



# RECOMMENDATION ENGINE

## DEMYSTIFIED

NEIGHBORHOOD METHODS  
COLLABORATIVE FILTERING

Alex Lin

Senior Architect

Intelligent Mining

# Outline

- Introduction
- User-oriented Collaborative Filtering
- Item-oriented Collaborative Filtering
- Challenges
- Best Practices



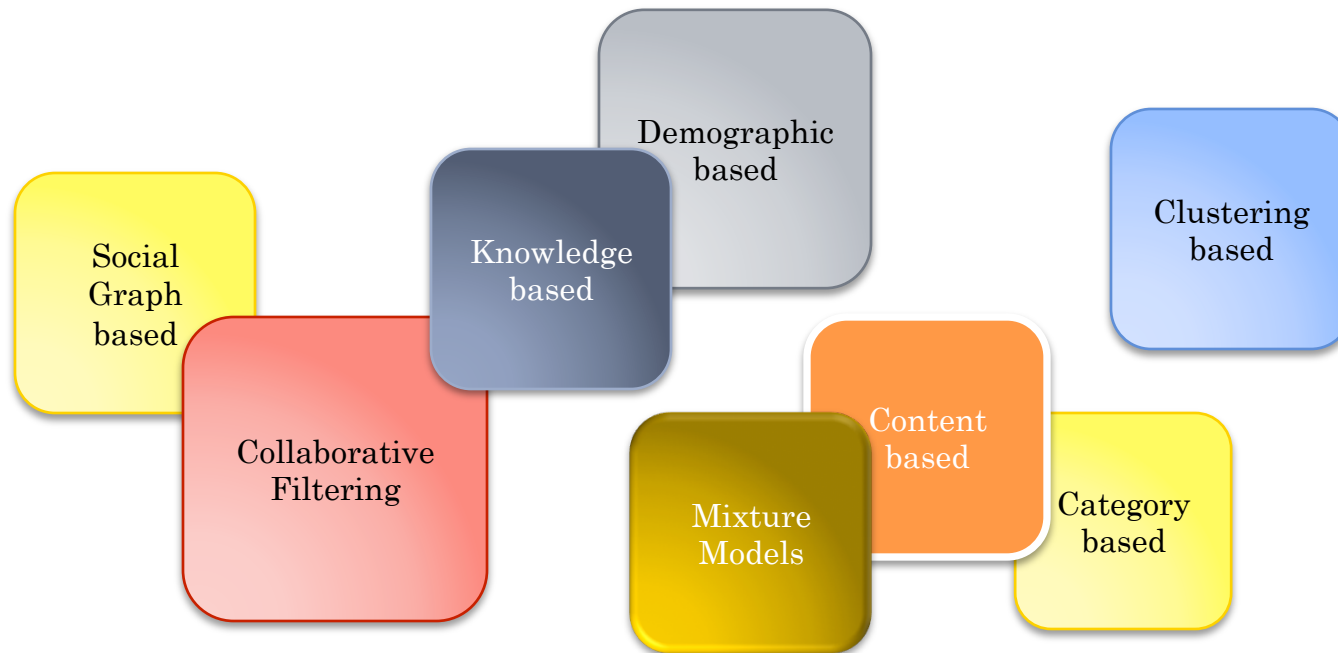
# Recommendation Engine

What is a Recommendation Engine (RE)?

- RE takes “observation” data and uses machine learning / statistical algorithms to predict outcomes or levels of interest.
- “Recommender systems form a specific type of information filtering (IF) technique that attempts to present information items (movies, music, books, news, images, web pages, etc.) that are likely of interest to the user.” – Wikipedia



