



## INTRODUCTION TO DATA ANALYTICS

# Recommender System

Dr. Rathachai Chawuthai

Department of Computer Engineering

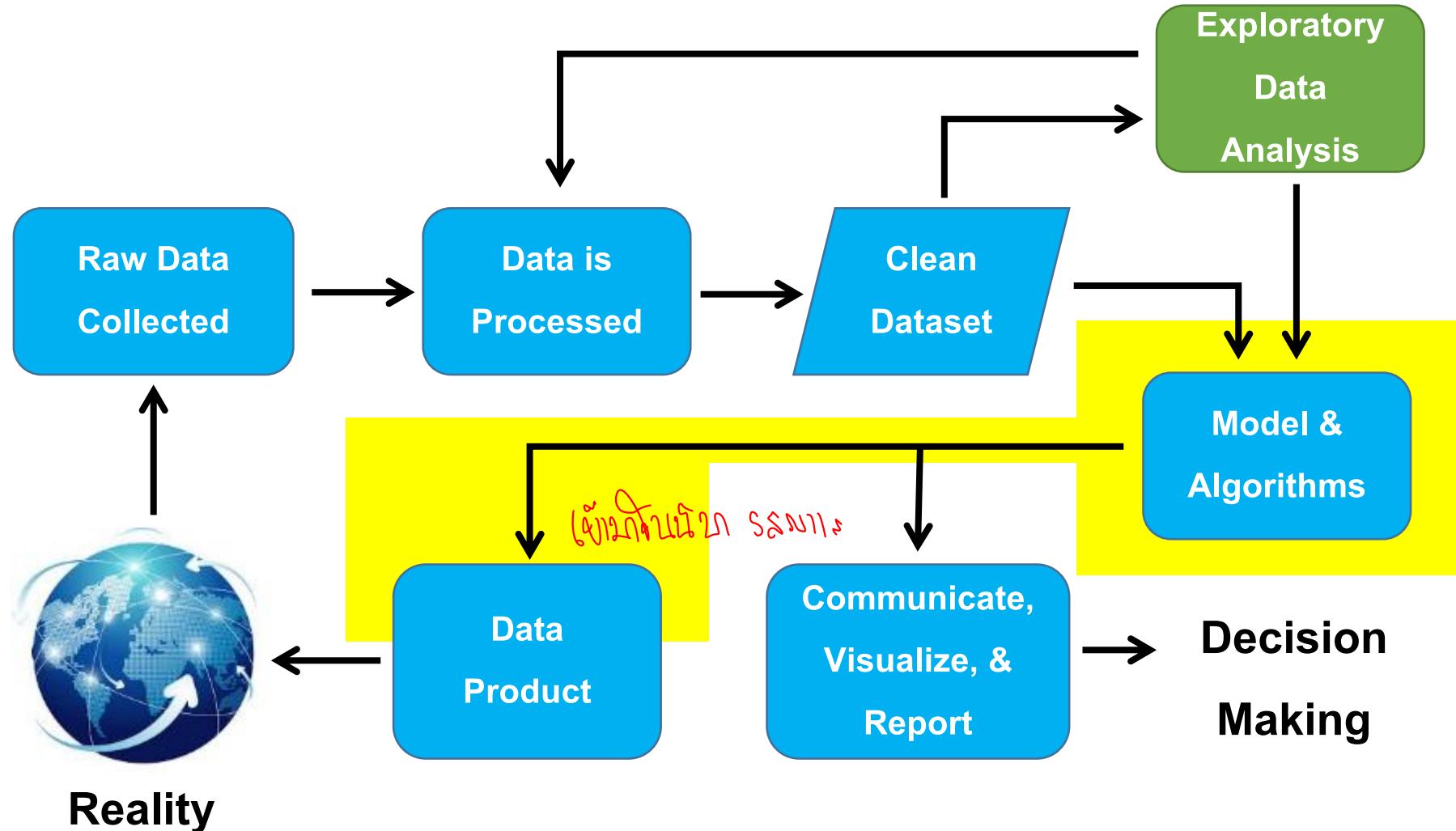
Faculty of Engineering

King Mongkut's Institute of Technology Ladkrabang

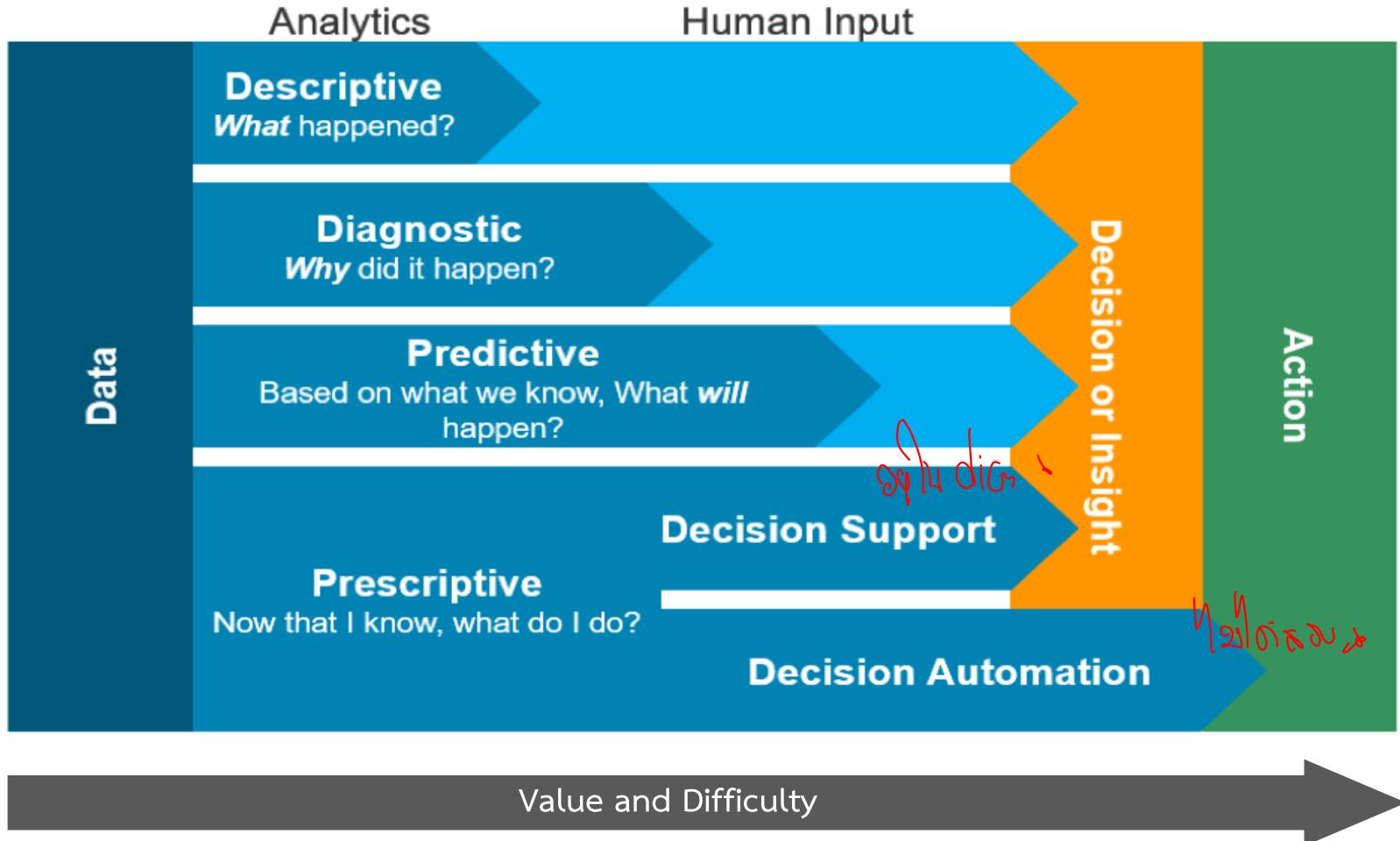
# Agenda

- Overview
- Evaluating Recommender Systems
- \* • Collaborative Recommendation  
*କାଲାପରିବହନ ପରିମାଣନ୍ତିରଣ*
- Content-based Recommendation
- Knowledge-based Recommendation
- Hybrid Recommender Approach

# Data Science Process



# Data Analytics



- Ref:
- Four types of analytics capability (Gartner, 2014)
  - (image) <https://www.healthcatalyst.com/closed-loop-analytics-method-healthcare-data-insights>

# Overview

# Recommender Systems

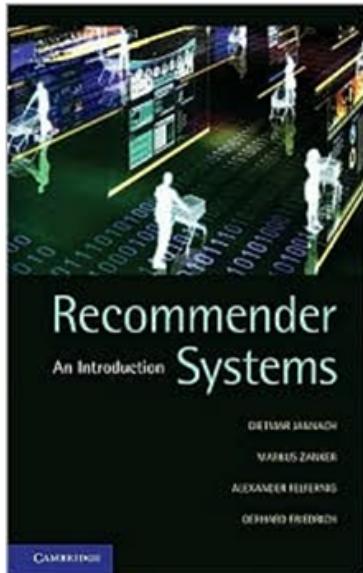
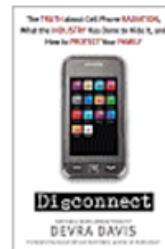


Table of Contents

## Customers who bought this also bought



## Recommender Systems: An Introduction

by [Dietmar Jannach](#), [Markus Zanker](#), [Alexander Felfernig](#), [Gerhard Friedrich](#)

### AVERAGE CUSTOMER RATING:

★★★★★ ([Be the first to review](#))



Registrieren, um sehen zu können, was  
deinen Freunden gefällt.

### FORMAT:

Hardcover

NOOKbook (eBook) - not available

[Tell the publisher you want this in NOOKbook format](#)

### NEW FROM BN.COM

\$65.00 List Price

**\$52.00** Online Price

(You Save 20%)

**Add to Cart**

### NEW & USED FROM OUR

New starting at **\$56.46** (You S

Used starting at **\$51.98** (You :

**See All Prices**

ମୟାନ୍ ହୁଏ ବେଳେ → ଯେ ପାଇଁ କେଣ୍ଟିଲେଖନ୍ତିରେ

# Recommender Systems

## Application areas

You may also like



Jack & Jones  
JAMIE - Polo shirt - orange  
£21.00  
Free delivery & returns

### ALTERNATIVE PRODUCTS

Beko Washing Machine

Code: WMB81431LW

**£269.99**

Zanussi Washing Machine

Code: ZWH6130P

**£269.99**

Blomberg Washing Machine

Code: WNF6221

**£299.99**

You may also like

gadget



★★★★☆ (109)



★★★★★ (53)



★★★★☆ (33)

### Related hotels...



Hotel 41

★★★★ 1,170 Reviews

London, England

Show Prices

Read Commented Recommended



Germany Just Rejected The Idea That The European Bailout Fund Would Buy Spanish Debt



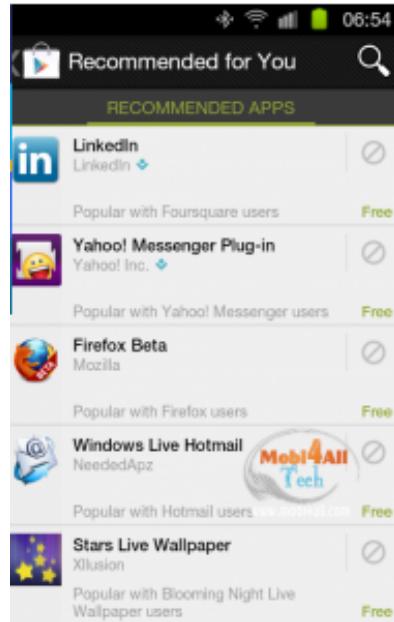
There Is Almost No Gold In The Olympic Gold Medal

MOST POPULAR RECOMMENDED

How to Break NRA's Grip on Politics: Michael R. Bloomberg +

Growth in U.S. Slows as Consumers Restrain Spending +

# In the Social Web



Jobs you may be interested in Beta Email Alerts | See More »

-  **Technical Sales Manager - Europe**  
Thermal Transfer Products - Home office X
-  **Senior Program Manager (f/m)**  
Johnson Controls - Germany-NW-Burscheid X

Groups You May Like More »

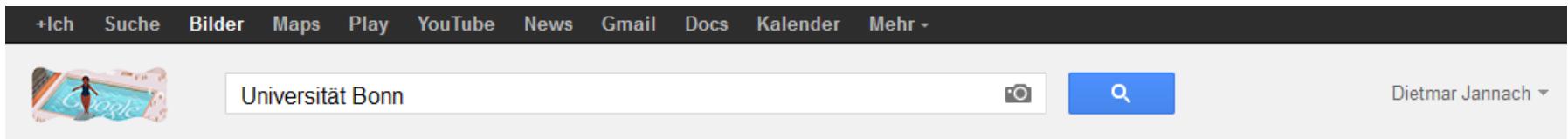
-  **Advances in Preference Handling**  
[Join](#)
-  **FP7 Information and Communication Technologies (ICT)**  
[Join](#)
-  **The Blakemore Foundation**  
[Join](#)



# Even more ...

---

- Personalized search



- "Computational advertising"

A screenshot of a search results page featuring personalized computational advertising. On the left, there is a profile picture of a woman and a blue button with a white outline. Next to it, the text reads 'gefällt Danubius Hotels in Bad Bük und Bad Sárvár.' Below this, there is another profile picture of a couple in a pool and a blue button with a white outline. Next to it, the text reads 'Danubius Hotels in Bad Bük und Bad Sárvár' followed by a 'Gefällt mir' button. On the right, there are two ads. The first ad is for 'Gewinne Deine Insel' from bahamas.canusa.de, featuring a small green island in the ocean. The second ad is for 'Citroën Austria', featuring a red Citroën car. Both ads include promotional text in German.

