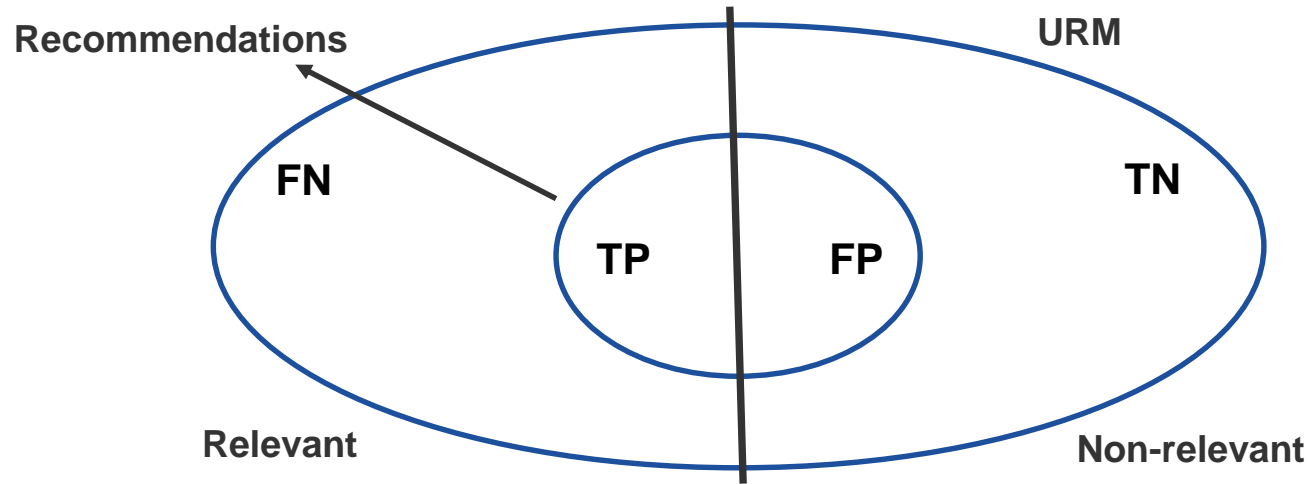
Off-line Evaluation: Classification Metrics



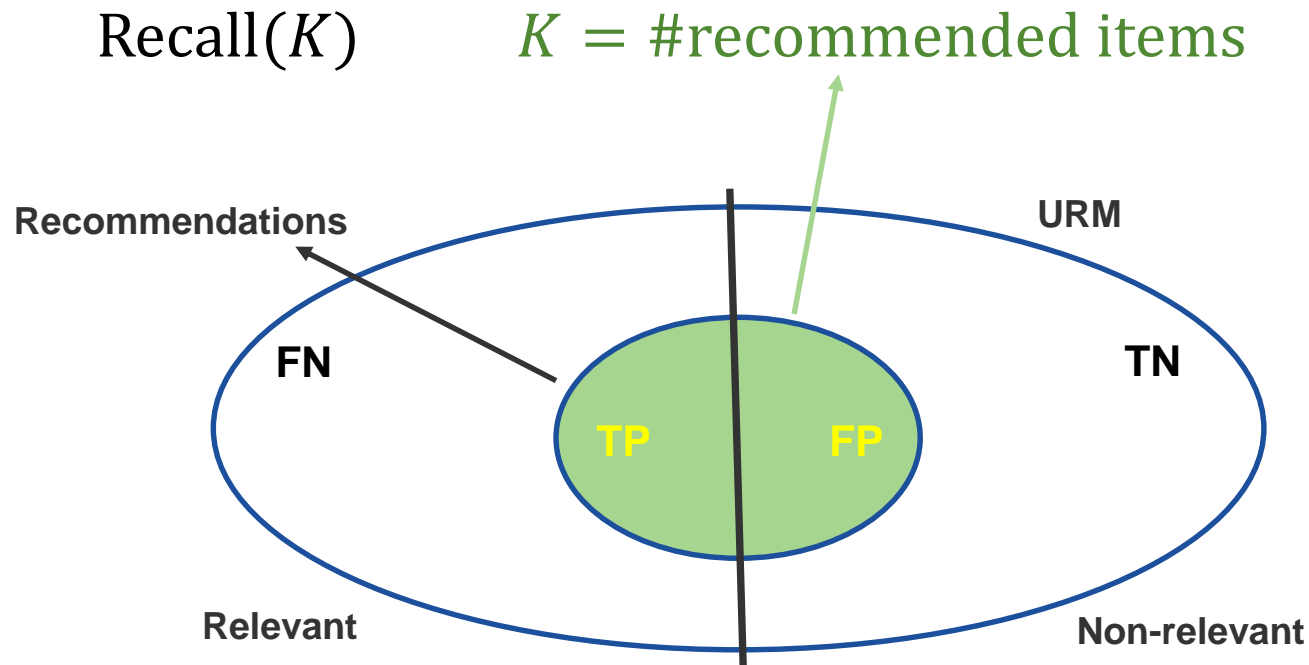
Off-line Evaluation: Classification Metrics



Off-line Evaluation: Classification Metrics

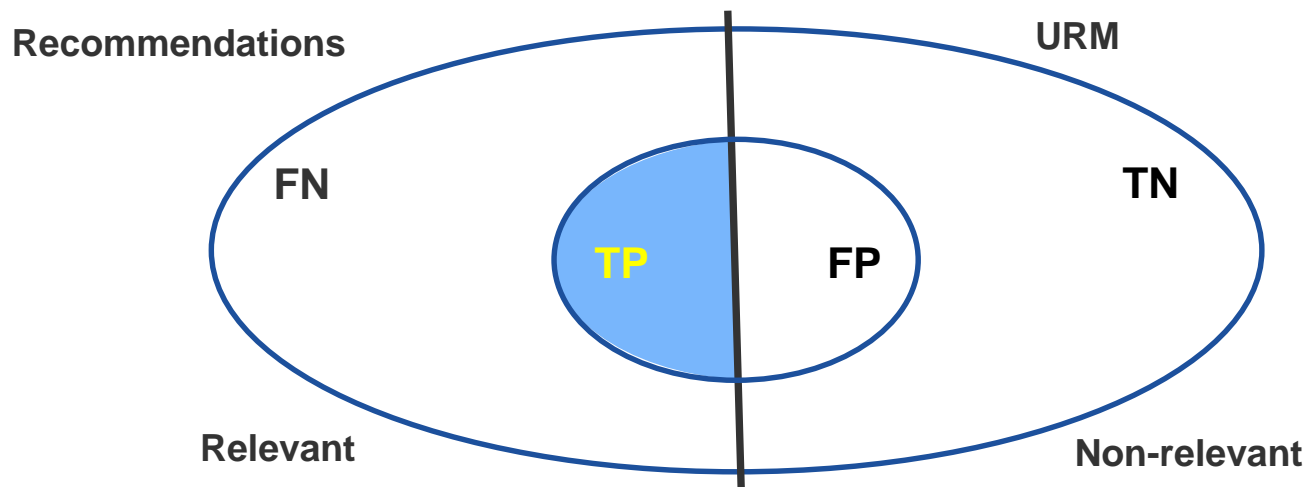


Off-line Evaluation: Classification Metrics



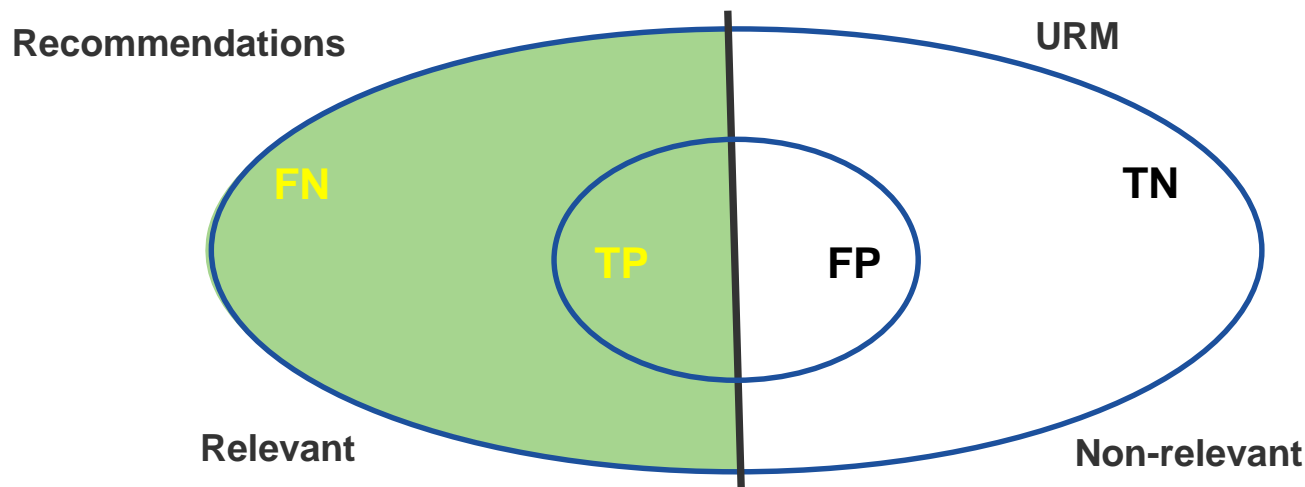
Off-line Evaluation: Classification Metrics

$$\text{Recall}(K) = \frac{\text{\#relevant recommended items}}{\text{\#relevant items}}$$