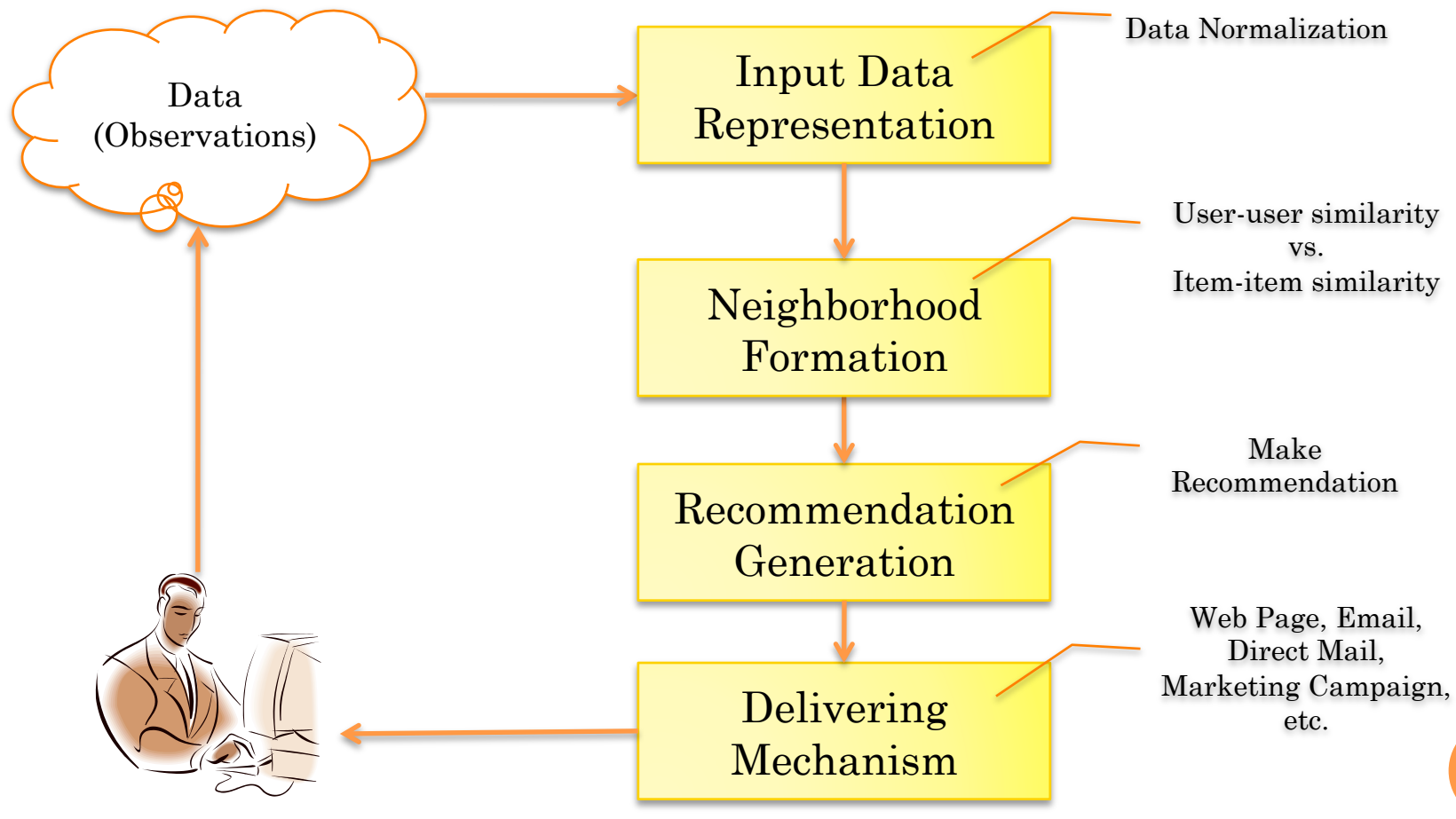
# Recommendation Engine



- This presentation will focus on Neighborhood-based Collaborative Filtering
  - User-oriented method
  - Item-oriented method



# Neighborhood-based Collaborative Filtering



# Outline

- Introduction
- User-oriented Collaborative Filtering
- Item-oriented Collaborative Filtering
- Challenges
- Best Practices



# User-oriented CF

- Input Data Representation: Users / Items matrix.

		users											
		1	2	3	4	5	6	7	8	9	10	...	n
items	1	1		1			1				1		
	2							1	1	1			
	3	1	1		1				1	1			
	4		1			1			1	1			
	⋮			1				1					
	m												

- Cell value “1” means user purchased the item.
- Data Normalization is not shown on this slide.



# User-oriented CF

- Neighborhood Formation:
  - Find the k most like-minded users in the system.

		users											
		1	2	3	4	5	6	7	<u>8</u>	9	10	...	n
items	1	1	1	1			1				1		
	2							1	1	1			
	3	1	1		1				1	1			1
	4		1			1			1	1			
	⋮			1				1					
	m				1					1			





# User-oriented CF

- Neighborhood Formation:
  - Find the  $k$  most like-minded users in the system.

		users											
		1	2	3	4	5	6	7	8	9	10	...	n
items	1	1	1	1			1				1		
	2							1	1	1			
	3	1	1		1				1	1			1
	4		1			1			1	1			
	⋮			1				1					
	m				1					1			