# Why using Recommender Systems?

የኢንተርፕራይድ ወጥቶ

- **Value for the customer**
  - Find things that are interesting
  - Narrow down the set of choices
  - Help me explore the space of options
  - Discover new things
  - Entertainment
  - ...
- **Value for the provider**
  - Additional and probably unique personalized service for the customer
  - Increase trust and customer loyalty
  - Increase sales, click trough rates, conversion etc.
  - Opportunities for promotion, persuasion
  - Obtain more knowledge about customers
  - ...

# Real-world check

---

- **Myths from industry**
  - Amazon.com generates X percent of their sales through the recommendation lists ( $30 < X < 70$ )
  - Netflix (DVD rental and movie streaming) generates X percent of their sales through the recommendation lists ( $30 < X < 70$ )
- **There must be some value in it**
  - See recommendation of groups, jobs or people on LinkedIn
  - Friend recommendation and ad personalization on Facebook
  - Song recommendation at last.fm
  - News recommendation at Forbes.com (plus 37% CTR)
- **Academia**
  - A few studies exist that show the effect
    - increased sales, changes in sales behavior

# Recommender systems

---

- RS seen as a function

- ① Given:

- User model (e.g. ratings, preferences, demographics, situational context)
  - Items (with or without description of item characteristics)

- Find:

- Relevance score. Used for ranking. → ព័ត៌មាននេះរាយការណ៍ប្រើប្រាស់នៅក្នុងនៃទាំងអស់ item

- Finally:

- Recommend items that are assumed to be relevant

- But:

- Remember that relevance might be context-dependent
  - Characteristics of the list itself might be important (diversity)

# Paradigms of recommender systems

---

Recommender systems reduce information overload by estimating relevance



Recommendation component

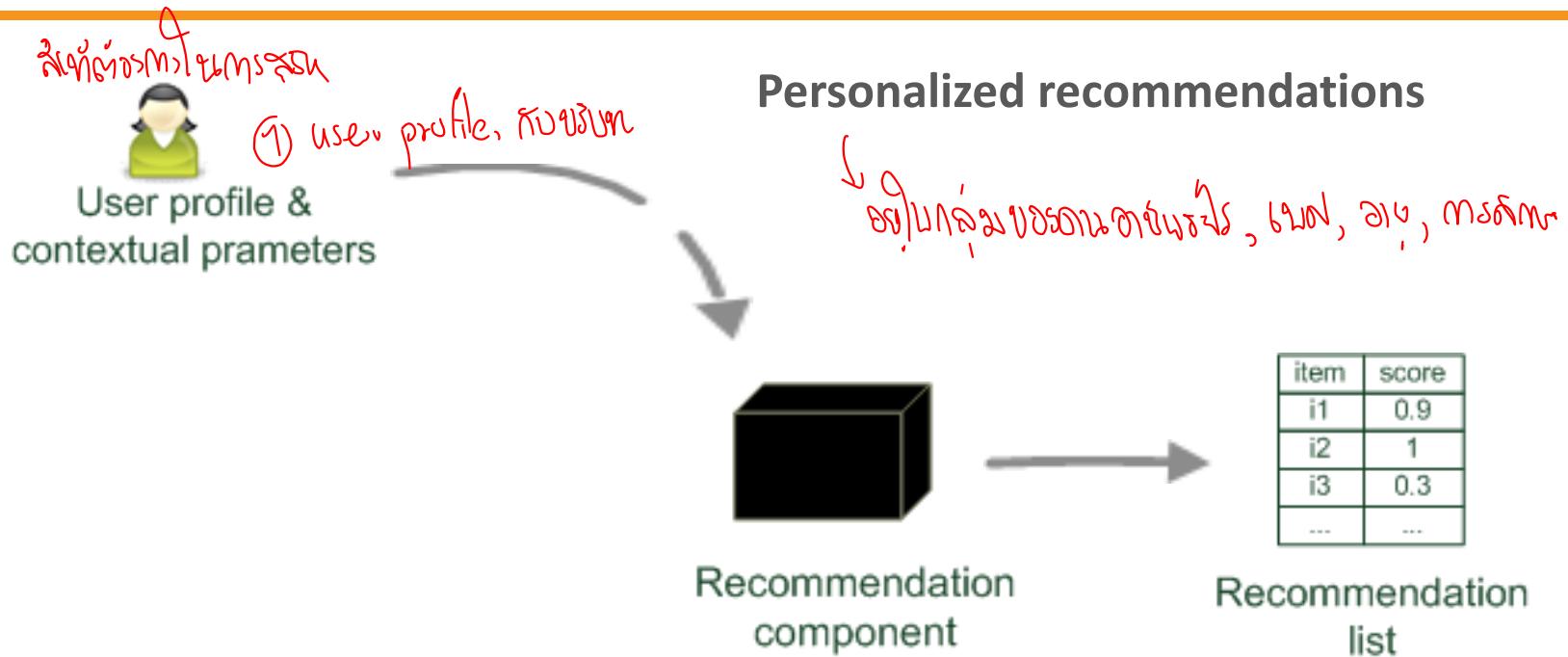


item	score
i1	0.9
i2	1
i3	0.3
...	...