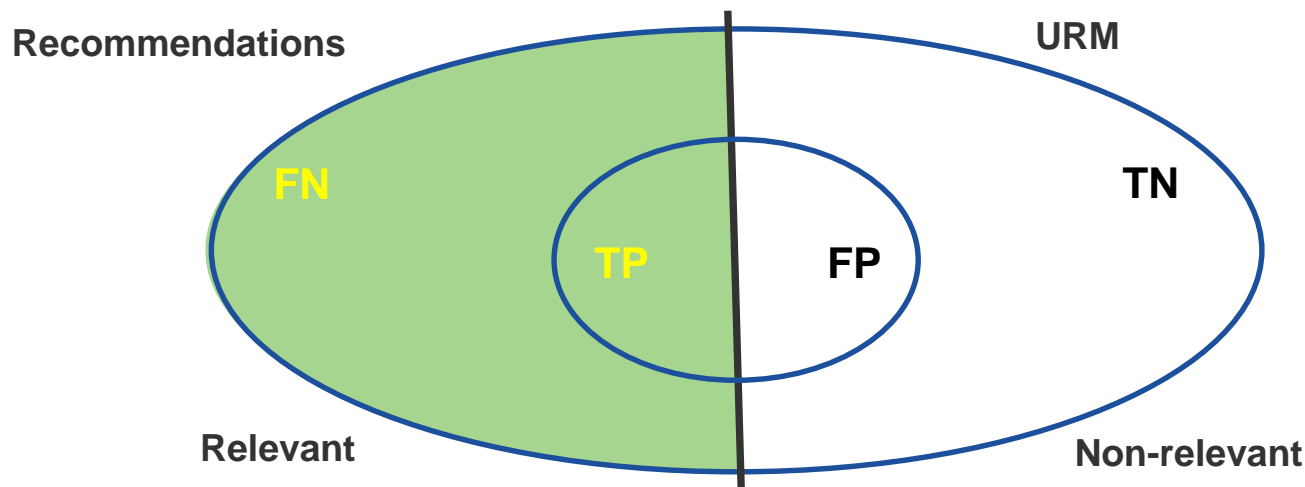
Off-line Evaluation: Classification Metrics

$$\text{Recall}(K) = \frac{\text{\#relevant recommended items}}{\text{\#tested relevant items}}$$



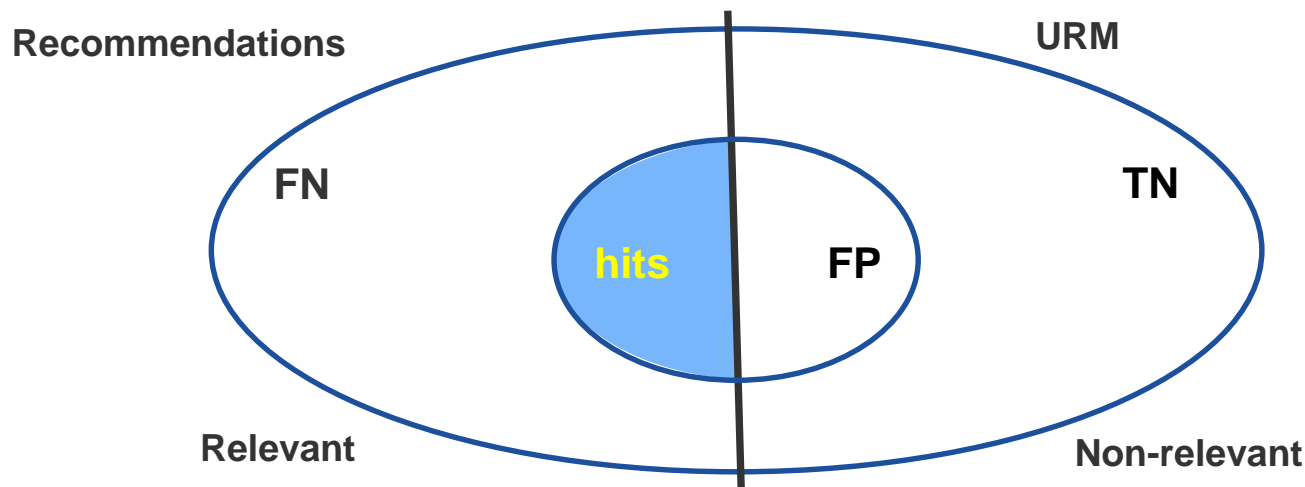
Off-line Evaluation: Classification Metrics

$$\text{Recall}(K) = \frac{\text{\#relevant recommended items}}{\text{\#tested relevant items}} = \frac{TP}{FN + TP}$$



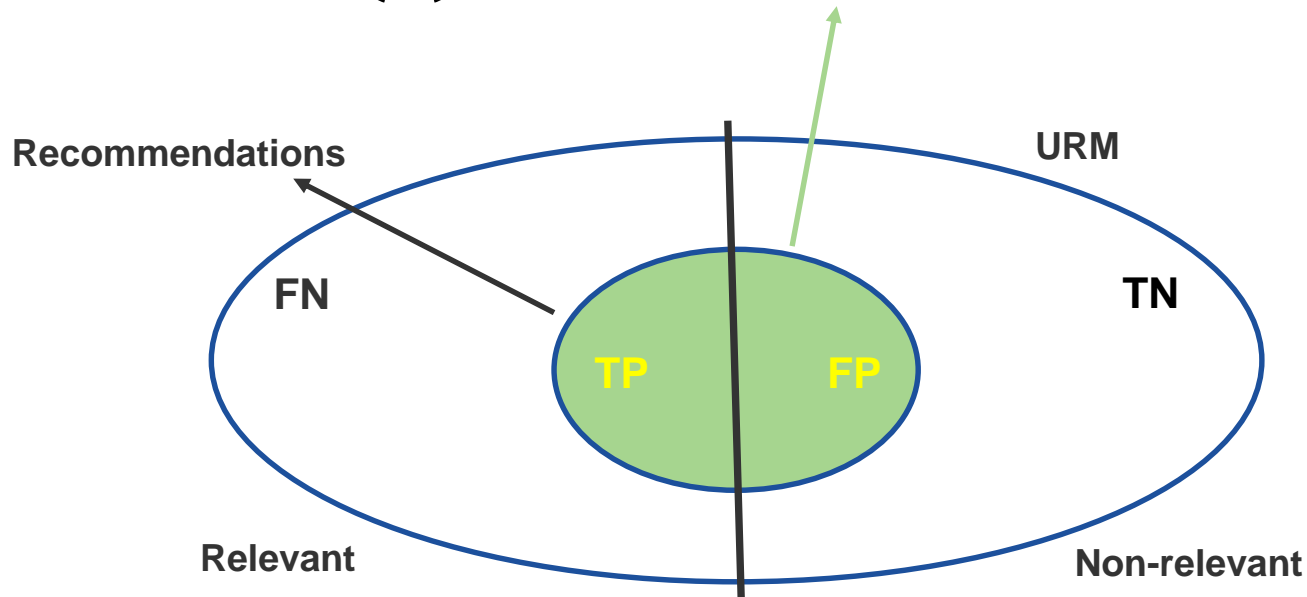
Off-line Evaluation: Classification Metrics

$$\text{Recall}(K) = \frac{\text{\#relevant recommended items}}{\text{\#tested relevant items}} = \frac{\text{\#hits}}{\text{FN} + \text{TP}}$$



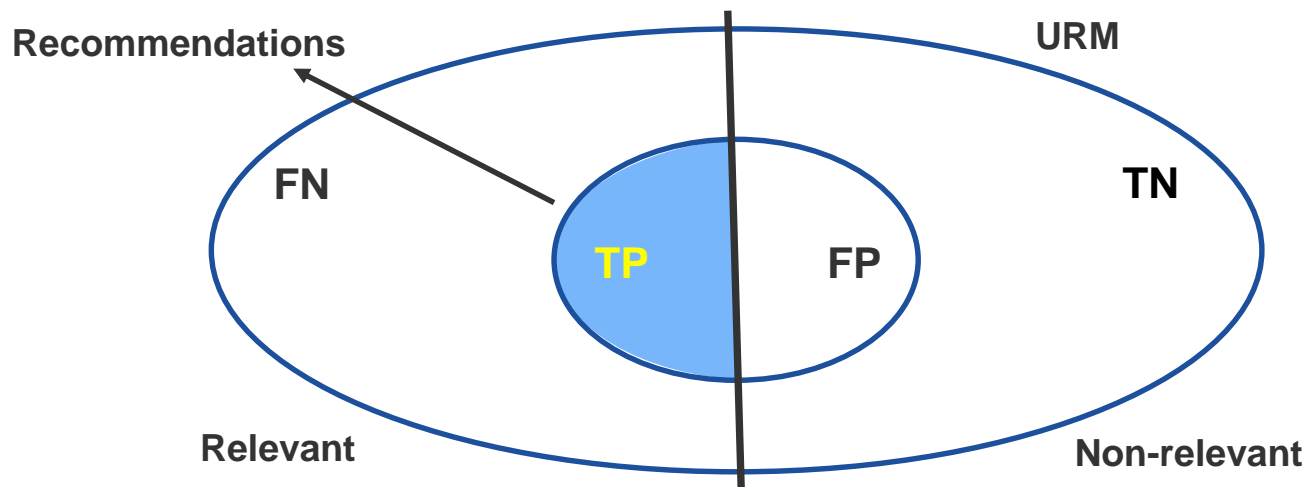
Off-line Evaluation: Classification Metrics

$$\text{Precision}(K) = \frac{TP}{TP + FP} \quad K = \# \text{recommended items}$$