- Identify  $U_9$  and  $U_2$  are similar to  $U_8$



# User-oriented CF

- Recommendation Generation:

	users											
	1	<b>2</b>	3	4	5	6	7	<u>8</u>	<b>9</b>	10	...	n
items	1	1	<b>1</b>	1		1				1		
	2						1	1	1			
	3	1	1		1			1	1			1
	4		1		1			1	1			
	⋮			1			1					
	m			1					<b>1</b>			

- Identify  $I_1$  and  $I_9$  are not yet purchased by  $U_8$



# User-oriented CF

- Recommendation Generation:

		users											
		1	2	3	4	5	6	7	8	9	10	...	n
items	1	1	1	1			1		0.7		1		
	2							1	1	1			
	3	1	1		1				1	1			1
	4		1			1			1	1			
	⋮			1				1					
	m				1				0.9	1			

- Predict by taking weighted sum



# User-oriented CF

## Practical Implementation

- Compute and store all user-user similarities.

- Cosine similarity:  $sim(u,v) = \cos(\vec{u}, \vec{v}) = \frac{\vec{u} \bullet \vec{v}}{\|\vec{u}\|_2 * \|\vec{v}\|_2}$

- Find N items that will be most likely purchased by user **u**.

- Find **k** most similar users to **u**, save to  $U_{sim}$
- Get all items purchased by  $U_{sim}$ , save to  $I_{candidate}$
- Remove unavailable items in  $I_{candidate}$
- Get all items purchased by **u**, save to  $I_{purchased}$
- Take  $I_{candidate} - I_{purchased} = I_{recmd}$
- Re-order items in  $I_{recmd}$  based on sum of user-user similarity

$$pred(u,i) = \frac{\sum_{v \in k-similarUser(u)} userSim(u,v) * r_{vi}}{\sum_{v \in k-similarUser(u)} userSim(u,v)}$$



# Outline

- Introduction
- User-oriented Collaborative Filtering
- **Item-oriented Collaborative Filtering**
- Challenges
- Best Practices



# Item-oriented CF

- Input Data Representation: Users / Items matrix.

		users											
		1	2	3	4	5	6	7	8	9	10	...	n
items	1	1		1			1				1		
	2							1	1	1			
	3	1	1		1				1	1			
	4		1			1			1	1			
	⋮			1				1					
	m												

- Cell value “1” means user purchased the item.
- Data Normalization is not shown on this slide.



# Item-oriented CF

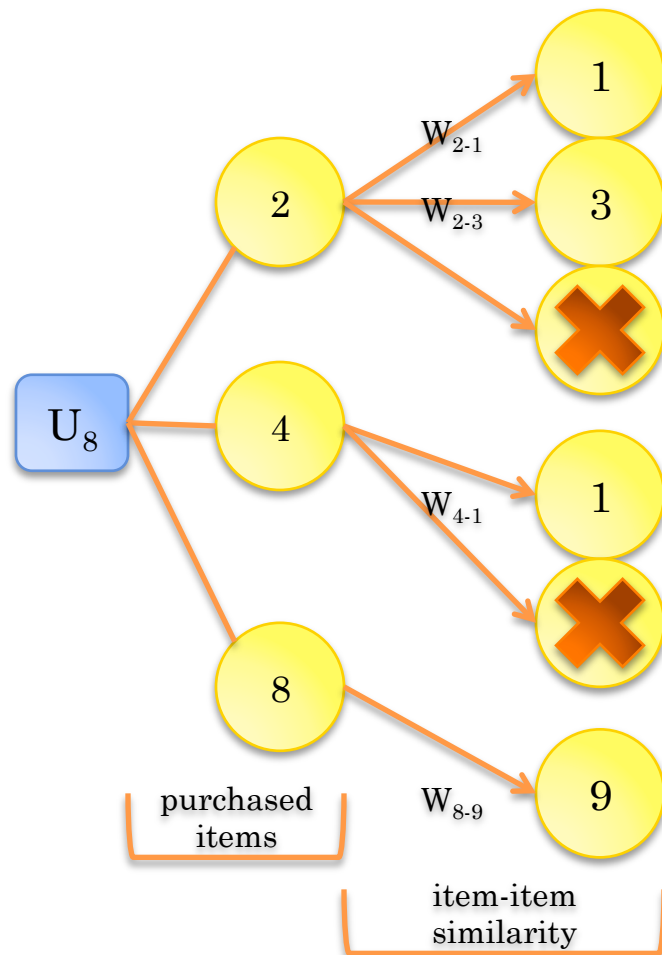
- Neighborhood Formation:
  - Find the  $k$  items that have the most similar user vectors

		users											
		1	2	3	4	5	6	7	8	9	10	...	n
items	1	1		1			1				1		
	<b>2</b>		1					1	1	1			
	3	1			1				1	1			
	<b><u>4</u></b>		1			1		1	1	1			
	⋮			1				1					
	m				1					1			



# Item-oriented CF – cont.

- Recommendation Generation



Predict by taking weighted sum

TopN Recmd. for  $U_8$  :  $\{1, 9, 3\}$





# Item-oriented CF

## Practical Implementation

- Compute and store all item-item similarities.