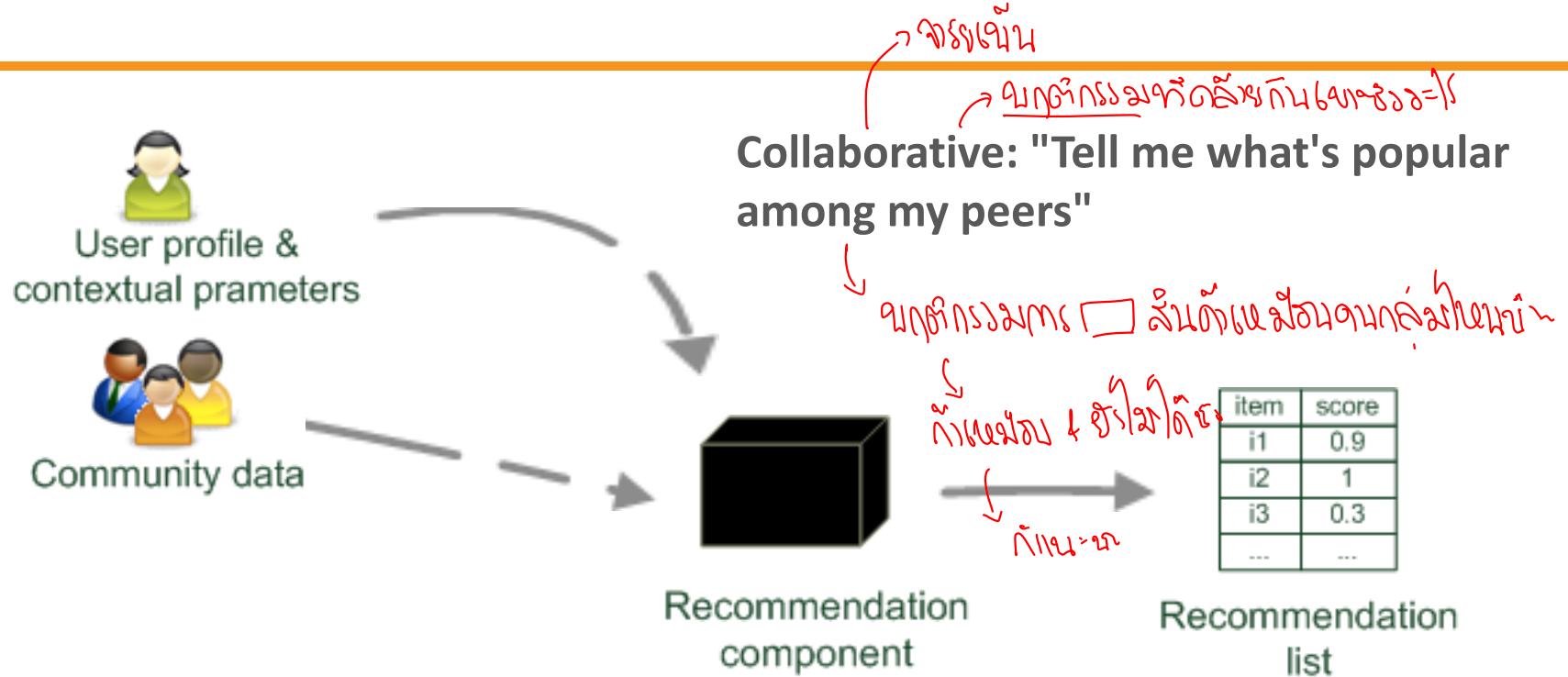
Recommendation list

↳  
ανανεώστε user\_x από την πλήρωση στο item\_n για γενικές

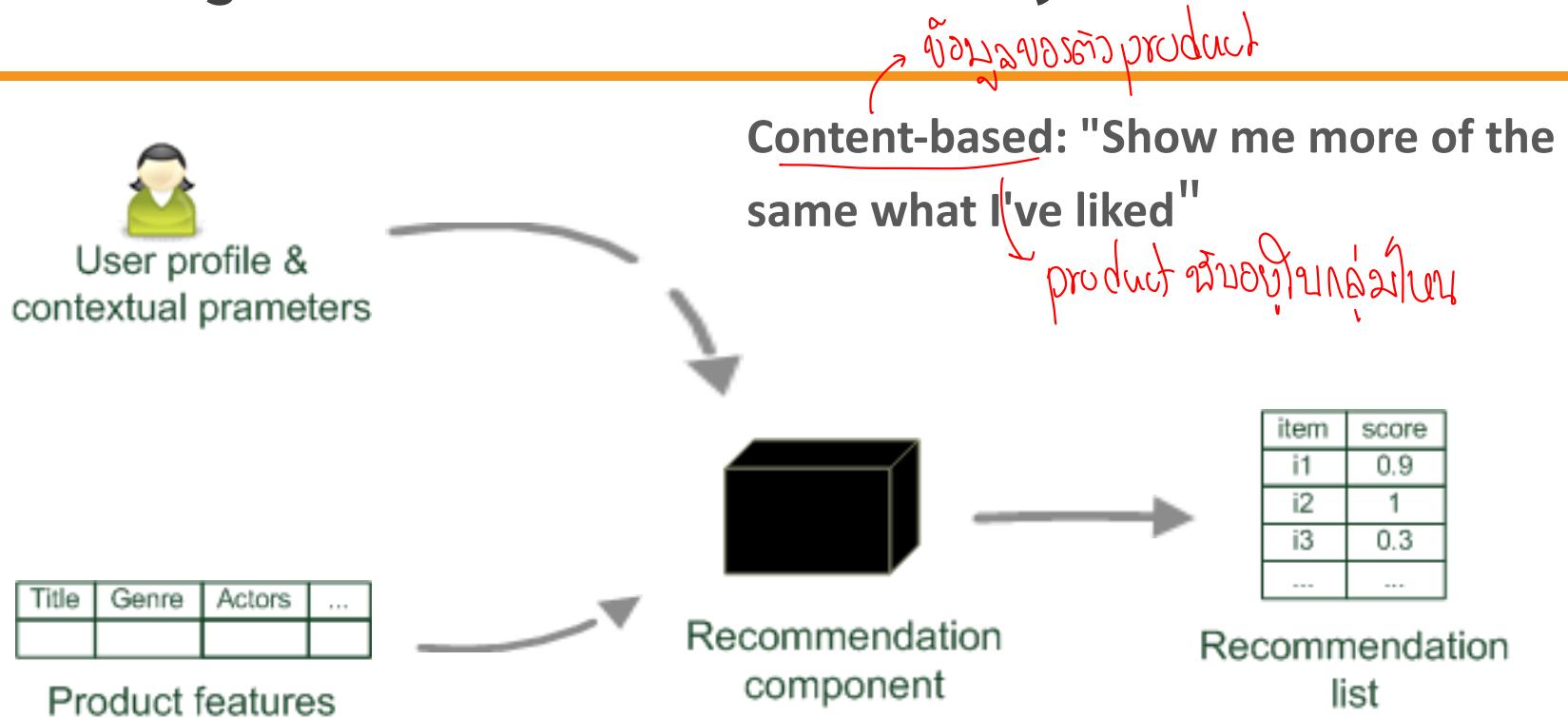
# Paradigms of recommender systems



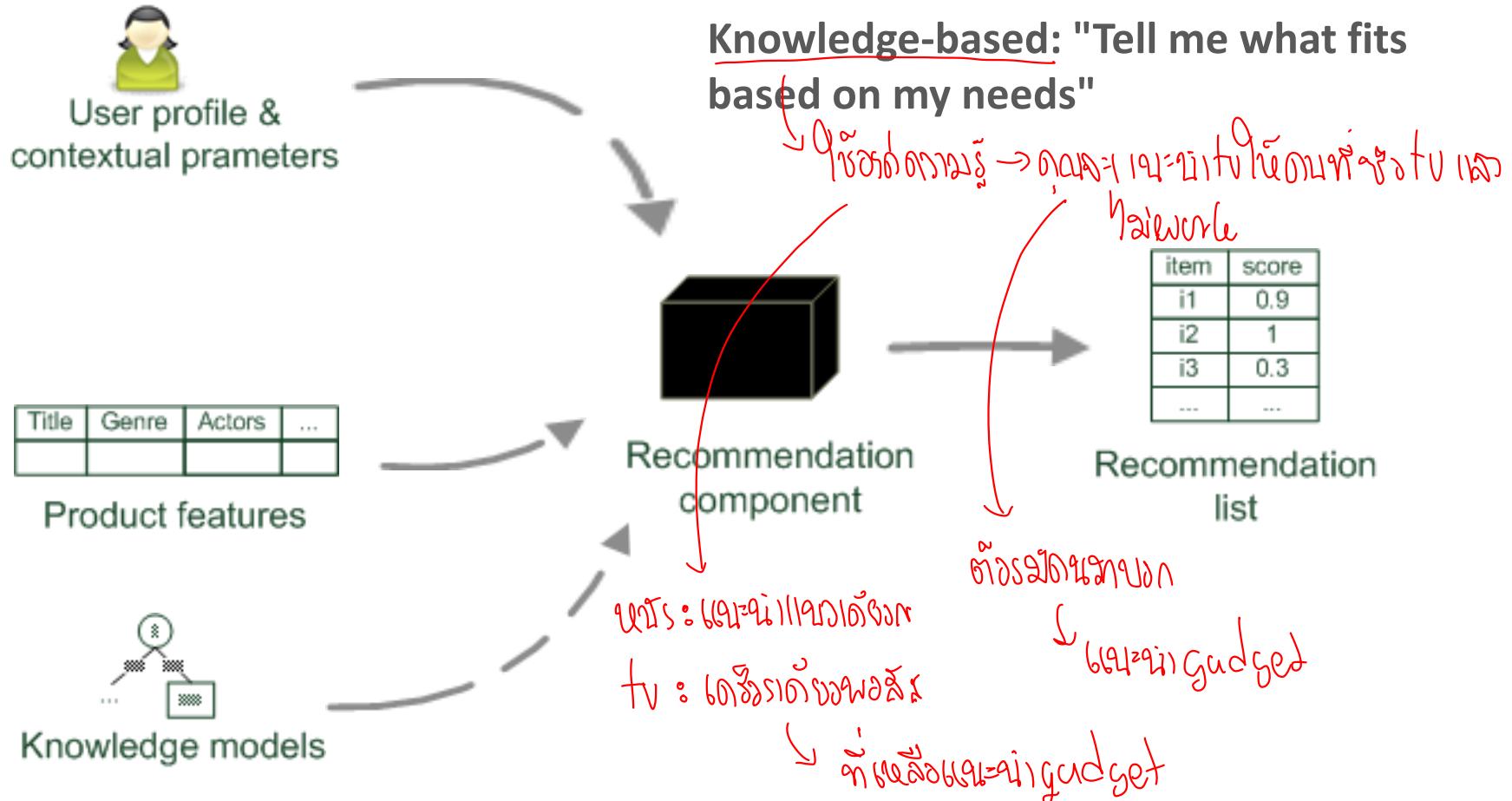
# Paradigms of recommender systems



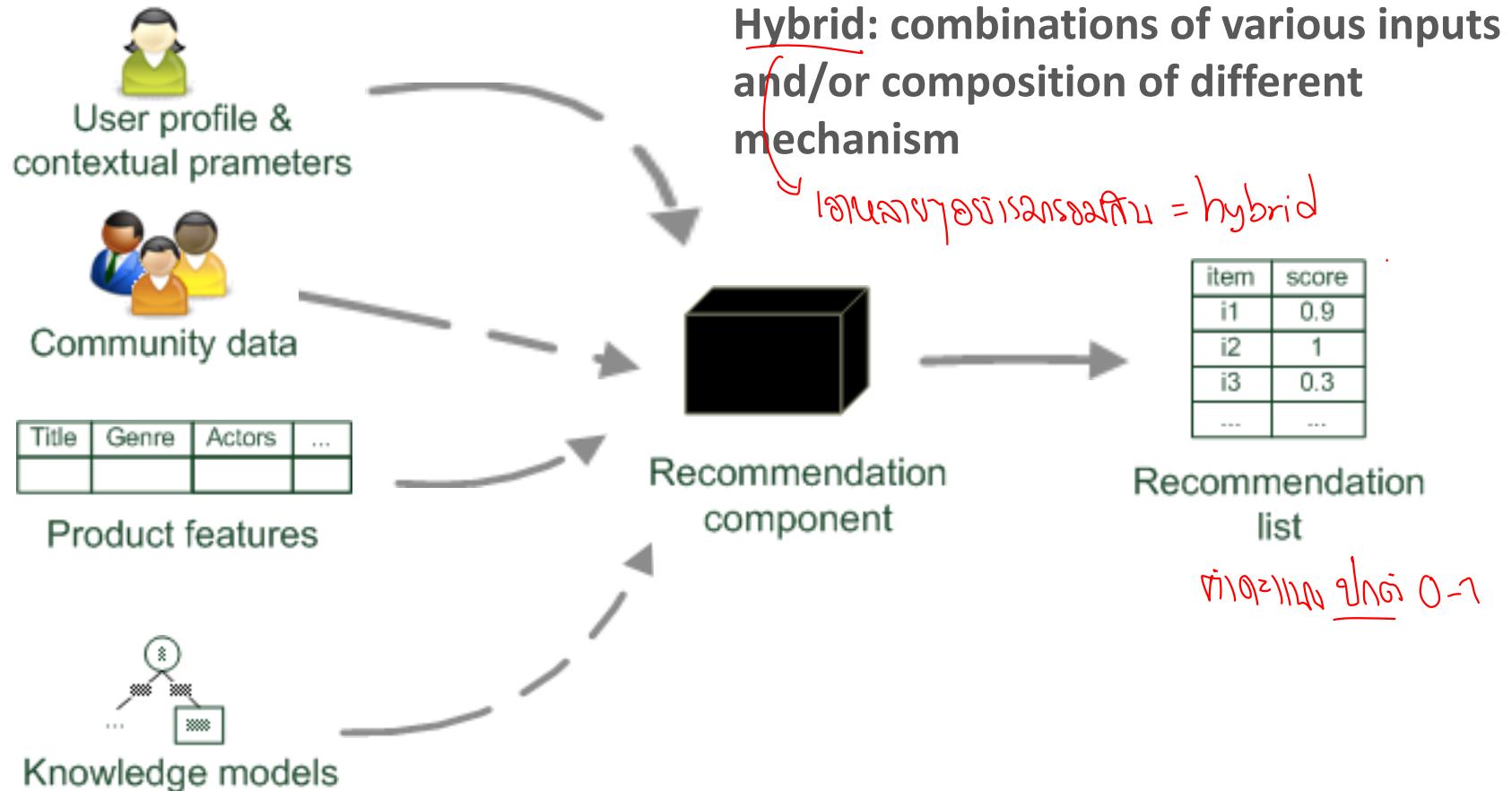
# Paradigms of recommender systems



# Paradigms of recommender systems



# Paradigms of recommender systems



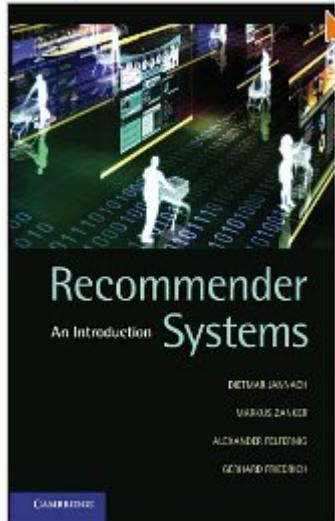
# Recommender systems: basic techniques

	Pros  ចិត្ត	Cons  ចិត្ត
<b>Collaborative</b> ការងារដែលបានរាយជាប្រព័ន្ធដែលបានរាយជាប្រព័ន្ធ	No knowledge-engineering effort, serendipity of results, learns market segments ការងារដែលបានរាយជាប្រព័ន្ធដែលបានរាយជាប្រព័ន្ធ	Requires some form of rating feedback, cold start for new users and new items ត្រូវមានសម្រាប់សម្រាប់អ្នកប្រើប្រាស់ និងផ្លូវការថ្មី
<b>Content-based</b> ការងារដែលបានរាយជាប្រព័ន្ធដែលបានរាយជាប្រព័ន្ធ	No community required, comparison between items possible ការងារដែលបានរាយជាប្រព័ន្ធដែលបានរាយជាប្រព័ន្ធ	Content descriptions necessary, cold start for new users, no surprises ត្រូវព័ត៌មានពីផ្លូវការថ្មី និងផ្លូវការថ្មី
<b>Knowledge-based</b>	Deterministic recommendations, assured quality, no cold-start, can resemble sales dialogue ការងារដែលបានរាយជាប្រព័ន្ធដែលបានរាយជាប្រព័ន្ធ	Knowledge engineering effort to bootstrap, basically static, does not react to short-term trends ត្រូវការងារដែលបានរាយជាប្រព័ន្ធដែលបានរាយជាប្រព័ន្ធ

# Evaluating Recommender Systems

ກວດສອບ ປະເມີນຫຼາຍ.

# Recommender Systems in e-Commerce



- One Recommender Systems research question?
  - What should be in that list?

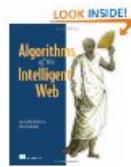
↳  
↳  
↳  
↳  
↳



## Customers Who Bought This Item Also Bought



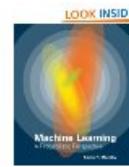
[Recommender Systems Handbook](#)  
Francesco Ricci  
Hardcover  
**\$167.73**



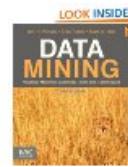
[Algorithms of the Intelligent Web](#)  
Haralambos Marmanis  
**★★★★★ (14)**  
Paperback  
**\$26.76**



[Programming Collective Intelligence: ...](#)  
Toby Segaran  
**★★★★★ (91)**  
Paperback  
**\$25.20**

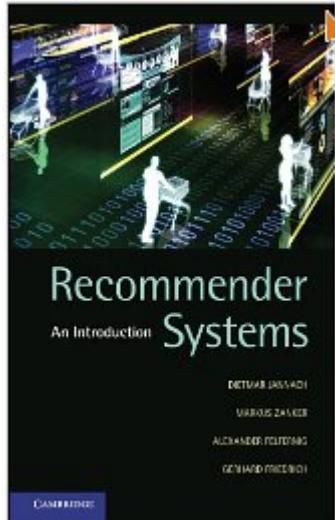


[Machine Learning: A Probabilistic Approach](#)  
Kevin P. Murphy  
**★★★★★ (15)**  
Hardcover  
**\$81.00**

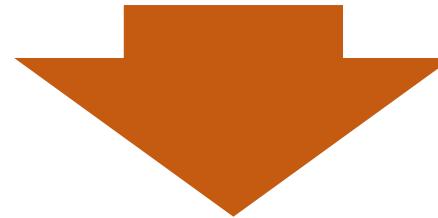


[Data Mining: Practical Machine Learning Tools and Techniques](#)  
Ian H. Witten  
**★★★★★ (29)**  
Paperback  
**\$42.61**

# Recommender Systems in e-Commerce



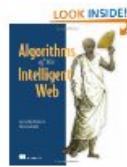
- Another question both in research and practice
  - How do we know that these are good recommendations? ວິທີໄຮລ່ສູນເຊີ່ວຍ



## Customers Who Bought This Item Also Bought



[Recommender Systems Handbook](#)  
Francesco Ricci  
Hardcover  
\$167.73



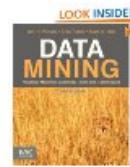
[Algorithms of the Intelligent Web](#)  
Haralambos Marmanis  
★★★★★ (14)  
Paperback  
\$26.76



[Programming Collective Intelligence: ...](#)  
Toby Segaran  
★★★★★ (91)  
Paperback  
\$25.20