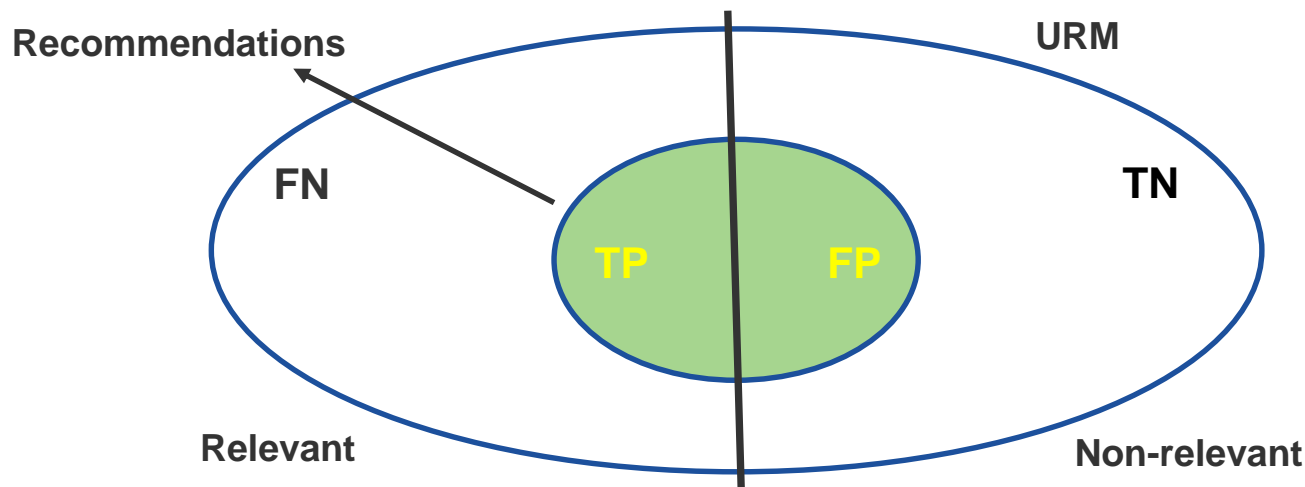
Off-line Evaluation: Classification Metrics

$$\text{Precision}(K) = \frac{\text{\#relevant recommended items}}{\text{\#recommended items}}$$



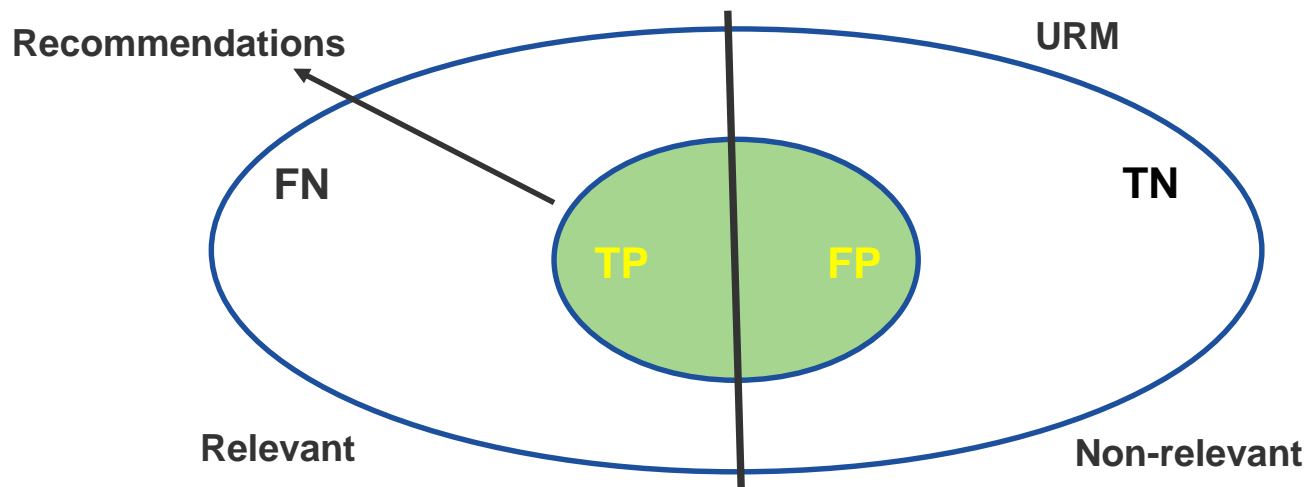
Off-line Evaluation: Classification Metrics

$$\text{Precision}(K) = \frac{\text{\#relevant recommended items}}{\text{\#all recommended items}}$$



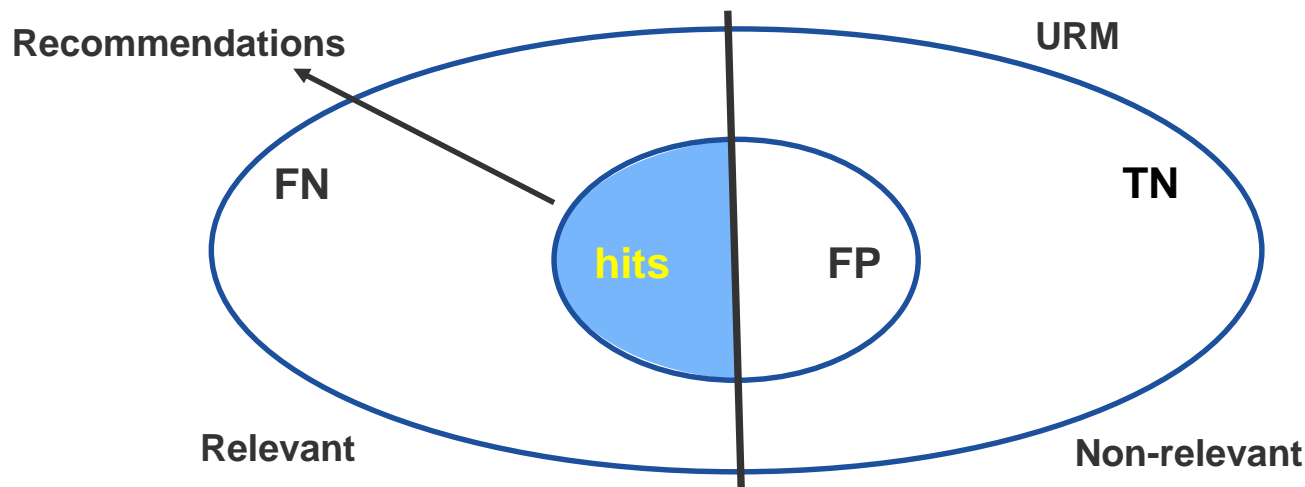
Off-line Evaluation: Classification Metrics

$$\text{Precision}(K) = \frac{\text{\#relevant recommended items}}{\text{\#all recommended items}} = \frac{TP}{FP + TP}$$