- Cosine similarity:  $sim(a,b) = \cos(\vec{a}, \vec{b}) = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\|_2 * \|\vec{b}\|_2}$

- Find N items that will be most likely purchased by user **u**.
  - Get all items purchased by **u**, save to  $I_{\text{purchased}}$
  - For each item in  $I_{\text{purchased}}$ , find **k** most similar items save them to  $I_{\text{candidate}}$
  - Remove unavailable items in  $I_{\text{candidate}}$
  - Get all items purchased by **u**, save to  $I_{\text{purchased}}$
  - Take  $I_{\text{candidate}} - I_{\text{purchased}} = I_{\text{recmd}}$
  - Re-order items in  $I_{\text{recmd}}$  based on

$$pred(u,i) = \frac{\sum_{j \in purchasedItems(u)} itemSim(i,j) * r_{uj}}{\sum_{j \in purchasedItems(u)} itemSim(i,j)}$$



# Outline

- Introduction
- User-oriented Collaborative Filtering
- Item-oriented Collaborative Filtering
- **Challenges**
- Best Practices



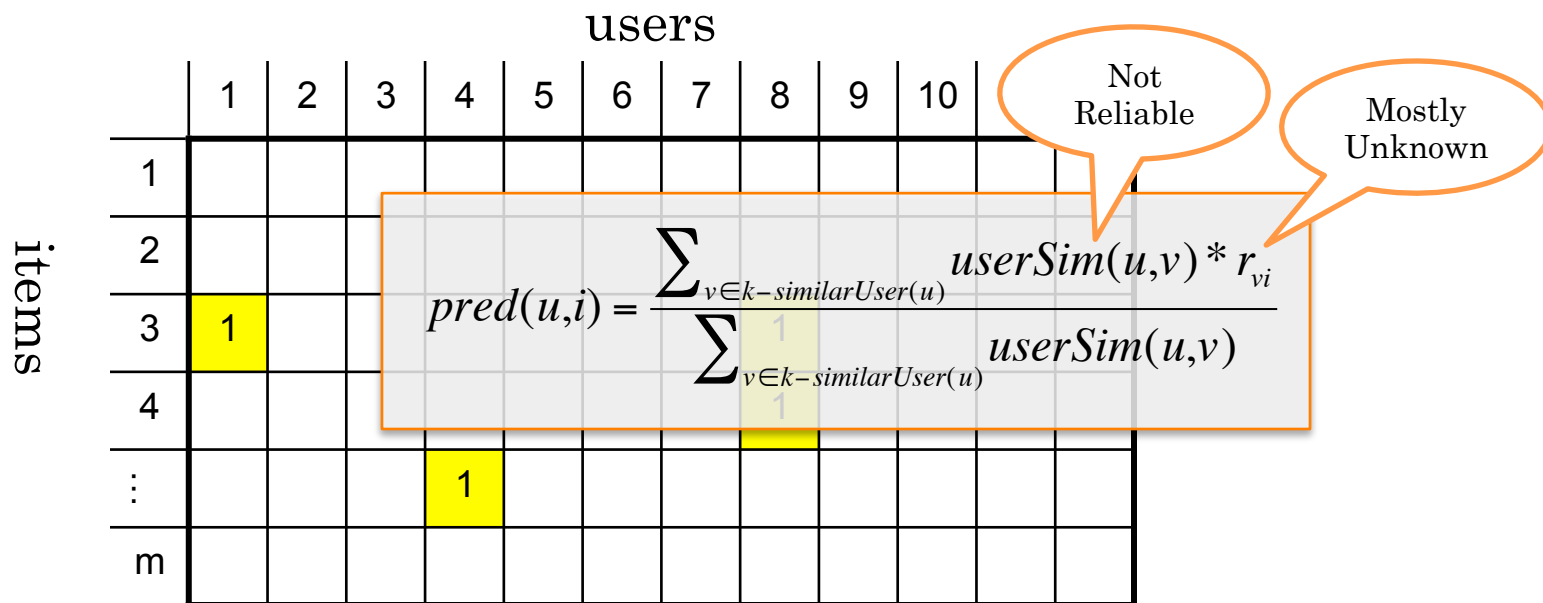
# The Challenges you will face

- Data Sparsity Issue
- Cold Start Problem
- Curse of Dimensionality
- Scalability



# Data Sparsity Issue

- Missing values in the Users / Items matrix.