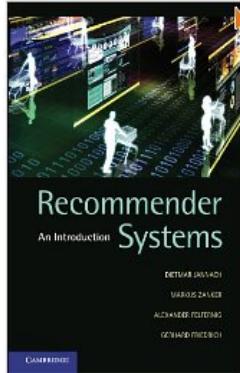


[Machine Learning: A Probabilistic Approach](#)  
Kevin P. Murphy  
★★★★★ (15)  
Hardcover  
\$81.00



[Data Mining: Practical Machine Learning Tools and Techniques](#)  
Ian H. Witten  
★★★★★ (29)  
Paperback  
\$42.61

# Recommender Systems in e-Commerce



- This might lead to ...
    - What is a good recommendation?
    - What is a good recommendation **strategy**?
    - What is a good recommendation strategy  
for my business?
- Handwritten notes in red:  
↳ මෙයින්ගේ ප්‍රාග්ධනය සඳහා  
↳ ප්‍රතිච්චිත ප්‍රාග්ධනය  
↳ නොමැත්තුවක්



## Customers Who Bought This Item Also Bought

 Recommender Systems Handbook Francesco Ricci Hardcover \$167.73	 Algorithms of the Intelligent Web Haralambos Marmanis ★★★★★ (14) Paperback \$26.76	 Programming Collective Intelligence: ... Toby Segaran ★★★★★ (91) Paperback \$25.20	 Machine Learning: A Probabilistic Approach Kevin P. Murphy ★★★★★ (15) Hardcover \$81.00	 Data Mining: Practical Machine Learning Tools and Techniques Ian H. Witten ★★★★★ (29) Paperback \$42.61
--	--	--	---	---

# What is a good recommendation?

What are the measures in practice? → ຕົກສ້າງ



ອຳນວຍດົກ

ຮັບຮູ້ອະນຸມາ

- Total sales numbers
- Promotion of certain items
- ...
- Click-through-rates
- Interactivity on platform
- ...
- Customer return rates
- Customer satisfaction and loyalty



Best Sellers

①

# MAE

---

- Mean Absolute Error

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

የኩስ ተወስኗል ነው ይጠናና ነው

# MAE

(for rating)

Qaw  
Isruwarrating

Movie	Actual Rating ( $y_i$ )
A	4
B	3
C	
D	5
E	2
F	

# MAE

(for rating)

សមតុល្យបានឱ្យលាងទូទៅ

Movie	Actual Rating ( $y_i$ )	Predicted Rating ( $\hat{y}_i$ )
A	4	4
B	3	4
C		3
D	5	5
E	2	3
F		2

# MAE

(for rating)

Movie	Actual Rating ( $y_i$ )	Predicted Rating ( $\hat{y}_i$ )	$ y_i - \hat{y}_i $
A	4	4	0
B	3	4	1
C		3	
D	5	5	0
E	2	3	1
F		2	

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

# MAE

(for rating)

ឧប្បជ្ជនៃការពិនិត្យ

Movie	Actual Rating ( $y_i$ )	Predicted Rating ( $\hat{y}_i$ )	$ y_i - \hat{y}_i $
A	4	4	0
B	3	4	1
C		3	
D	5	5	0
E	2	3	1
F		2	

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$\frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|}$$

$1+1 = 2$   
 $2/4 = 0.5$  //

# MAE

(for buying)

ສຳເນົາກຮັບແພັນໄມ້ຫຼຸດ

ຊື່ລັດ A

Product	is Bought ( $y_i$ )
A	1
B	1
C	0
D	1
E	1
F	0

# MAE

(for buying)

Product	is Bought ( $y_i$ )	predicting Score ( $\hat{y}_i$ )	$ y_i - \hat{y}_i $
A	1	0.8	0.2
B	1	0.6	0.4
C	0	0.3	0.3
D	1	0.5	0.5
E	1	0.8	0.2
F	0	0.2	0.2

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$\frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|} = 1.8 / 6 = 0.3$$

②

## Top-K Precision → ດາວໂຫຼນ

ເລືອດ່ານຸ້າທີ່ຕ່າງໆຍັງຍັງມາເລີ່ມຈະກຳປົກປິດໃນວິນິວງ ເລືອນຂອງ top k ຊົ່ວໂມງ  
ແລືວກໍ່ໄຟ້ໃນຊັບແລ້ວຢືນສໍາກໍາພວຍໃນຕາມ  
ເບີ້ມີສູງຍິນເລືອດລູກແຜ່ນ

- Pickup  $K$  items from the recommended list
- Count how many actual items that users interact
  - If there are  $N$  items,

$$\text{Top - K Precision} = \frac{N}{K}$$

ກະທົດລົດກັບ business ຈົດງາງ  
ດາວວັດທີ  $K$  ຕາມມັງງານຕາມເປົ້າ  
ເຊີ້ມປັດ (ແບ່ງຂົດສິນຕີ)  
ກົ່ຽວມັດກັບຕື່ມາ

# Top-K Precision

recommend function

$\text{rec}(u, i)$

Item	Score
A	0.59
B	0.66
C	0.41
D	0.51
E	0.67
F	0.09
G	0.13
H	0.49
I	0.66
J	0.64
K	0.13
L	0.81
M	0.40
N	0.06

# Top-K Precision

---

rec(u,i)	
Item	Score
A	0.59
B	0.66
C	0.41
D	0.51
E	0.67
F	0.09
G	0.13
H	0.49
I	0.66
J	0.64
K	0.13
L	0.81
M	0.40
N	0.06

← Bought Item

← Bought Item

← Bought Item

← Bought Item

# Top-K Precision

---

rec(u,i)	
Item	Score
L	0.81
E	0.67
I	0.66
B	0.66
J	0.64
A	0.59
D	0.51
H	0.49
C	0.41
M	0.40
K	0.13
G	0.13
F	0.09
N	0.06

(70)

**Bought Item**

**Bought Item**

**Bought Item**

**Bought Item**

**Bought Item**

# Top-K Precision

rec(u,i)	
Item	Score
L	0.81
E	0.67
I	0.66
B	0.66
J	0.64
A	0.59
D	0.51
H	0.49
C	0.41
M	0.40
K	0.13
G	0.13
F	0.09
N	0.06

TOP 5

$$\text{Top - 5 Precision} = \frac{N}{K}$$

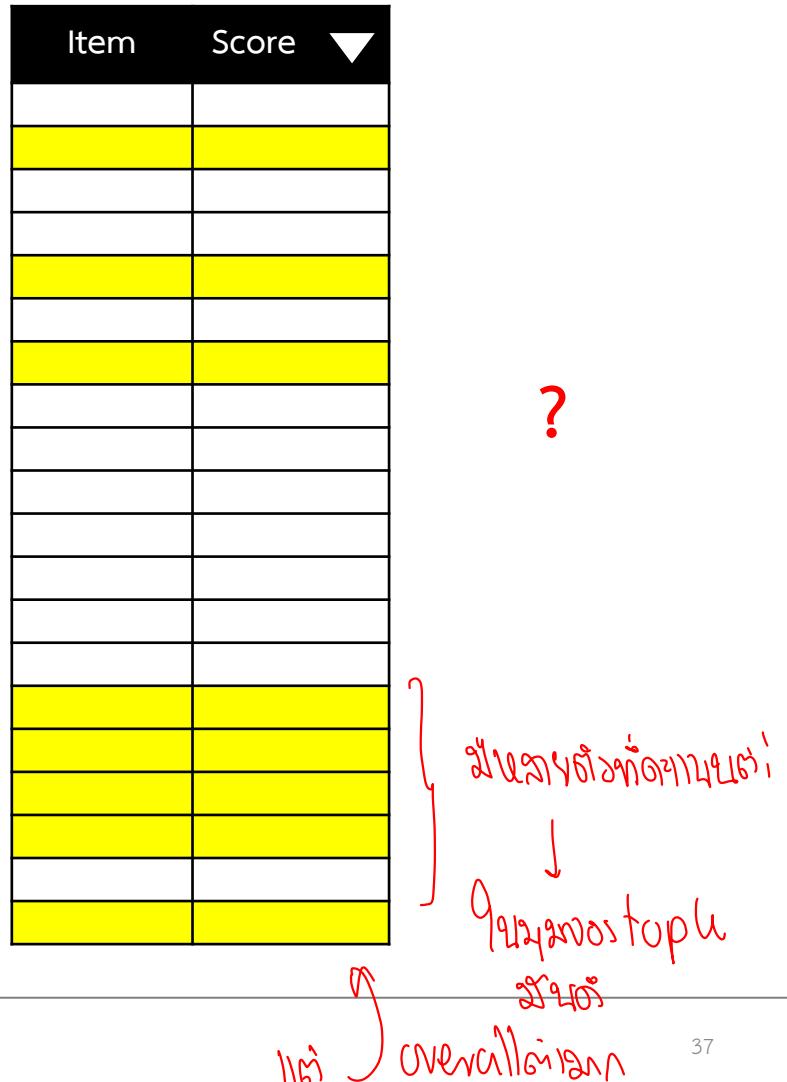
$= \frac{3}{5}$

Vasilevski  
/  
Rouhani

# Top-K Precision

---

What happened if



# Area Under ROC Curve (AUC)

- Area Under ROC Curve for Recommender System

$$AUC = \frac{n' + 0.5n''}{n}$$

ជាមួយកំពង់ទាំង test ដូចជាប្រធានមាត្រា  
 ចាត់recom កំណត់  
 ជាមួយកំពង់ទាំងអស់  
 និងនានាកំណត់  
 ចាត់កំពង់recom

For ***n'*** comparisons,

- n'*** is number of times when the tested items have higher score than the recommended items.
- n''*** is number of times when the tested items have same score as the recommended items.

# AUC

Item	Score ▼
E	0.8
C	0.7
B	0.6
A	0.6
D	0.5
F	0.3
G	0.2

rec(u,i) ໄອເຂົາໄສຮຽນເລື່ອ

- All items are sorted by their score.
- Comparisons are
  - C compares to E, A, D, F and G (5 times)  
ເຖິງພຶກສັງຫຼັບທີ່ໄລ້ຮູ້
  - B compares to E, A, D, F and G (5 times)
- There are  $5+5 = 10$  comparisons
- So,  $n = 10$  ເດືອນປະກາດໃຫຍ່ປົກກວດ ສຳ x ສຳ

# AUC

Item	Score
E	0.8
C	0.7
B	0.6
A	0.6
D	0.5
F	0.3
G	0.2

rec(u,i)

Focus on tested items

- C compares to E, A, and D
  - C > A, D, F, G *ដែលចាប់តាំងពីការសំរាប់ទាំងអស់* (4 times)
- B compares to E, A, and D
  - B = A *ដែលមានតម្លៃគ្មាន* (1 time) -  $n''$
  - B > D, F, G (3 time)
- There are  $4+3=7$  times that tested items have higher score than recommended items, so  $n' = 7$
- There are 1 times that tested items have the same score as recommended items, so  $n'' = 1$

# AUC

rec(u,i)

Item	Score
E	0.8
C	0.7
B	0.6
A	0.6
D	0.5
F	0.3
G	0.2

For  $n$  comparisons,

- $n'$  is number of times when the tested items have higher score than the recommended items.
- $n''$  is number of times when the tested items have same score as the recommended items.

$$AUC = \frac{n' + 0.5n''}{n}$$

*↑  
n=10 မျှ AUC > 0.8 ခုရှေ့စ်  
↓  
ပါ့ပါ ကျော် စိုက်*

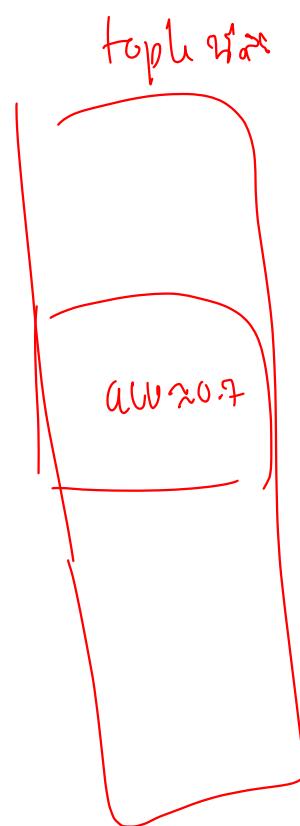
$$= \frac{7 + 0.5 \times 1}{10} = \frac{7.5}{10} = 0.75$$

0.5 is a random recommendation

# AUC ?

top 45%

Item	Score	▼



AUC = ?

AUC ≈ 1

top 51%

Item	Score	▼

AUC = ?

AUC ≈ 0

top 44%

Item	Score	▼

AUC = ?