Off-line Evaluation: Classification Metrics

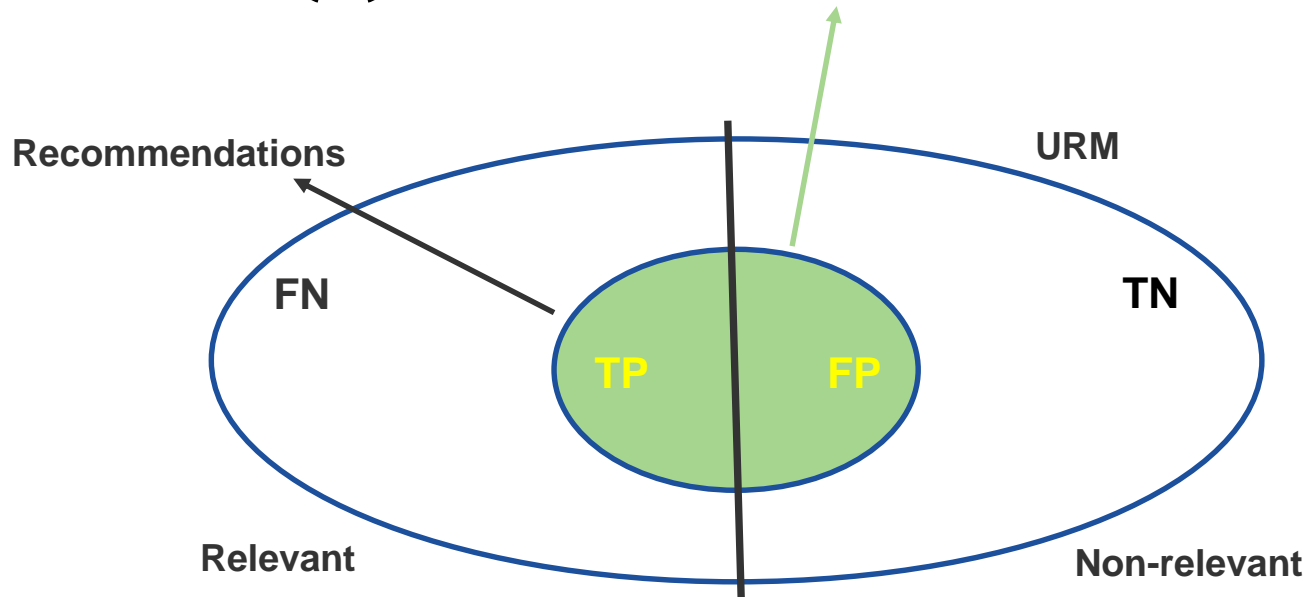
$$\text{Precision}(K) = \frac{\text{\#relevant recommended items}}{\text{\#all recommended items}} = \frac{\text{\#hits}}{\text{FP} + \text{TP}}$$



Off-line Evaluation: Classification Metrics

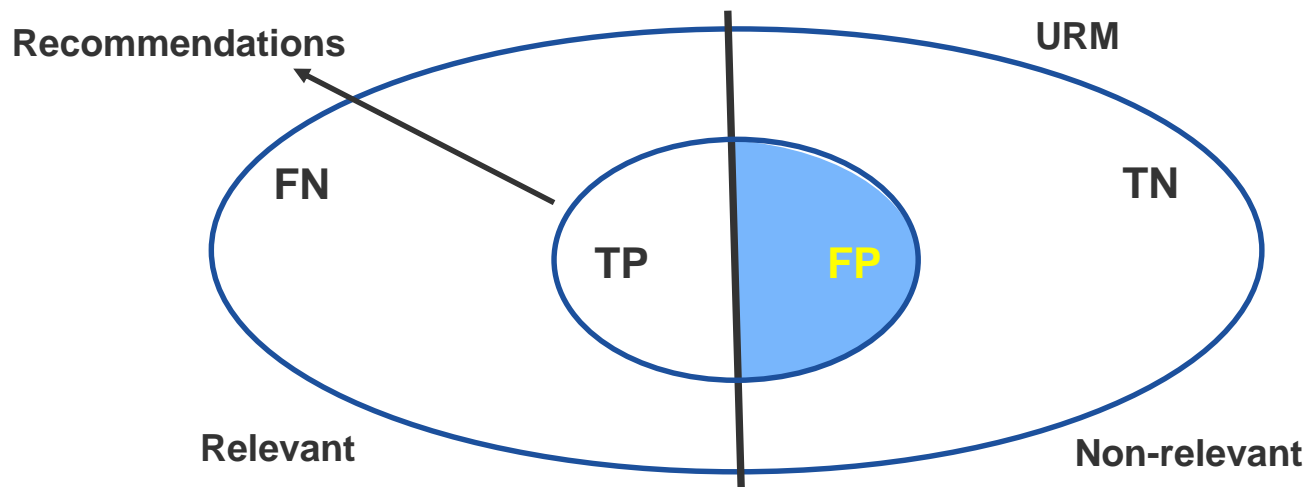
Fallout(K) =

K = #recommended items



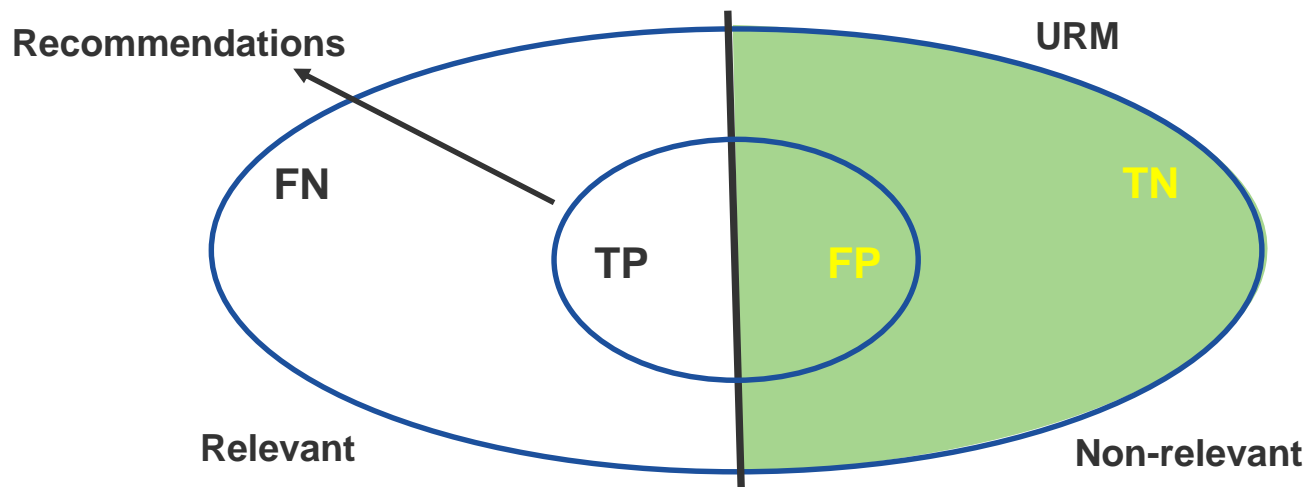
Off-line Evaluation: Classification Metrics

$$\text{Fallout}(K) = \frac{\text{\#non relevant recommended items}}{\text{\#non relevant items}}$$



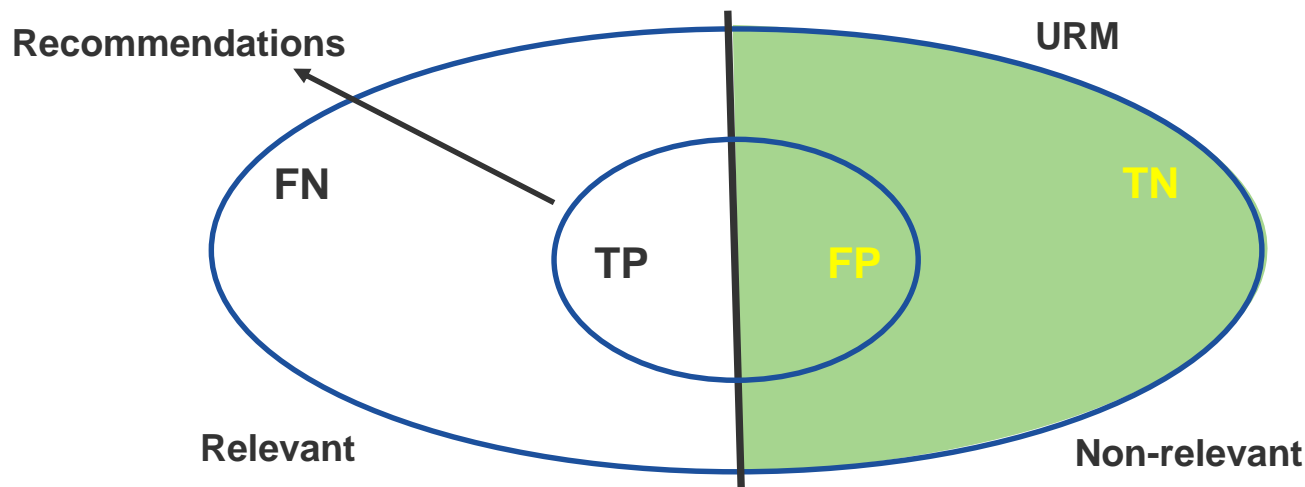
Off-line Evaluation: Classification Metrics

$$\text{Fallout}(K) = \frac{\text{\#non relevant recommended items}}{\text{\#all non relevant items}}$$