- Netflix Prize data set: 98.82% of cells are blank
- Typical e-commerce txn data set can be 10-100 time more sparse than Netflix Prize data set !!



# Cold Start Problem

- It occurs when new item or new user is added to the data matrix.

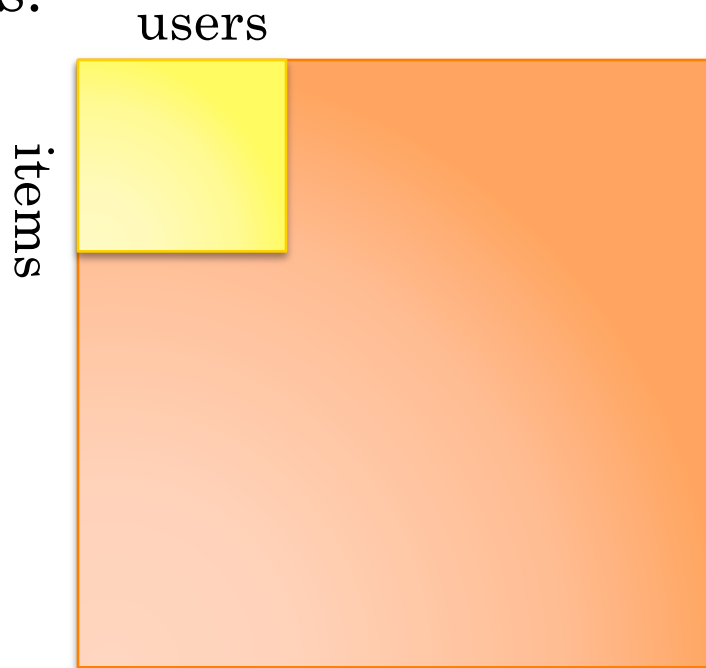
		users											
		1	2	3	4	5	6	7	8	9	10	...	n
items	1	1		1			1				1		
	2							1	1	1			
	3	1	1		1				1	1			1
	4		1			1			1	1			
	⋮			1				1					
m													

- RE does not have enough knowledge about this new user or this new item yet.
- Content-based REs can be incorporated to alleviate cold start problem.



# Curse of Dimensionality

- Adding more features (items or users) can increase the noise, and hence the error.
- There aren't enough observations to get good estimates.



# Scalability

- User neighborhood formation:  $O(n^2)$  for  $n$  users
- Item neighborhood formation:  $O(m^2)$  for  $m$  items
- When  $m$  (# of items)  $\ll n$  (# of users), item-based CF will be more efficient than user-based CF
- Ability to update neighborhood incrementally



# Outline

- Introduction
- User-oriented Collaborative Filtering
- Item-oriented Collaborative Filtering
- Challenges
- **Best Practices**





# Best Practices

- Understand the data thoroughly
- Define business objectives and conversion metrics judiciously
- Understand context and user intent
- Apply adaptive reinforcement learning
- Optimize RE using cost-based methods
- Be aware of data-shift issue
- Optimize marketing messages delivered with recommendation results