# Decide to Evaluate? ແລ້ວເຈັດອ່ານວ່າໄວ

Firstly,

- Use **Top-K Precision**
  - check the accuracy of top recommended items
  - close to the real situation when the number of recommend items are limited

Secondly,

- Use **AUC**
  - check the overall performance of an algorithm

ລາສີບຕາມວິດີເປົ້ານັ້ນງວ່າໃຈດໍາຍມາຮູ້

# Collaborative Recommendation

# Collaborative Recommendation

---

- The most prominent approach to generate recommendations
  - used by large, commercial e-commerce sites
  - well-understood, various algorithms and variations exist
  - applicable in many domains (book, movies, DVDs, ..)
- Approach
  - use the "wisdom of the crowd" to recommend items
- Basic assumption and idea
  - Users give ratings to catalog items (implicitly or explicitly)
  - Customers who had similar tastes in the past, will have similar tastes in the future



# User-based Collaborative Filtering

	Item1	Item2	Item3	Item4	Item5
Alice	5	3	4	4	?
User1	3	1	2	3	3
User2	4	3	4	3	5
User3	3	3	1	5	4
User4	1	5	5	2	1

# User-based Collaborative Filtering

---

- Given an "active user" (Alice) and an item I not yet seen by Alice
- The *goal is to estimate Alice's rating for this item*, e.g., by
  - find a set of users (peers) who liked the same items as Alice in the past **and** who have rated item I
  - use, e.g. the average of their ratings to predict, if Alice will like item I
  - do this for all items Alice has not seen and recommend the best-rated

# User-based Collaborative Filtering

---

- Some first questions
  - How do we measure similarity?
  - How many neighbors should we consider?
  - How do we generate a prediction from the neighbors' ratings?

	Item1	Item2	Item3	Item4	Item5
Alice	5	3	4	4	?
User1	3	1	2	3	3
User2	4	3	4	3	5
User3	3	3	1	5	4
User4	1	5	5	2	1

# User-based Collaborative Filtering

ឧបតាមអនុលោះ

- A popular similarity measure in user-based CF: Pearson correlation

a, b : users

$r_{a,p}$  : rating of user a for item p

P : set of items, rated both by a and b

ឧបតាមវិធានសមត្ថភាព

$$sim(a, b) = \frac{\sum_{p \in P} (r_{a,p} - \bar{r}_a)(r_{b,p} - \bar{r}_b)}{\sqrt{\sum_{p \in P} (r_{a,p} - \bar{r}_a)^2} \sqrt{\sum_{p \in P} (r_{b,p} - \bar{r}_b)^2}}$$

Possible similarity values between -1 and 1;  $\bar{r}_a, \bar{r}_b$  = user's average ratings

	Item1	Item2	Item3	Item4	Item5
Alice	5	3	4	4	?
User1	3	1	2	3	3
User2	4	3	4	3	5
User3	3	3	1	5	4
User4	1	5	5	2	1

ការគាំទ្រនៃ Alice ទៅមេដែនល្អ 2 រាល់

sim(alice,user1) = 0.85

sim(alice,user2) = 0.70

sim(alice,user3) = 0.00

sim(alice,user4) = -0.79

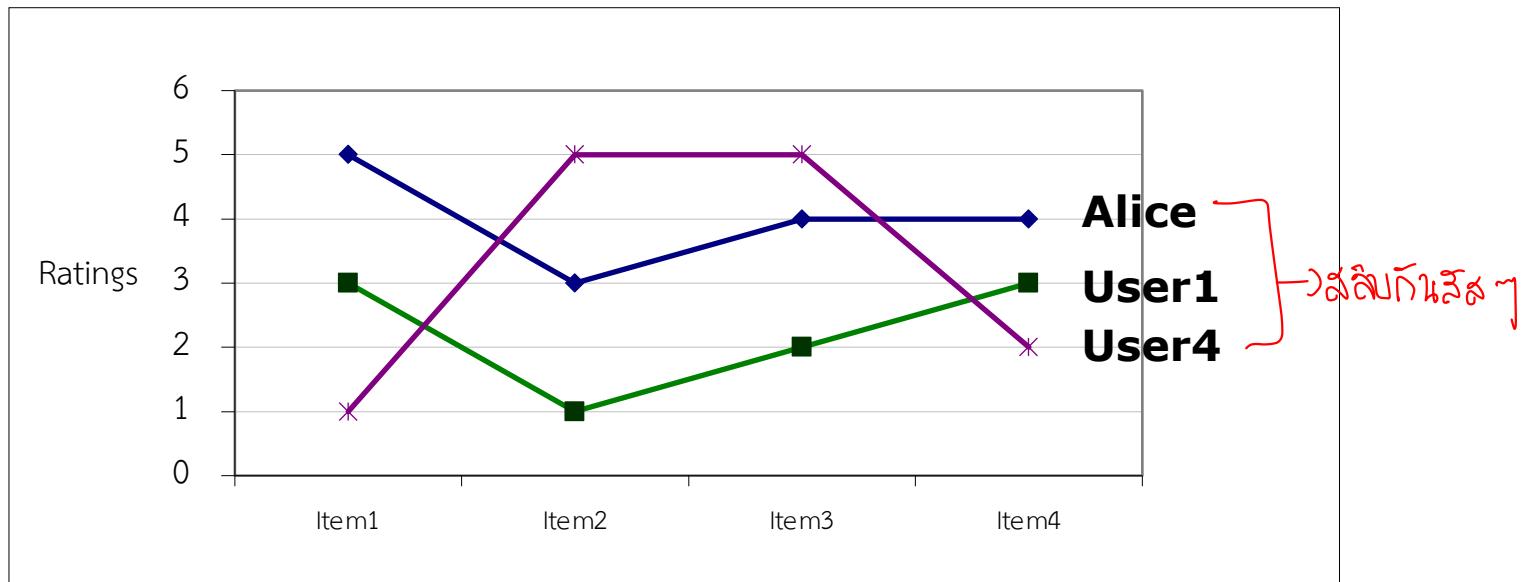
និង pearson នេះ

↓ នឹងពាក្យសារណ៍

# User-based Collaborative Filtering

## Pearson Correlation

- Takes differences in rating behavior into account



- Works well in usual domains, compared with alternative measures
  - such as cosine similarity

# User-based Collaborative Filtering

## Making Prediction

	Item1	Item2	Item3	Item4	Item5	
Alice	5	3	4	4	?	
User1	3	1	2	3	3	$\text{sim}(\text{alice}, \text{user1}) = 0.85$
User2	4	3	4	3	5	$\text{sim}(\text{alice}, \text{user2}) = 0.70$
User3	3	3	1	5	4	$\text{sim}(\text{alice}, \text{user3}) = 0.00$
User4	1	5	5	2	1	$\text{sim}(\text{alice}, \text{user4}) = -0.79$

user item rating ເຊີ່ມ  
user ອັນທຶນ ດາວວະນຸ້ມ rating ພົບ  
 $\text{pred}(a, p) = \bar{r}_a + \frac{\sum_{b \in N} \text{sim}(a, b) * (r_{b,p} - \bar{r}_b)}{\sum_{b \in N} \text{sim}(a, b)}$

$\frac{3 + 1 + 2 + 3 + 3}{5}$   
 $\frac{4 + 3 + 4 + 3 + 5}{5}$

$\text{pred}(\text{alice}, \text{item5}) = \frac{5 + 3 + 4 + 4}{4} + \frac{0.85 * (3 - 2.4) + 0.7 * (5 - 3.8)}{0.85 + 0.7} = 4.87$

# Improving the metrics / prediction function

---

- Not all neighbor ratings might be equally "valuable"
  - Agreement on commonly liked items is not so informative as agreement on controversial items
  - Possible solution: Give more weight to items that have a higher variance
- Value of number of co-rated items
  - Use "significance weighting", by e.g., linearly reducing the weight when the number of co-rated items is low
- Case amplification
  - Intuition: Give more weight to "very similar" neighbors, i.e., where the similarity value is close to 1.
- Neighborhood selection
  - Use similarity threshold or fixed number of neighbors

# Memory-based and model-based approaches

---

- User-based CF is said to be "memory-based"
  - the rating matrix is directly used to find neighbors / make predictions
  - does not scale for most real-world scenarios
  - large e-commerce sites have tens of millions of customers and millions of items
- Model-based approaches
  - based on an offline pre-processing or "model-learning" phase
  - at run-time, only the learned model is used to make predictions
  - models are updated / re-trained periodically
  - large variety of techniques used
  - model-building and updating can be computationally expensive

# Item-based Collaborative Filtering

② គណន់លែងការងារ item ទីផ្សារ

- Basic idea:
  - Use the similarity between items (and not users) to make predictions
- Example:
  - Look for items that are similar to Item5
  - Take Alice's ratings for these items to predict the rating for Item5

	Item1	Item2	Item3	Item4	Item5
Alice	5	3	4	4	?
User1	3	1	2	3	3
User2	4	3	4	3	5
User3	3	3	1	5	4
User4	1	5	5	2	1

# Item-based Collaborative Filtering

---

- Produces better results in item-to-item filtering
  - for some datasets, no consistent picture in literature
- Ratings are seen as vector in n-dimensional space
- Similarity is calculated based on the angle between the vectors
- Adjusted cosine similarity
  - take average user ratings into account, transform the original ratings
  - U: set of users who have rated both items a and b

$$sim(\vec{a}, \vec{b}) = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| * |\vec{b}|}$$

$$sim(a, b) = \frac{\sum_{u \in U} (r_{u,a} - \bar{r}_u)(r_{u,b} - \bar{r}_u)}{\sqrt{\sum_{u \in U} (r_{u,a} - \bar{r}_u)^2} \sqrt{\sum_{u \in U} (r_{u,b} - \bar{r}_u)^2}}$$

} 90% 와 같았습니다~

# Item-based Collaborative Filtering

(1) ຄົກລົງທີ່ມາຈະສຳເນົາ

	Item1	Item2	Item3	Item4	Item5	avg of user, or $\bar{r}_u$
Alice	5	3	4	4	?	4.0
User1	3	1	2	3	3	2.4
User2	4	3	4	3	5	3.8
User3	3	3	1	5	4	3.2
User4	1	5	5	2	1	2.8

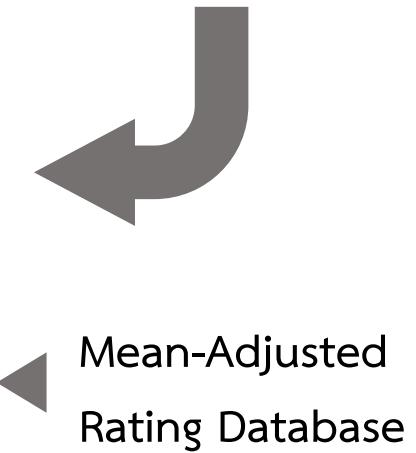
(2) ລົງດັບຕິດ

	Item1	Item2	Item3	Item4	Item5
Alice	$5 - 4.0$	$3 - 4.0$	$4 - 4.0$	$4 - 4.0$	?
User1	$3 - 2.4$	$1 - 2.4$	$2 - 2.4$	$3 - 2.4$	$3 - 2.4$
User2	$4 - 3.8$	$3 - 3.8$	$4 - 3.8$	$3 - 3.8$	$5 - 3.8$
User3	$3 - 3.2$	$3 - 3.2$	$1 - 3.2$	$5 - 3.2$	$4 - 3.2$
User4	$1 - 2.8$	$5 - 2.8$	$5 - 2.8$	$2 - 2.8$	$1 - 2.8$

# Item-based Collaborative Filtering

	Item1	Item2	Item3	Item4	Item5	avg of user, or $\bar{r}_u$
Alice	5	3	4	4	?	4.0
User1	3	1	2	3	3	2.4
User2	4	3	4	3	5	3.8
User3	3	3	1	5	4	3.2
User4	1	5	5	2	1	2.8

	Item1	Item2	Item3	Item4	Item5
Alice	1.0	-1.0	0.0	0.0	?
User1	0.6	-1.4	-0.4	0.6	0.6
User2	0.2	-0.8	0.2	-0.8	1.2
User3	-0.2	-0.2	-2.2	1.8	0.8
User4	-1.8	2.2	2.2	-0.8	-1.8



# Item-based Collaborative Filtering

Mean-Adjusted  
Rating Database



	Item1	Item2	Item3	Item4	Item5
Alice	1.0	-1.0	0.0	0.0	?
User1	0.6	-1.4	-0.4	0.6	0.6
User2	0.2	-0.8	0.2	-0.8	1.2
User3	-0.2	-0.2	-2.2	1.8	0.8
User4	1.8	2.2	2.2	-0.8	-1.8

ผู้ใช้ที่มีความคล้ายกัน item 1

# Item-based Collaborative Filtering

Mean-Adjusted  
Rating Database



	Item1	Item2	Item3	Item4	Item5
Alice	1.0	-1.0	0.0	0.0	?
User1	0.6	-1.4	-0.4	0.6	0.6
User2	0.2	-0.8	0.2	-0.8	1.2
User3	-0.2	-0.2	-2.2	1.8	0.8
User4	-1.8	2.2	2.2	-0.8	-1.8

$$sim(a, b) = \frac{\sum_{u \in U} (r_{u,a} - \bar{r}_u)(r_{u,b} - \bar{r}_u)}{\sqrt{\sum_{u \in U} (r_{u,a} - \bar{r}_u)^2} \sqrt{\sum_{u \in U} (r_{u,b} - \bar{r}_u)^2}}$$