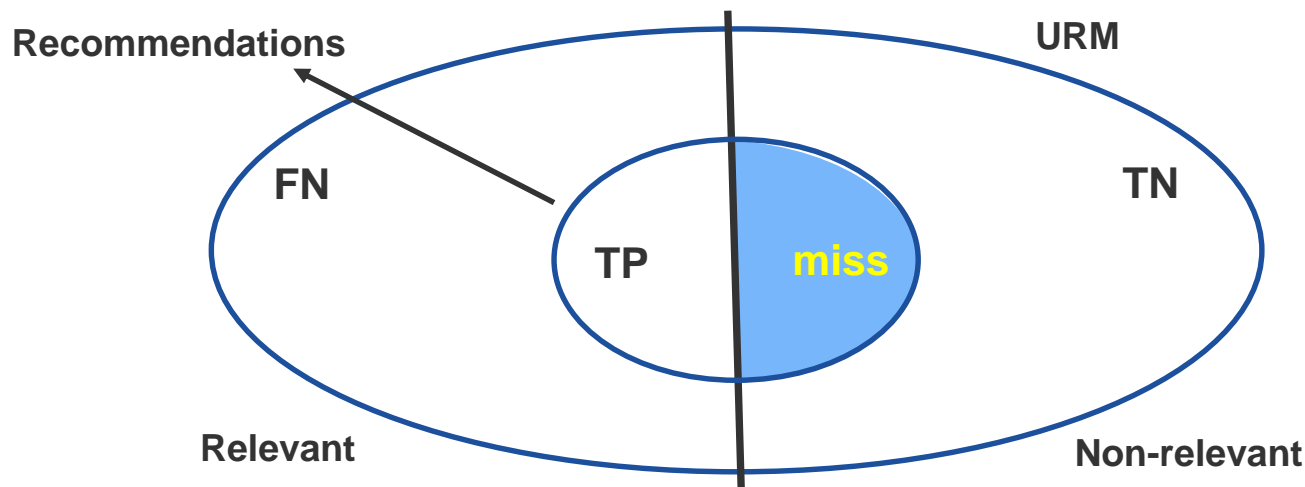
Off-line Evaluation: Classification Metrics

$$\text{Fallout}(K) = \frac{\text{\#non relevant recommended items}}{\text{\#all non relevant items}} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$



Off-line Evaluation: Classification Metrics

$$\text{Fallout}(K) = \frac{\text{\#non relevant recommended items}}{\text{\#all non relevant items}} = \frac{\text{\#miss}}{\text{FP} + \text{TN}}$$



Off-line Evaluation: Classification Metrics

Dataset

Ground truth

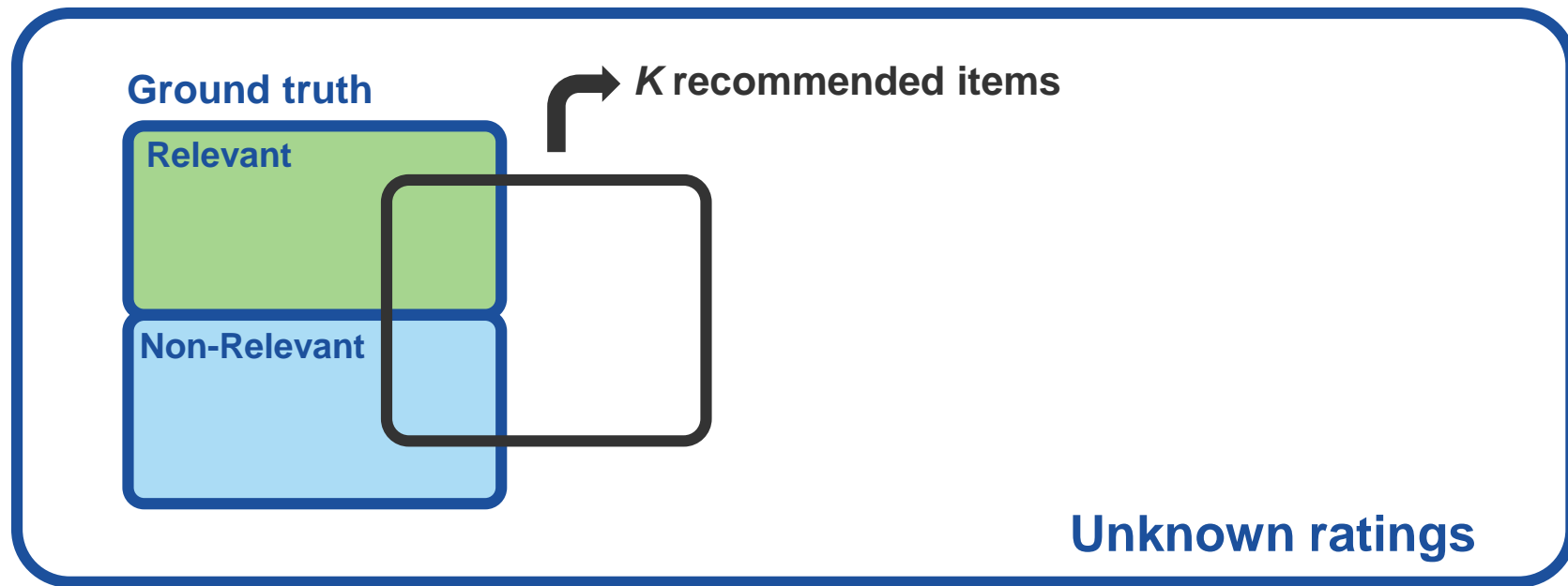
Relevant

Non-Relevant

Unknown ratings

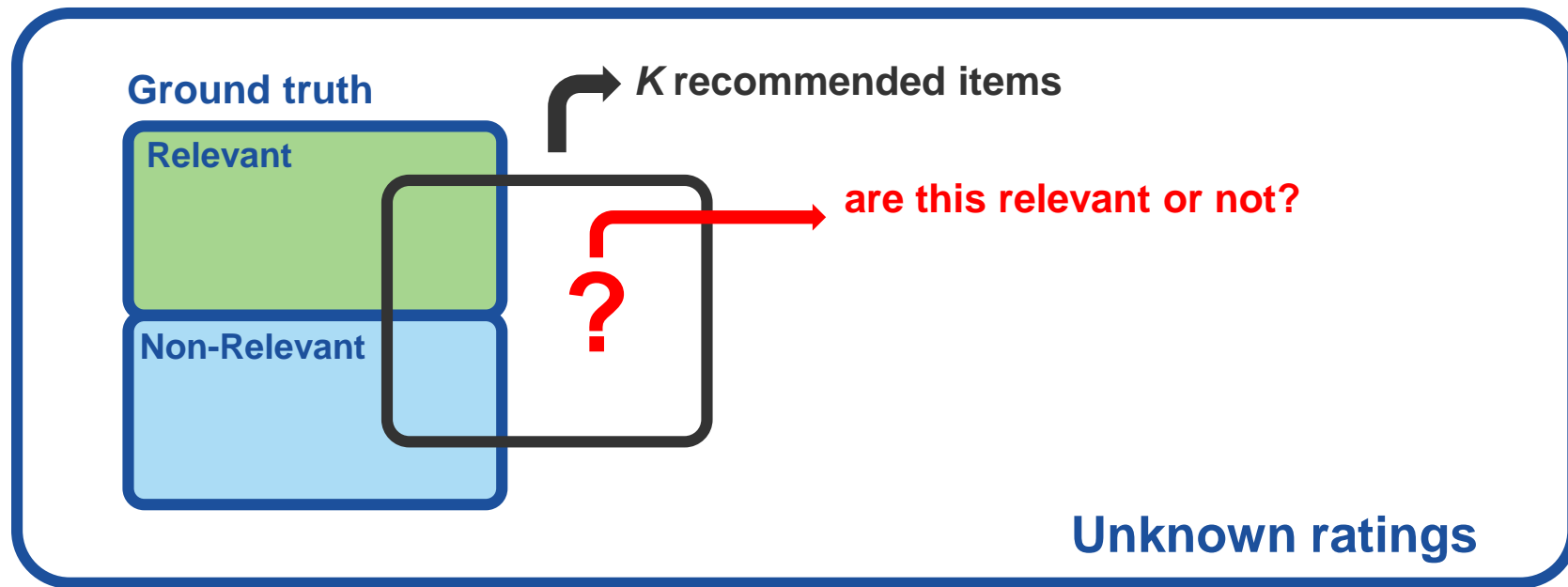
Off-line Evaluation: Classification Metrics