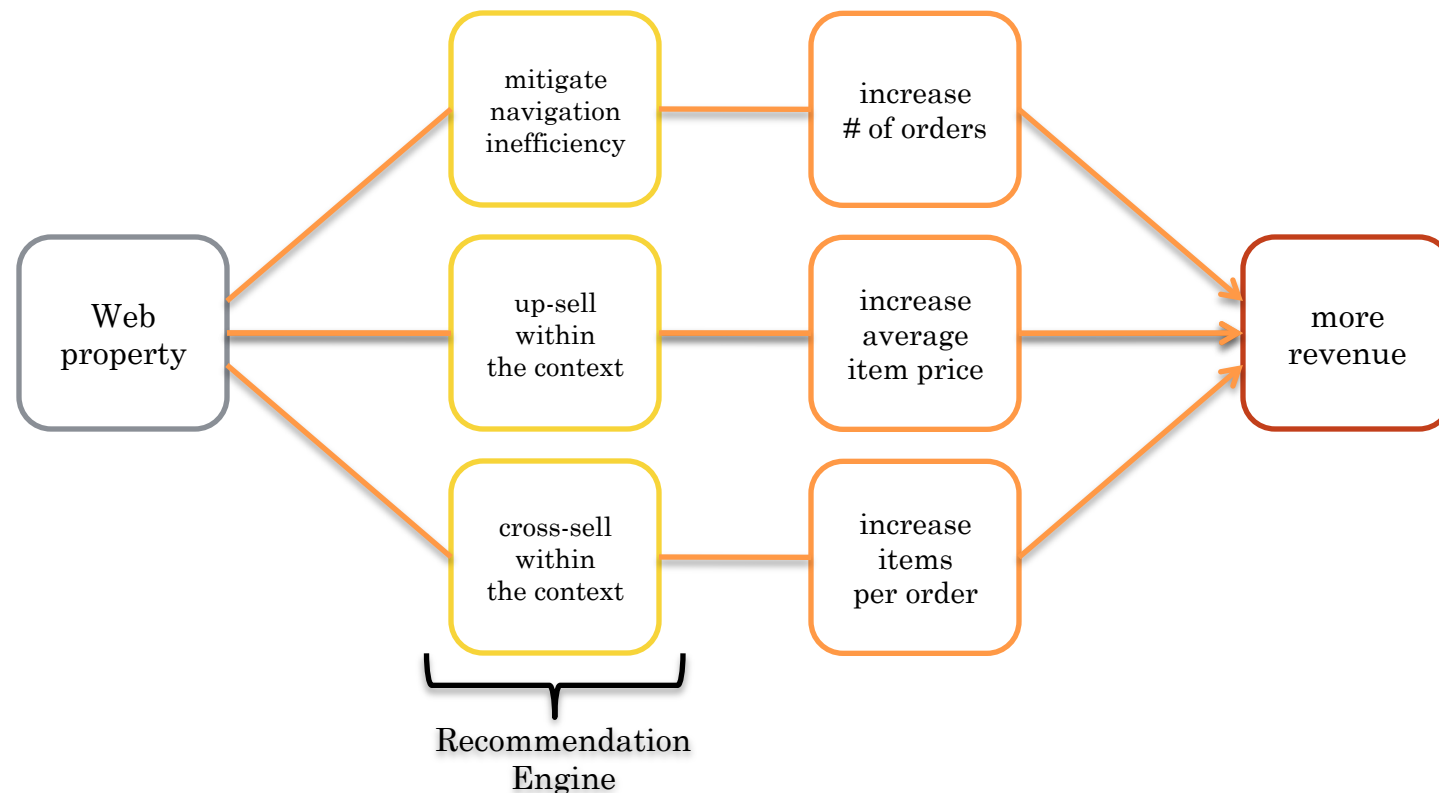
# Understand the data thoroughly

- What data are available?
- E-commerce data set typically contains:
  - Clickstream
  - Shopping cart / Saved Items / Wish list / Shared Item
  - Order / Return
  - User profile
  - User ratings
- How are these data points being collected?
- Is there pre-existing bias in the data? or leakage?
- Is the data related to what we want to predict?



# Define business objectives and conversion metrics judiciously

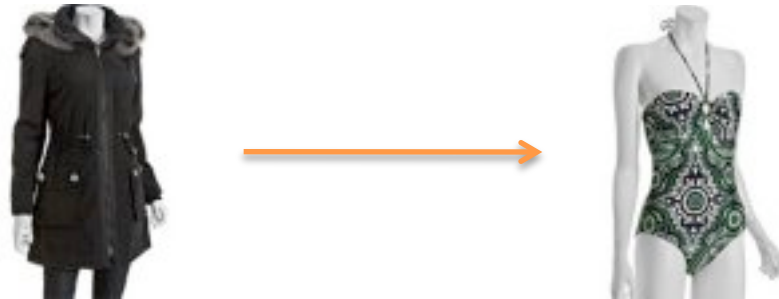
- Defining correct conversion metrics can be a competitive advantage.