ခုပေါင်းစုံ  
ရှုပါနီ

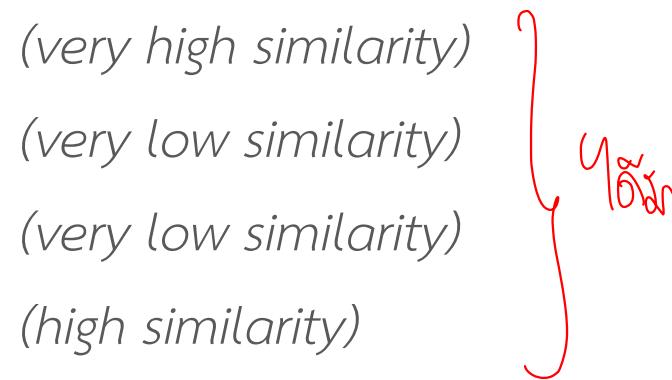
$$sim(i1, i5) = \frac{0.6 \times 0.6 + 0.2 \times 1.2 + (-0.2) \times 0.8 + (-1.8) \times (-1.8)}{\sqrt{0.6^2 + 0.2^2 + (-0.2)^2 + (-1.8)^2} \times \sqrt{0.6^2 + 1.2^2 + 0.8^2 + (-1.8)^2}} = 0.8$$

# Item-based Collaborative Filtering

---

- after calculation

- $sim(i5, i1) = +0.8$  (*very high similarity*)
- $sim(i5, i2) = -0.9$  (*very low similarity*)
- $sim(i5, i3) = -0.8$  (*very low similarity*)
- $sim(i5, i4) = +0.4$  (*high similarity*)



# Item-based Collaborative Filtering

- choose similar items

$\bullet \sim(i_5, i_1) = +0.8$

$\bullet \sim(i_5, i_2) = -0.9$

$\bullet \sim(i_5, i_3) = -0.8$

$\bullet \sim(i_5, i_4) = +0.4$

	Item1	Item2	Item3	Item4	Item5
Alice	5	3	4	4	?
User1	3	1	2	3	3
User2	4	3	4	3	5
User3	3	3	1	5	4
User4	1	5	5	2	1

item ที่ Alice ชอบ ก็ user คนนั้นเดียวกัน

$$pred(u, p) = \frac{\sum_{i \in ratedItem(u)} sim(i, p) \times r_{u,i}}{\sum_{i \in ratedItem(u)} sim(i, p)}$$

$\text{work with}$

$$pred(alice, i5) = \frac{0.8 \times 5 + 0.4 \times 4}{0.8 + 0.4} = \frac{5.6}{1.2} = 4.7$$



	Item5
Alice	4.7

# Item-based Collaborative Filtering

- choose similar items
  - $sim(i5, i1) = +0.8$
  - $sim(i5, i2) = -0.9$
  - $sim(i5, i3) = -0.8$
  - $sim(i5, i4) = +0.4$

	Item1	Item2	Item3	Item4	Item5
Alice	5	3	4	4	?
User1	3	1	2	3	3
User2	4	3	4	3	5
User3	3	3	1	5	4
User4	1	5	5	2	1

$$pred(u, p) = \frac{\sum_{i \in ratedItem(u)} sim(i, p) \times r_{u,i}}{\sum_{i \in ratedItem(u)} sim(i, p)}$$

$$pred(alice, i5) = \frac{0.8 \times 5}{0.8} = 5$$



	Item5
Alice	5

# Pre-processing for item-based filtering

---

- Item-based filtering does not solve the scalability problem itself
- Pre-processing approach by Amazon.com (in 2003)
  - Calculate all **pair-wise item similarities in advance**
  - The neighborhood to be used at run-time is typically rather small, because only items are taken into account which the user has rated
  - Item similarities are supposed to be more stable than user similarities
- Memory requirements
  - Up to  $N^2$  **pair-wise similarities** to be memorized ( $N$  = number of items) in theory
  - In practice, this is significantly lower (items with no co-ratings)
  - Further reductions possible
    - **Minimum threshold for co-ratings** (items, which are rated at least by  $n$  users)
    - **Limit the size of the neighborhood** (might affect recommendation accuracy)

# More on ratings

---

- Pure CF-based systems only **rely on the rating matrix**
- Explicit ratings
  - Most commonly used (1 to 5, 1 to 7 Likert response scales)
  - Research topics
    - "Optimal" granularity of scale; indication that 10-point scale is better accepted in movie domain
    - Multidimensional ratings (multiple ratings per movie)
  - Challenge
    - Users not always willing to rate many items; **sparse rating matrices**
    - How to stimulate users to rate more items?
- Implicit ratings
  - **clicks, page views, time spent** on some page, **demo downloads ...**
  - Can be used in addition to explicit ones; question of correctness of interpretation

# Data sparsity problems

---

- Cold start problem
  - How to recommend **new items**? What to recommend to **new users**?
- Straightforward approaches
  - Ask/force users to rate a set of items
  - Use another method (e.g., **content-based**, **demographic** or simply non-personalized) in the initial phase
- Alternatives
  - Use **better algorithms** (beyond nearest-neighbor approaches), for example:
    - In nearest-neighbor approaches, the set of sufficiently similar neighbors might be too small to make good predictions
    - Assume "**transitivity**" of neighborhoods

សំណើអាមេរិក

# SVD (Singular Value Decomposition)

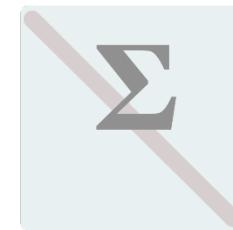
មាត្រូវការ



=



គ្រប់គ្រងការបញ្ចូល



គ្រប់គ្រងការបញ្ចូល



data matrix

left singular  
vectors

diagonal of  
singular  
values

right singular  
vectors

ត្រូវបានបញ្ចូលជាបន្ទាល់  
ប៉ុណ្ណោះ

# PCA

คุณลุงปูนเป็ด ขยัน

## Principal Component Analysis (PCA) ภาษาไทย Update Version



คุณลุง

<https://www.youtube.com/watch?v=1qZyWTTw0mI>



# SVD

---

- Example

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix} \times \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$

# Matrix Factorization

SVD

$$M_k = U_k \times \Sigma_k \times V_k^T$$

*compress column*

U <sub>k</sub>	Dim1	Dim2
Alice	0.47	-0.30
Bob	-0.44	0.23
Mary	0.70	-0.06
Sue	0.31	0.93

V <sub>k</sub> <sup>T</sup>	TERMINATOR	DIE HARD	TWINS	EAT PRAY LOVE	Pretty Woman
Dim1	-0.44	-0.57	0.06	0.38	0.57
Dim2	0.58	-0.66	0.26	0.18	-0.36

- Prediction:**  $\hat{r}_{ui} = \bar{r}_u + U_k(Alice) \times \Sigma_k \times V_k^T(EPL)$   
 $= 3 + 0.84 = 3.84$

ຈົກສ້າງຂອງ  $\Sigma_k$   $\rightarrow$  ຈົກສ້າງຂອງ  $V_k^T$

$\Sigma_k$	Dim1	Dim2
Dim1	5.63	0
Dim2	0	3.23

2008: Factorization meets the neighborhood: a multifaceted collaborative filtering model, Y. Koren, ACM SIGKDD

---

- Merges neighborhood models with latent factor models
- Latent factor models
  - good to capture weak signals in the overall data
- Neighborhood models
  - good at detecting strong relationships between close items
- Combination in one prediction single function
  - Local search method such as stochastic gradient descent to determine parameters
  - Add penalty for high values to avoid over-fitting

$$\hat{r}_{ui} = \mu + b_u + b_i + p_u^T q_i$$

$$\min_{p^*, q^*, b^*} \sum_{(u,i) \in K} (r_{ui} - \mu - b_u - b_i - p_u^T q_i)^2 + \lambda (\|p_u\|^2 + \|q_i\|^2 + b_u^2 + b_i^2)$$

# Collaborative Filtering: Pros & Cons

---

- Pros:  *sin yertosms lib*
  - well-understood, works well in some domains, no knowledge engineering required
- Cons:  *nosm community, user, sparsity*
  - requires user community, sparsity problems, no integration of other knowledge sources, no explanation of results

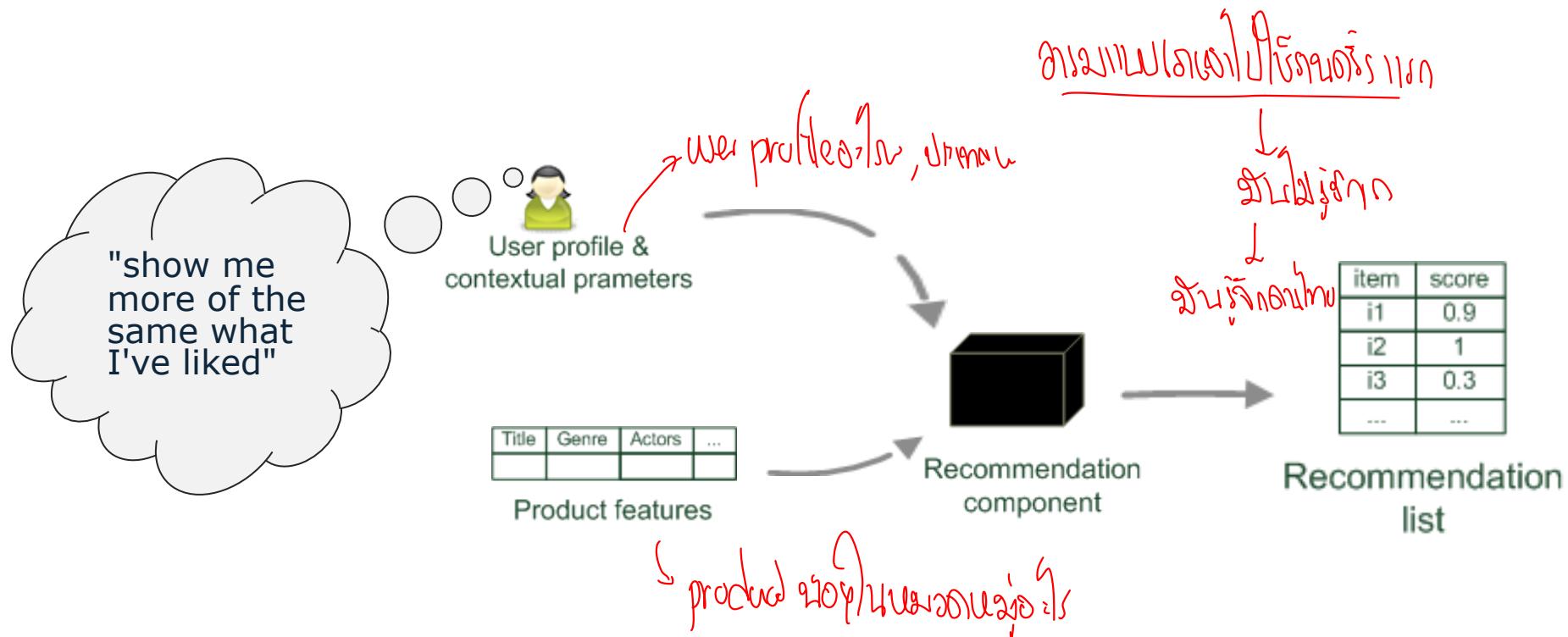
# Content-based Recommendation

# Content-based recommendation

---

- While CF – methods do not require any information about the items,
  - it might be reasonable to exploit such information; and
  - recommend fantasy novels to people who liked fantasy novels in the past
- What do we need:
  - some information about the available items such as the genre ("content")
  - some sort of *user profile* describing what the user likes (the preferences)
- The task:
  - learn user preferences
  - locate/recommend items that are "similar" to the user preferences

# Content-based recommendation



# What is the "content"?

---

- Most CB-recommendation techniques were applied to recommending text documents.
  - Like web pages or newsgroup messages for example.
- Content of items can also be represented as text documents.
  - With textual descriptions of their basic characteristics.
  - Structured: Each item is described by the same set of attributes

Title	Genre	Author	Type	Price	Keywords
The Night of the Gun	Memoir	David Carr	Paperback	29.90	Press and journalism, drug addiction, personal memoirs, New York
The Lace Reader	Fiction, Mystery	Brunonia Barry	Hardcover	49.90	American contemporary fiction, detective, historical
Into the Fire	Romance, Suspense	Suzanne Brockmann	Hardcover	45.90	American fiction, murder, neo-Nazism

# Content representation and item similarities

## ▪ Item representation

Title	Genre	Author	Type	Price	Keywords
The Night of the Gun	Memoir	David Carr	Paperback	29.90	Press and journalism, drug addiction, personal memoirs, New York
The Lace Reader	Fiction, Mystery	Brunonia Barry	Hardcover	49.90	American contemporary fiction, detective, historical
Into the Fire	Romance, Suspense	Suzanne Brockmann	Hardcover	45.90	American fiction, murder, neo-Nazism

- Simple approach
  - Compute the similarity of an unseen item with the user profile based on the keyword overlap (e.g. using the **Jaccard index**)

$$\frac{|\text{keywords}(A) \cap \text{keywords}(B)|}{|\text{keywords}(A) \cup \text{keywords}(B)|}$$

↑  
index จัคการ์ด

# Jaccard Index

जैकर्ड इंडेक्स

$$Jaccard\ Index = \frac{|keywords(A) \cap keywords(B)|}{|keywords(A) \cup keywords(B)|}$$

- e.g.
  - Book A contains “Thailand, Bangkok, Market, Tourist, Temple, Cuisine”.
  - Book B contains “Thailand, Bangkok, Temple, Museum”.

$$Jaccard\ Index = \frac{3}{7}$$

# Limitations → பிரைவர் (பார்களுக்கு உத்திரவுகள் கூடாது)

---

- Keywords alone may not be sufficient to judge quality/relevance of a document or web page
  - up-to-date-ness, usability, aesthetics, writing style
  - content may also be limited / too short
  - content may not be automatically extractable (multimedia)
- Ramp-up phase required
  - Web 2.0: Use other sources to learn the user preferences
- Overspecialization
  - Or: too similar news items
- Pure content-based systems are rarely found in commercial environments

# Knowledge-based Recommendation

# Knowledge-Based

content base ការងារ

ផ្លូវ យោងទៅ gadget ដែលបាន → រាយការណ៍ទី 100% នៃ item



Knowledge-based: "Tell me what fits based on my needs"

Title	Genre	Actors	...