Dataset



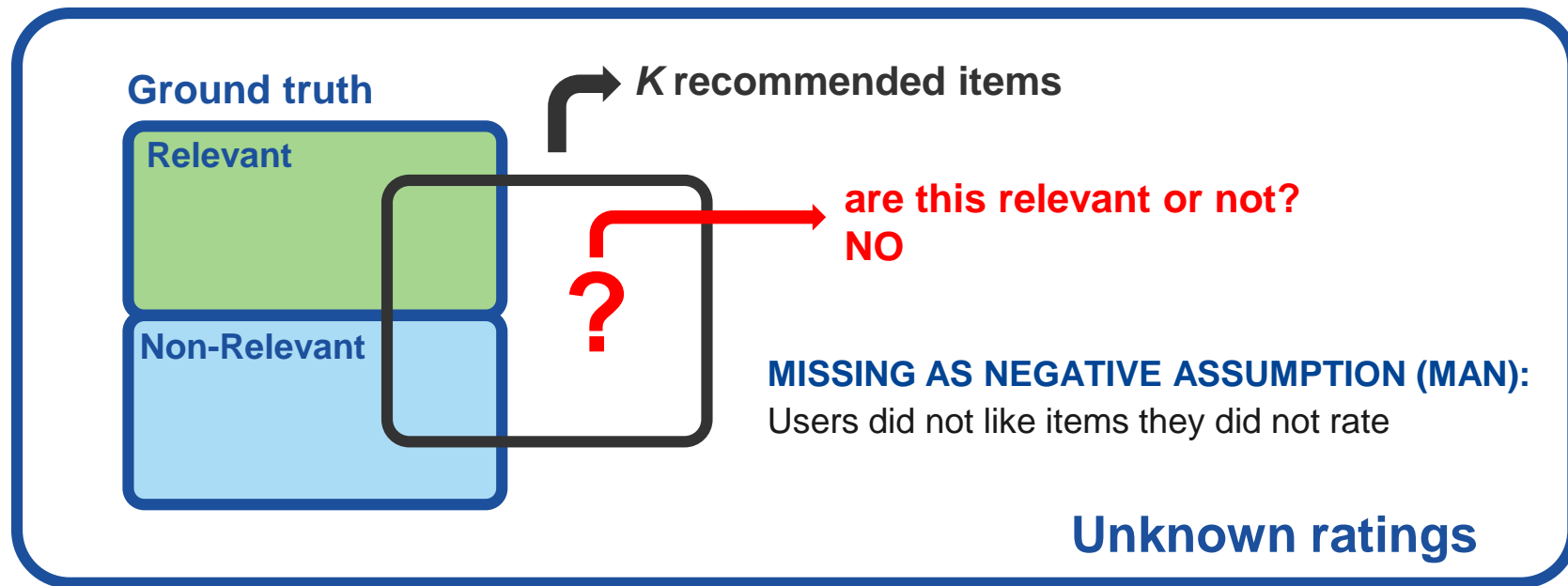
Off-line Evaluation: Classification Metrics

Dataset



Off-line Evaluation: Classification Metrics

Dataset



All-Missing-As-Negative (AMAN) hypothesis

- all missing ratings are irrelevant
- underestimate the true precision computed on the (unknown) complete data

Harald Steck, *Training and testing of RSs on data missing not at random*. In KDD '10

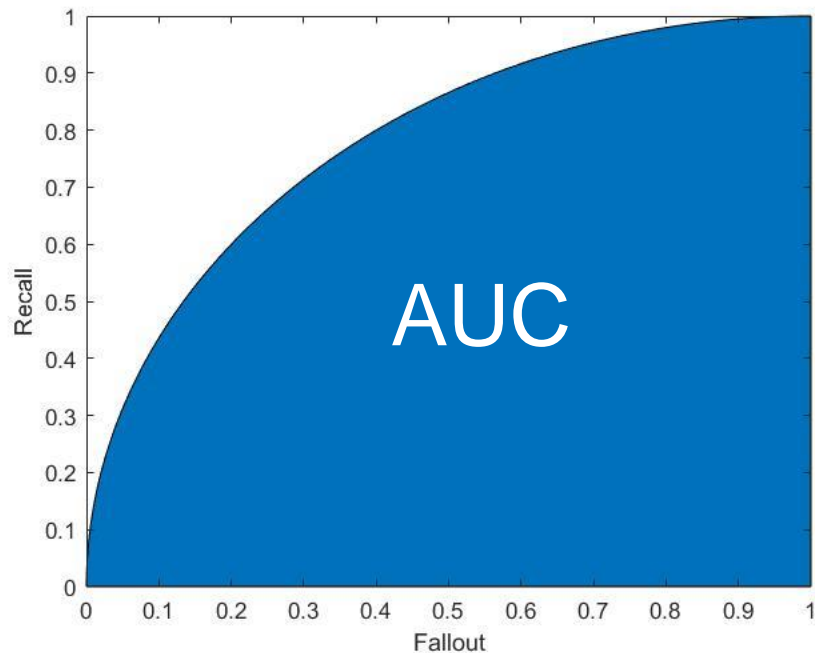
Missing-Not-At-Random (MNAR) hypothesis

- non-relevant ratings are missing with a higher probability than relevant ratings
- (nearly) unbiased estimate of recall on the (unknown) complete data
- much milder than assuming that
 - all the ratings are missing at random (MAE and RMSE)
 - all missing ratings are irrelevant (Precision)

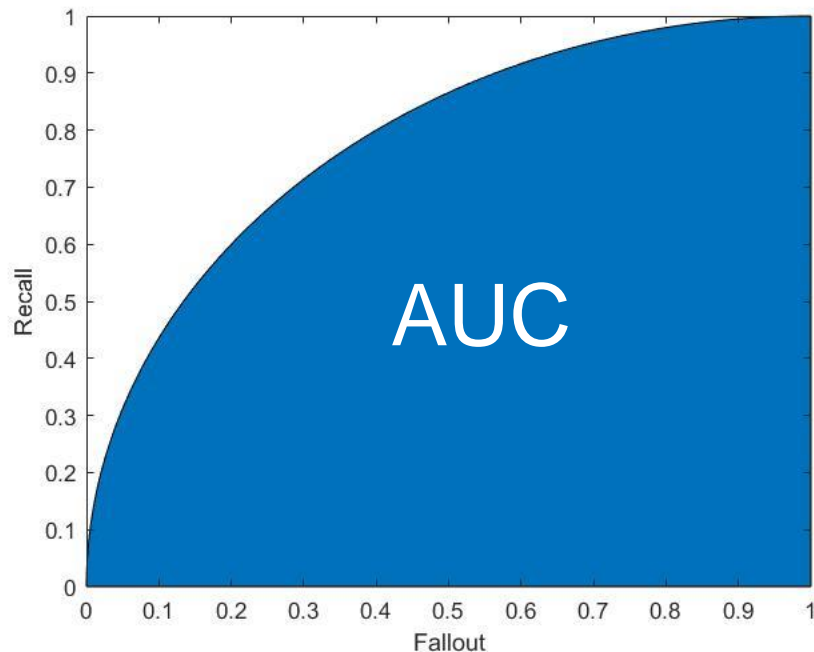
Ranking Metrics



ROC curve (area under curve)

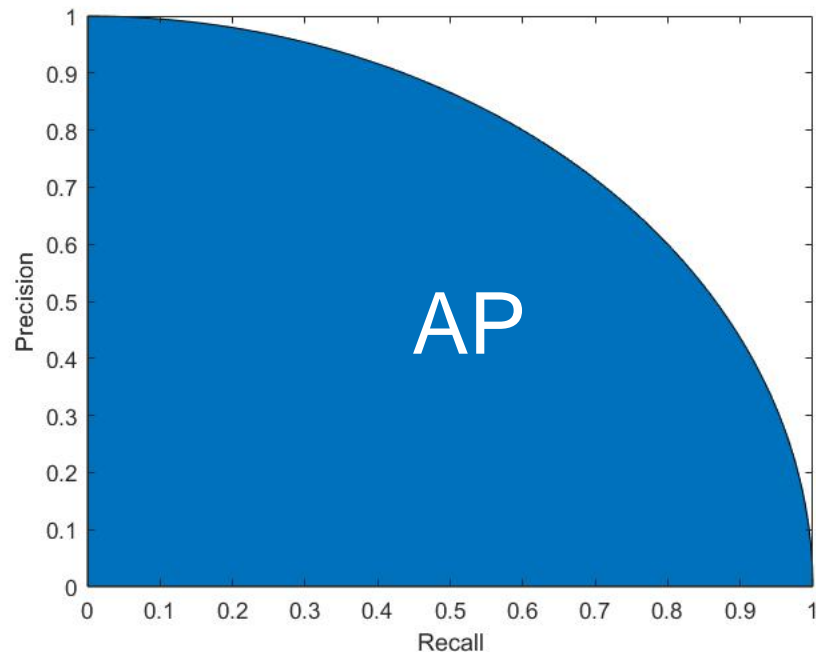


ROC curve (area under curve)