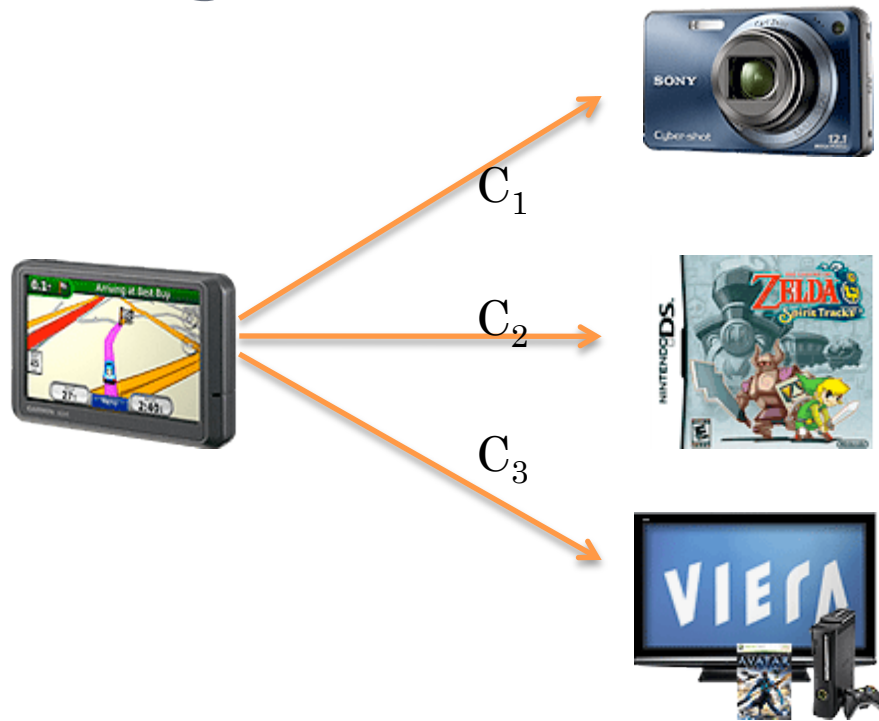
# Understand context and user intent

- Context should be considered when RE making recommendation.

Month: December  
Temperature: 45°F



# Apply adaptive reinforcement learning



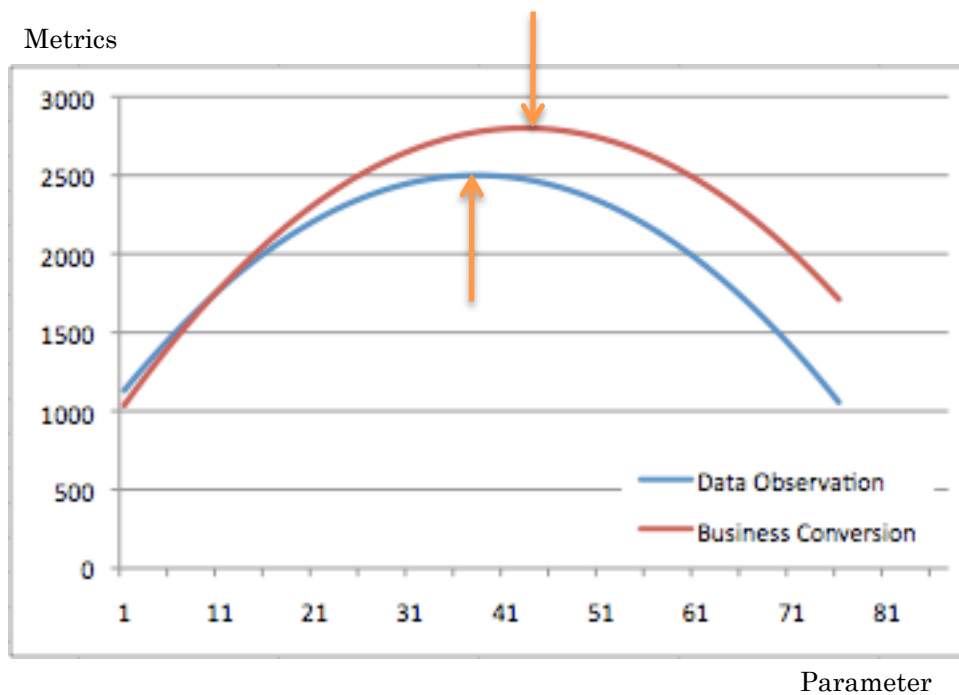
Incorporating clickstream adaptive reinforcement

$$\text{pred}(u,i) = \frac{\sum_{j \in \text{purchasedItems}(u)} \text{itemSim}(i,j) * r_{uj}}{\sum_{j \in \text{purchasedItems}(u)} \text{itemSim}(i,j)} + W_i C_i$$