Product features



Recommendation component

item	score
i1	0.9
i2	1
i3	0.3
...	...

Recommendation list



Knowledge models

ផ្សេងៗអ្នកមើល និងចំណាំជាបុន្ណោះ

# Why do we need knowledge-based recommendation?

---

- Products with low number of available ratings



- Time span plays an important role
  - five-year-old ratings for computers
  - user lifestyle or family situation changes
- Customers want to define their requirements explicitly
  - "the color of the car should be black"

# Knowledge-based recommender systems

①

- Constraint-based *ក្រឡិនសម្រាប់ផ្តល់ព័ត៌មានទៅ item*
- based on explicitly defined set of recommendation rules
- fulfill recommendation rules

②

- Case-based *បានរាយដូចជាលទ្ធផល*
- based on different types of similarity measures
- retrieve items that are similar to *specified requirements*
- Both approaches → *conversational* recommendation process
  - *users specify* the requirements
  - systems try to identify solutions
  - if no solution can be found, *users change requirements*

# Constraint-based recommendation problem

---

- Select items from this catalog that match the user's requirements

id	price(€)	mpix	opt-zoom	LCD-size	movies	sound	waterproof
P <sub>1</sub>	148	8.0	4×	2.5	no	no	yes
P <sub>2</sub>	182	8.0	5×	2.7	yes	yes	no
P <sub>3</sub>	189	8.0	10×	2.5	yes	yes	no
P <sub>4</sub>	196	10.0	12×	2.7	yes	no	yes
P <sub>5</sub>	151	7.1	3×	3.0	yes	yes	no
P <sub>6</sub>	199	9.0	3×	3.0	yes	yes	no
P <sub>7</sub>	259	10.0	3×	3.0	yes	yes	no
P <sub>8</sub>	278	9.1	10×	3.0	yes	yes	yes

- User's requirements can, for example, be
  - "the price should be lower than 300 €"
  - "the camera should be suited for sports photography"

# Summary

---

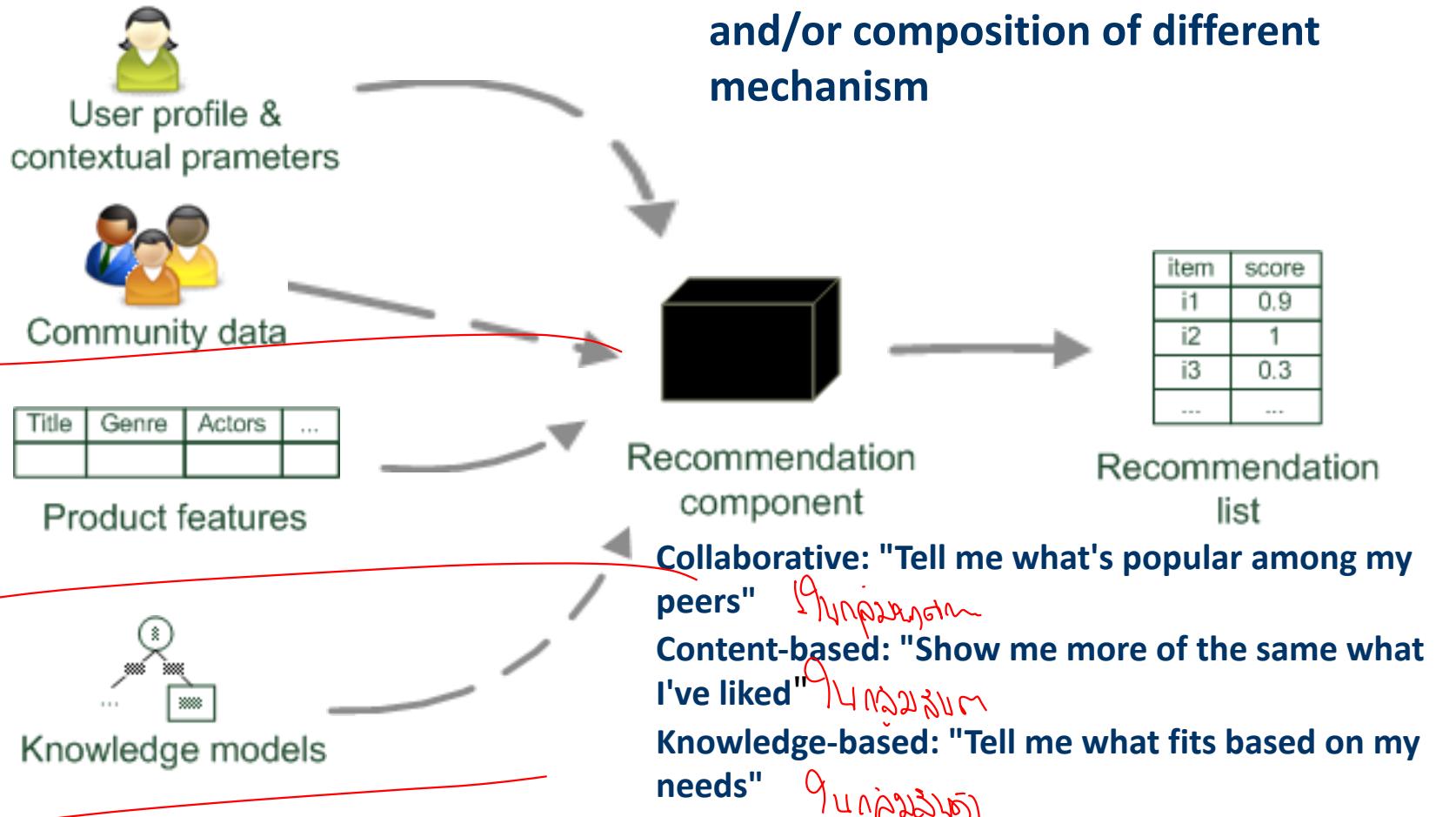
- Knowledge-based recommender systems
  - constraint-based
  - case-based
- Limitations
  - cost of knowledge acquisition
    - from domain experts
    - from users
    - from web resources
  - accuracy of preference models
    - very fine granular preference models require many interaction cycles
    - collaborative filtering models preference implicitly
  - independence assumption can be challenged
    - preferences are not always independent from each other

ទិន្នន័យទាមព័ត៌មាន  
xml, sql, resourc

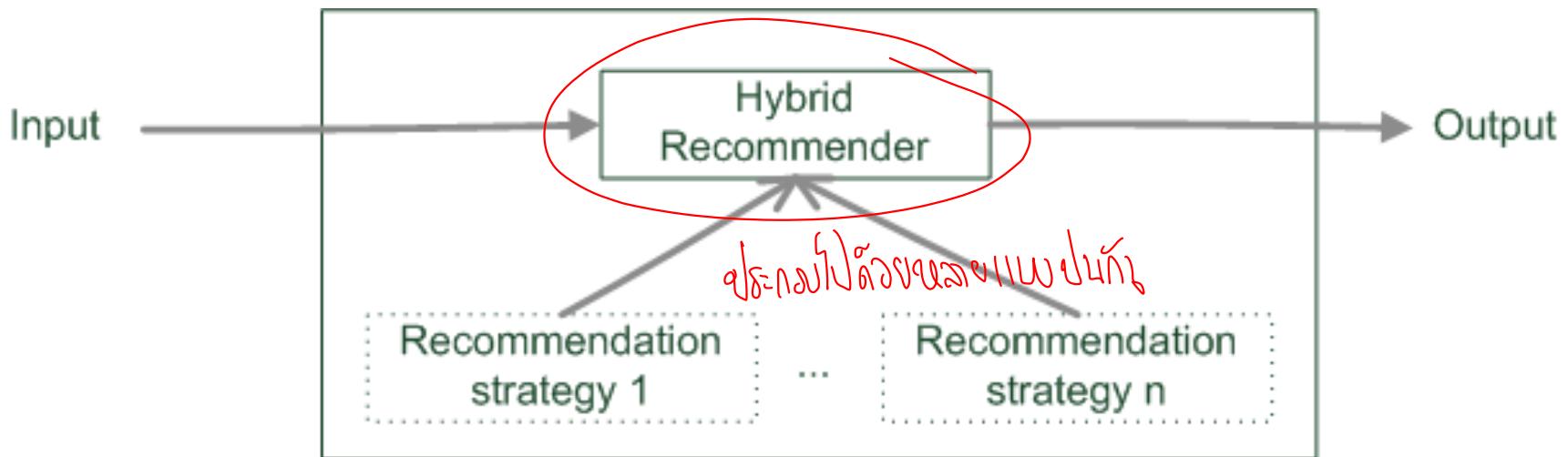
# Hybrid Recommender Approach

ແອັດເລືອກເຫຼົ່າໃຈ ທີ່  
ມີຜົນໄດ້ທາງໆ ສິ້ນສົ່ງ  
③ ແພວິບແນ

# Hybrid Recommender



# Abstract



# Weighted

---

- The scores (or votes) of several recommendation techniques are combined together to produce a single recommendation.

(ໃຊ້ໄວ້ deep learning  
ເພື່ອຮັດງານ)

**Recommendation Score =**

$$w1 \times \text{Score\_from\_Algorithm\_1} + w2 \times \text{Score\_from\_Algorithm\_2} + \dots + wN \times \text{Score\_from\_Algorithm\_N}$$

$\sum_{i=1}^n w_i = 1$

## ② ແຈກ ຢຸບພາຍໃຈທີ່ໃຊ້ algorithm (ເລືອດ) \*

# Switching

---

- The system switches between recommendation techniques depending on the current situation.

```
If (data is AAA ) {  
    use Algorithm_A ;  
}  
else if (data is BBB ){  
    use Algorithm_B ;  
}  
else{ use Algorithm_C ;
```