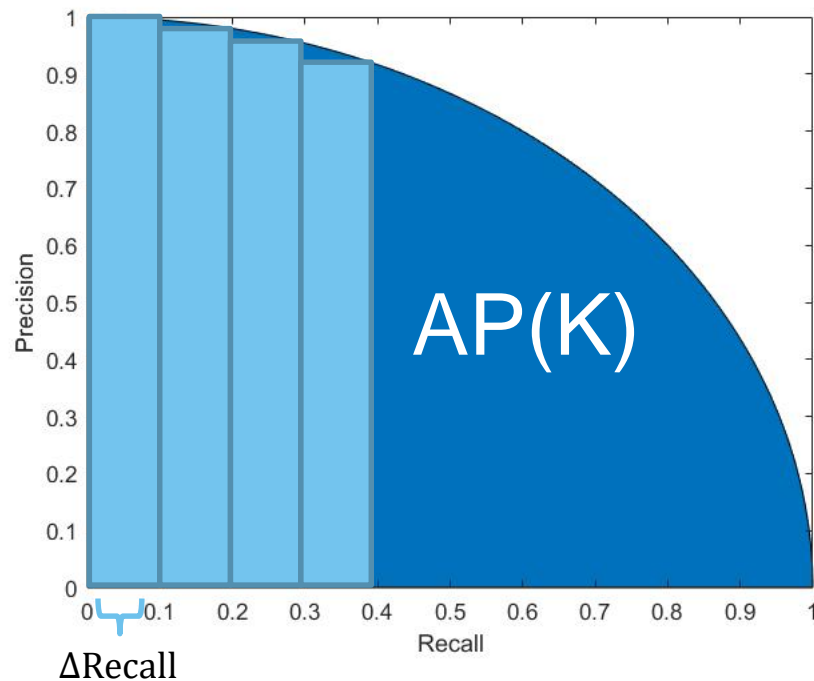


$$AUC = \sum_k \text{Recall}(k) * \Delta \text{Fallout}$$

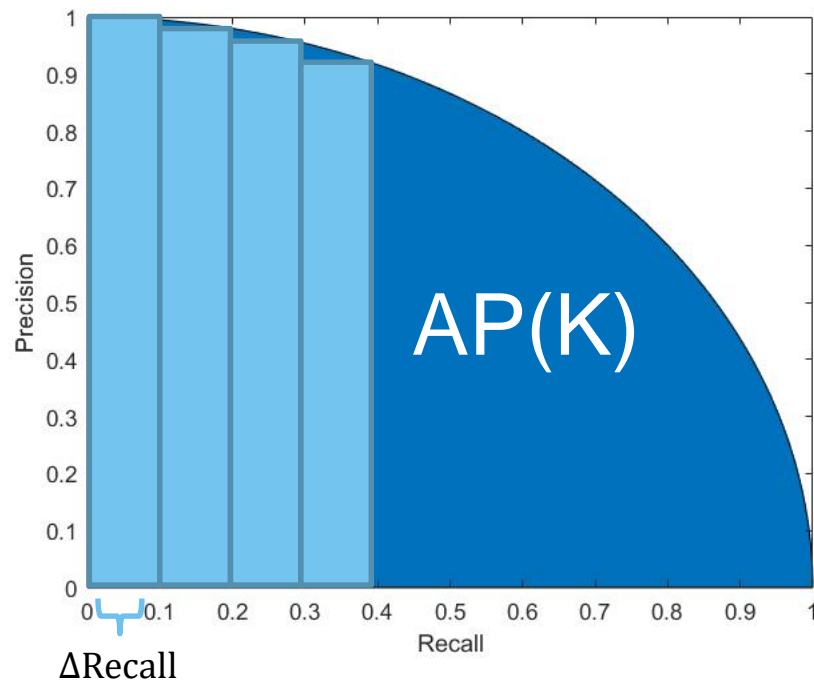
Mean Average Precision



Mean Average Precision

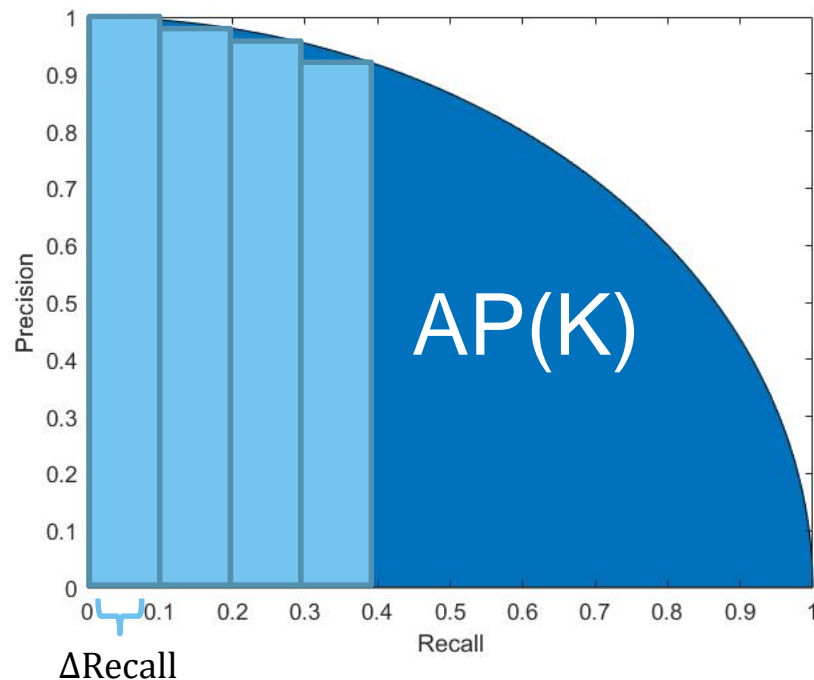


Mean Average Precision



$$AP = \sum_k \text{Precision}(k) * \Delta \text{Recall}$$

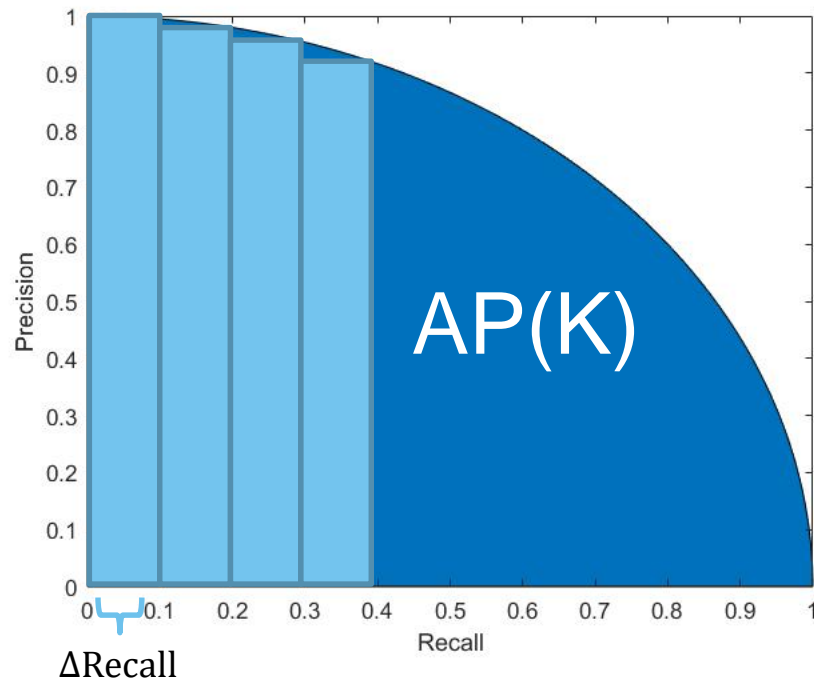
Mean Average Precision



$$AP = \sum_k \text{Precision}(k) * \Delta \text{Recall}$$

$$\Delta \text{Recall} = \text{Recall}(k) - \text{Recall}(k - 1)$$

Mean Average Precision



$$AP = \sum_k \text{Precision}(k) * \Delta \text{Recall}$$

$$\text{MAP} = \frac{\sum_u AP_u}{\# \text{users}}$$

Average Reciprocal Hit-Rank

Weighted version of recall

$$ARHR = \frac{\sum_i \frac{1}{\text{rank}(i)}}{\text{\#tested relevant items}}$$

Average Reciprocal Hit-Rank

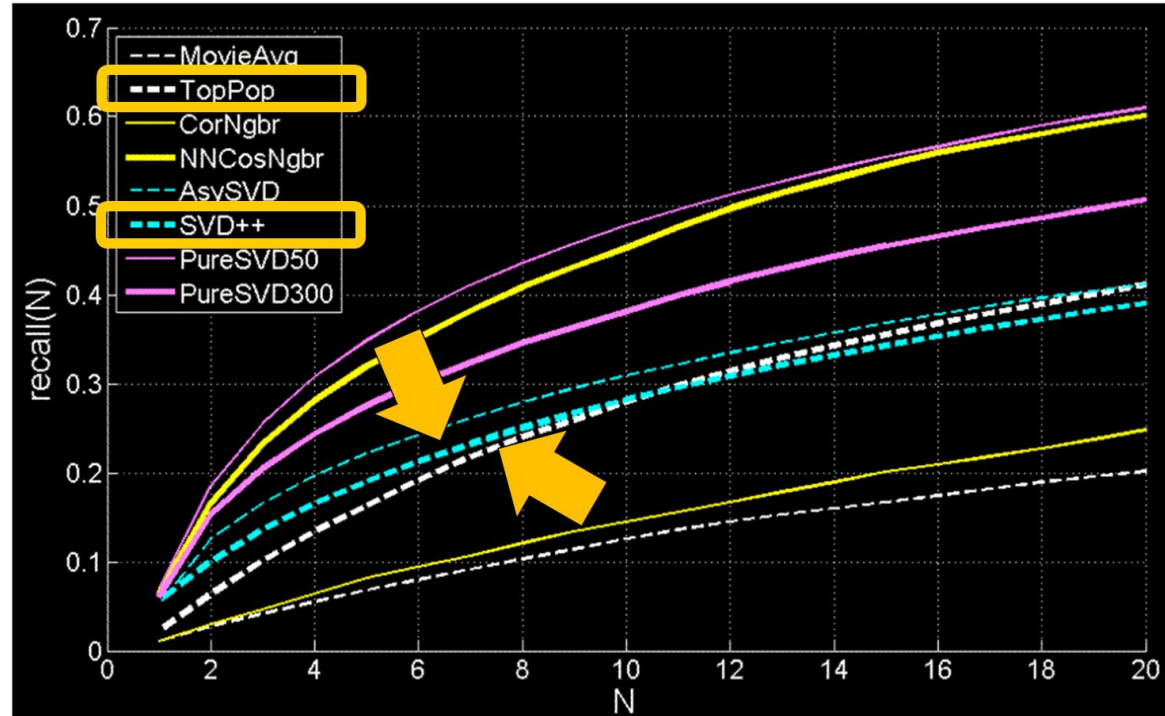
Weighted version of recall

$$ARHR = \frac{\sum_i \frac{1}{\text{rank}(i)}}{\text{\#tested relevant items}}$$