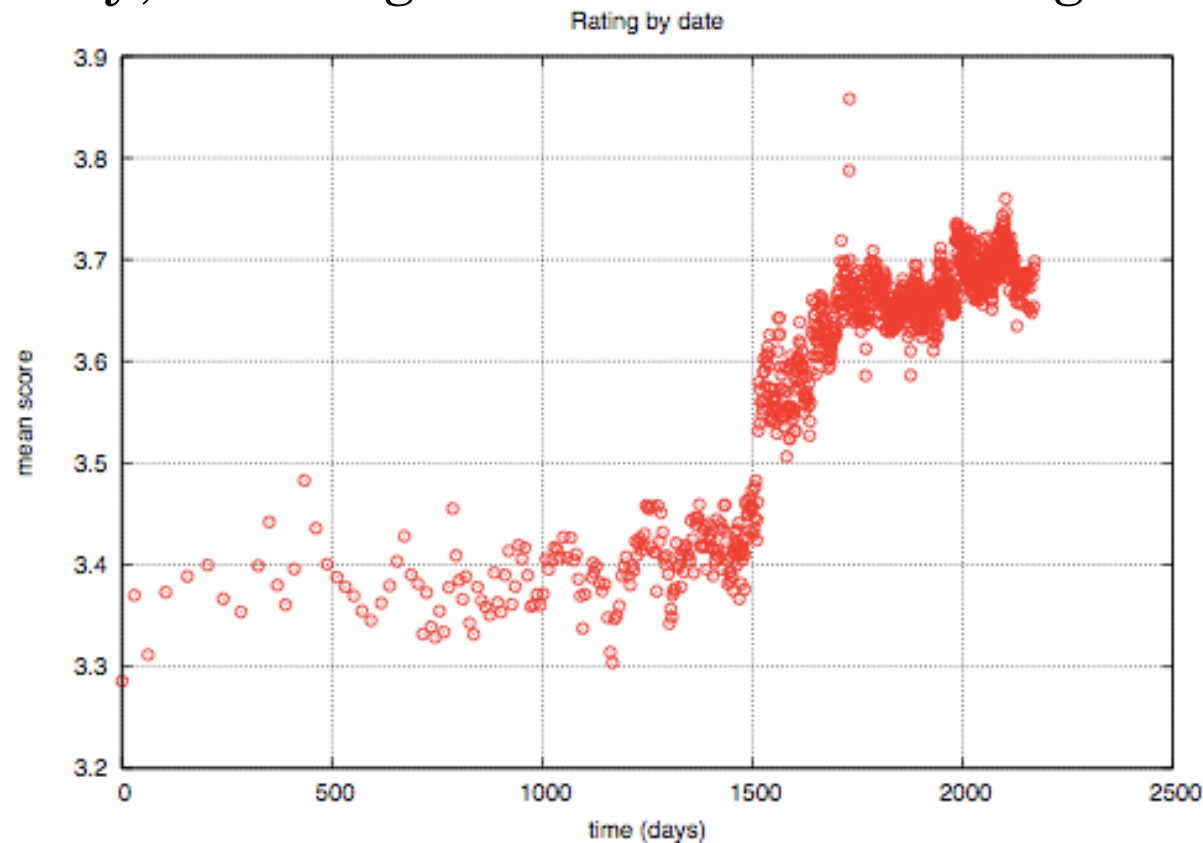
# Optimize RE using cost-based models

- Cost-based engine optimization



# Be aware of the data-shift issue

- Data collection UI changes will influence data significantly, creating artificial data shifting



## Netflix prize data set

Y. Koren, "Collaborative Filtering with Temporal Dynamics," Proc. 15th ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining (KDD 09), ACM Press, 2009, pp. 447-455.12

# Optimize marketing message delivered with recommendation result

## ○ It's How You Say It

### What Do Customers Ultimately Buy After Viewing This Item?



**51%** buy the item featured on this page:

Canon Rebel XS 10.1MP Digital SLR Camera with EF-S 18-55mm

Add to cart to see price.

**Complete Your Series**



**17%** buy

[Canon Digital Rebel XSi 12.2 MP Digital SLR Camera with EF-S 18-55mm](#)

[Click to see price](#)

**Customers Who Bought Items in Your Recent History Also Bought:**



**16%** buy

[Canon EF-S 55-250mm f/4.0-5.6 IS Telephoto Zoom Lens for Canon Digit](#)

~~\$231.00~~



**8%** buy

[Canon](#)

[Click](#)

### Today's Recommendations For You

Here's a daily sample of items recommended for you. Click here to [see all recommendations](#).

### Explore similar items

### Frequently Bought Together

Customers buy this item with [Transcend 8 GB SDHC Class 6 Flash Memory Card TS8GSDHC6](#)



+



**Price For Both:** To see our price, add these items to your cart. [View](#)



Add both to Cart



Add both to Wish List

[Show availability and shipping details](#)

Screenshots from Amazon.com







# RECOMMENDATION ENGINE DEMYSTIFIED

Alex Lin

Intelligent Mining

Email: [alin@intelligentmining.com](mailto:alin@intelligentmining.com)

Twitter: DKALab