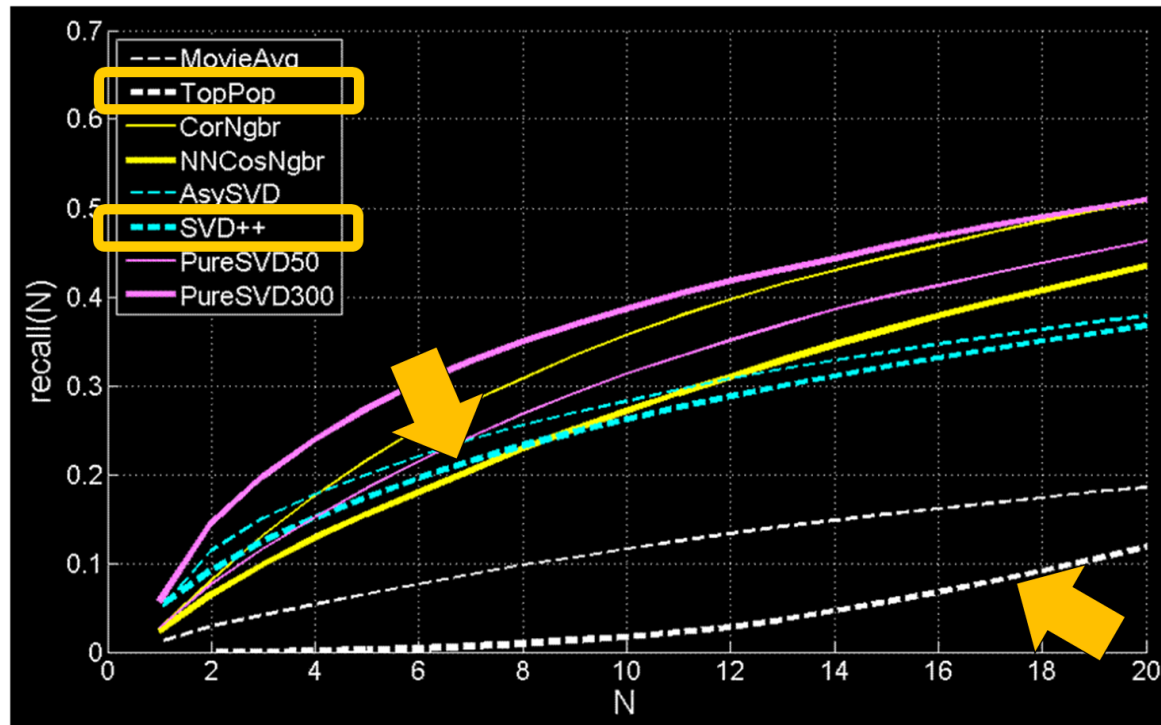
i: **relevant** item recommended to the user

rank(i): position of item *i* in the list of recommendations

Netflix dataset: recall



Netflix dataset: recall on long tail



Evaluation: are we really making much progress?

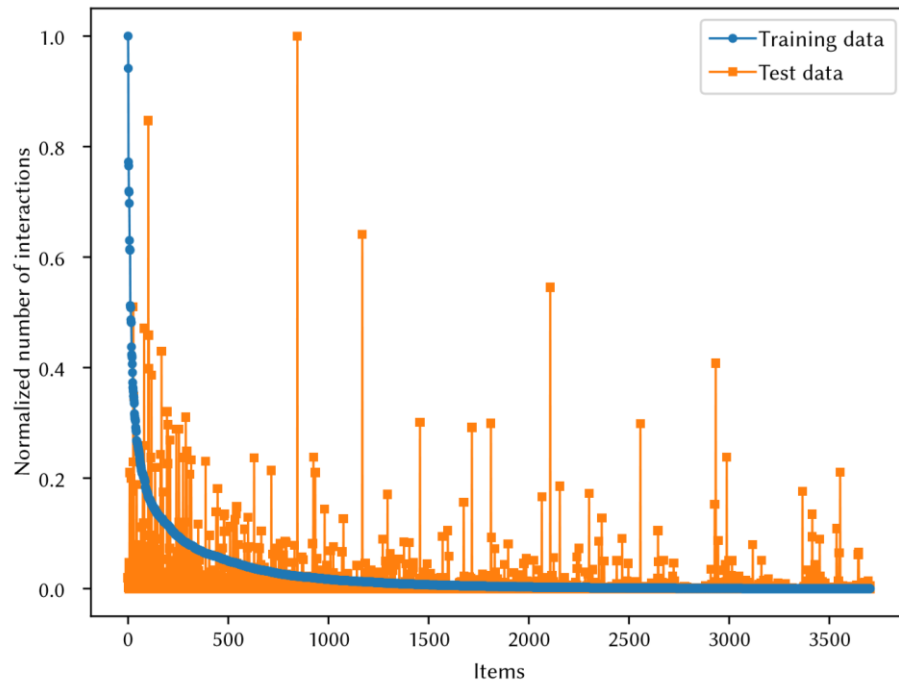
Conference	Rep. / Non-rep.	Reproducible
KDD	3/4 (75%)	[17], [23], [48]
RecSys	1/7 (14%)	[53]
SIGIR	1/3 (30%)	[10]
WWW	2/4 (50%)	[14], [24]
Total	7/18 (39%)	

Non-reproducible: KDD: [43], RecSys: [41], [6], [38], [44], [21], [45], SIGIR: [32], [7], WWW: [42], [11]

Evaluation: are we really making much progress?

	CiteULike-a			
	HR@5	NDCG@5	HR@10	NDCG@10
TopPopular	0.1803	0.1220	0.2783	0.1535
UserKNN	0.8213	0.7033	0.8935	0.7268
ItemKNN	0.8116	0.6939	0.8878	0.7187
$P^3\alpha$	0.8202	0.7061	0.8901	0.7289
$RP^3\beta$	0.8226	0.7114	0.8941	0.7347
CMN	0.8069	0.6666	0.8910	0.6942

Evaluation: are we really making much progress?



Evaluating Non Accuracy Metrics

ESSIR 2019 – Recommender Systems – Paolo Cremonesi



POLITECNICO
MILANO 1863

Measuring Diversity

$$\text{diversity} = \frac{\sum_{i,j} 1 - \text{similarity}(i, j)}{N(N - 1)}$$

Measuring Diversity

$$\text{diversity} = \frac{\sum_{i,j} 1 - \text{similarity}(i, j)}{N(N - 1)}$$

N: total number of items

i, j : considered items

Measuring Novelty

$$\text{Novelty} = \frac{\text{\#relevant and **unknown** recommended items}}{\text{\#recommended relevant items}}$$

Measuring Novelty

$$\text{Novelty} = \frac{\text{\#relevant and **unknown** recommended items}}{\text{\#recommended relevant items}}$$

$$\text{novelty} \approx 1/\text{popularity}$$

Measuring Novelty

$$\text{Novelty} = \frac{\text{\#relevant and **unknown** recommended items}}{\text{\#recommended relevant items}}$$

$$\text{novelty} \approx 1/\text{popularity}$$

$$\text{novelty} = \frac{\sum_{i \in \text{hits}} \log_2 \left(\frac{1}{\text{popularity}(i)} \right)}{\text{\#hits}}$$

Measuring Novelty

$$\text{Novelty} = \frac{\text{\#relevant and **unknown** recommended items}}{\text{\#recommended relevant items}}$$

$$\text{novelty} \approx 1/\text{popularity}$$

$$\text{novelty} = \frac{\sum_{i \in \text{hits}} \log_2 \left(\frac{1}{\text{popularity}(i)} \right)}{\text{\#hits}}$$

$\text{popularity}(i) =$
% of users who rated item i

Paolo Cremonesi
paolo.cremonesi@polimi.it