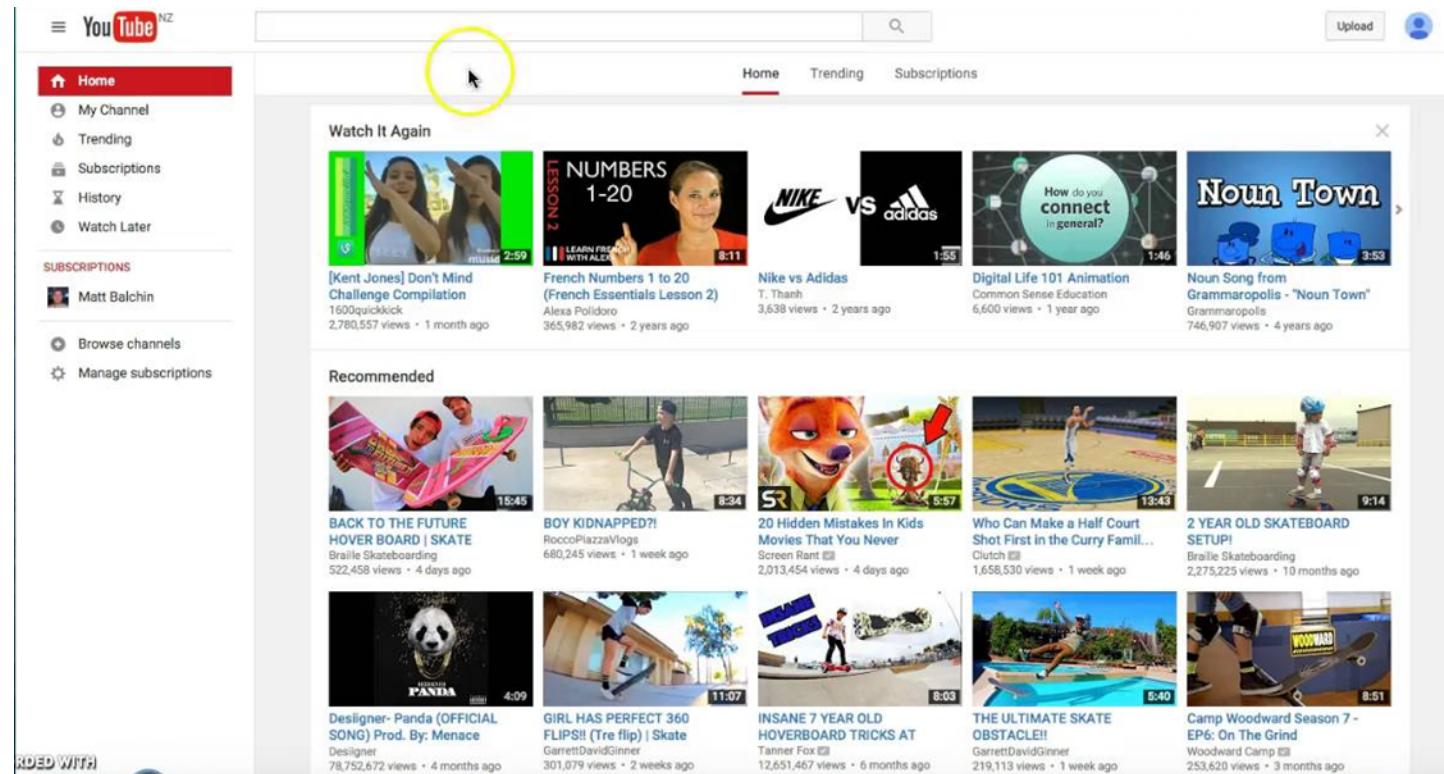
③ ເລັກສະໝັກຮວມກົງໂຕ (ເສີມທີ່ນຳນົດ)

Mixed

↳ ໄສ້ດູວ່າຕະຫຼາມນີ້ຢືນໃຈຂອງ algorithm ທ່ານ

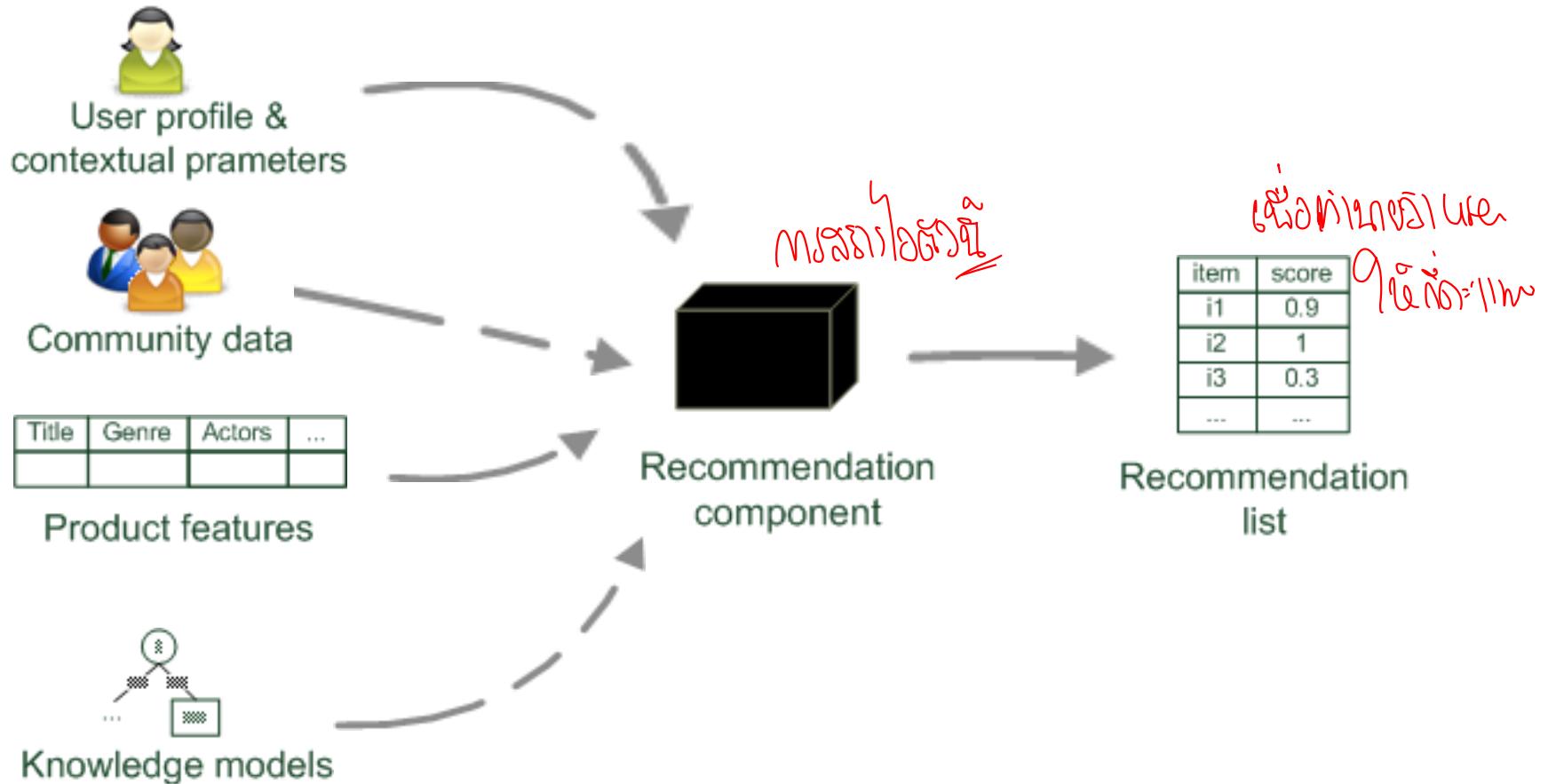
top  
ກິດ → SOL

- Recommendations from several different recommenders are presented at the same time

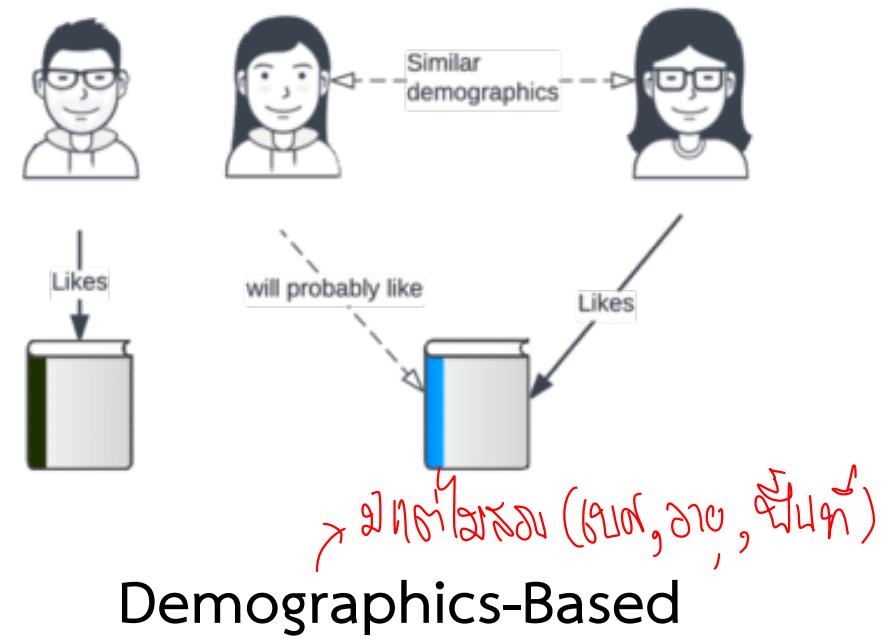
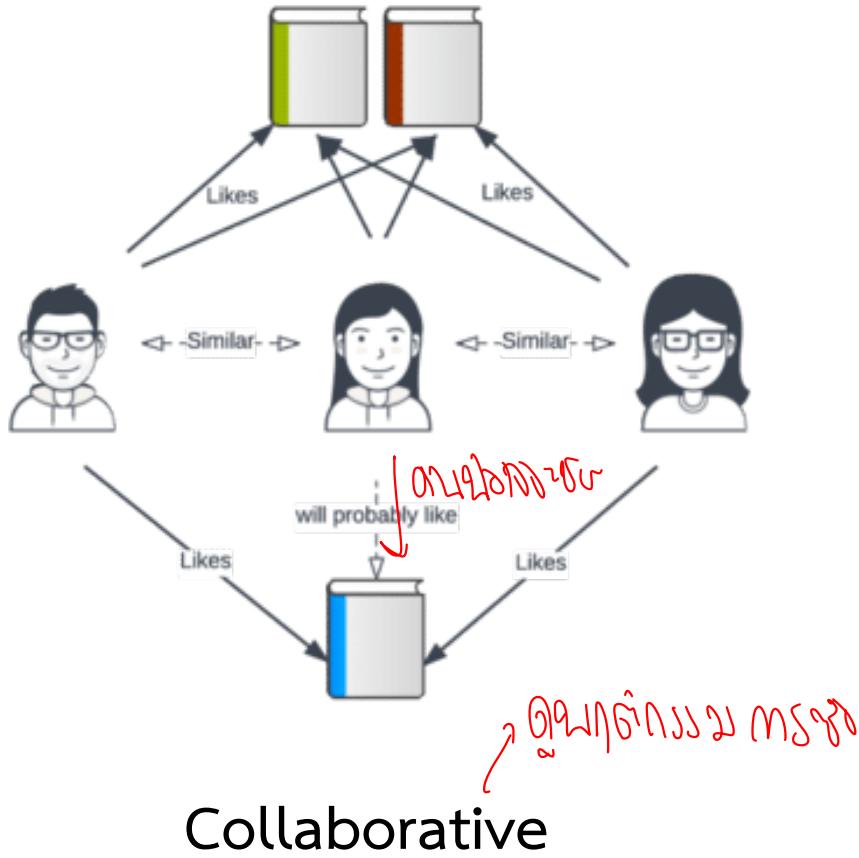


# Summary

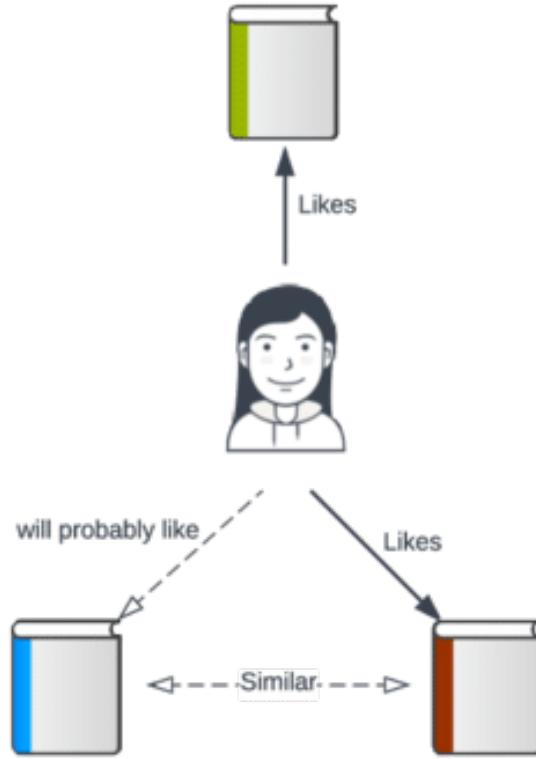
# Summary



# Summary

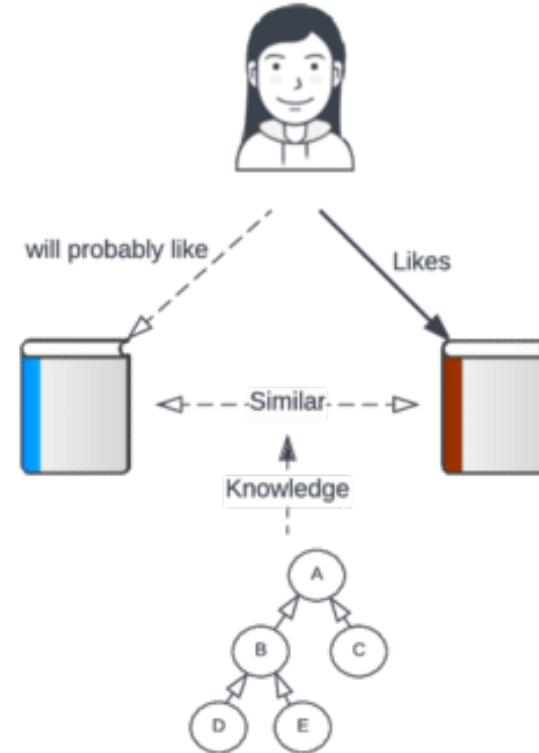


# Summary



Content-Based

ទូរស័ព្ទការបង្កើតបច្ចុប្បន្ន



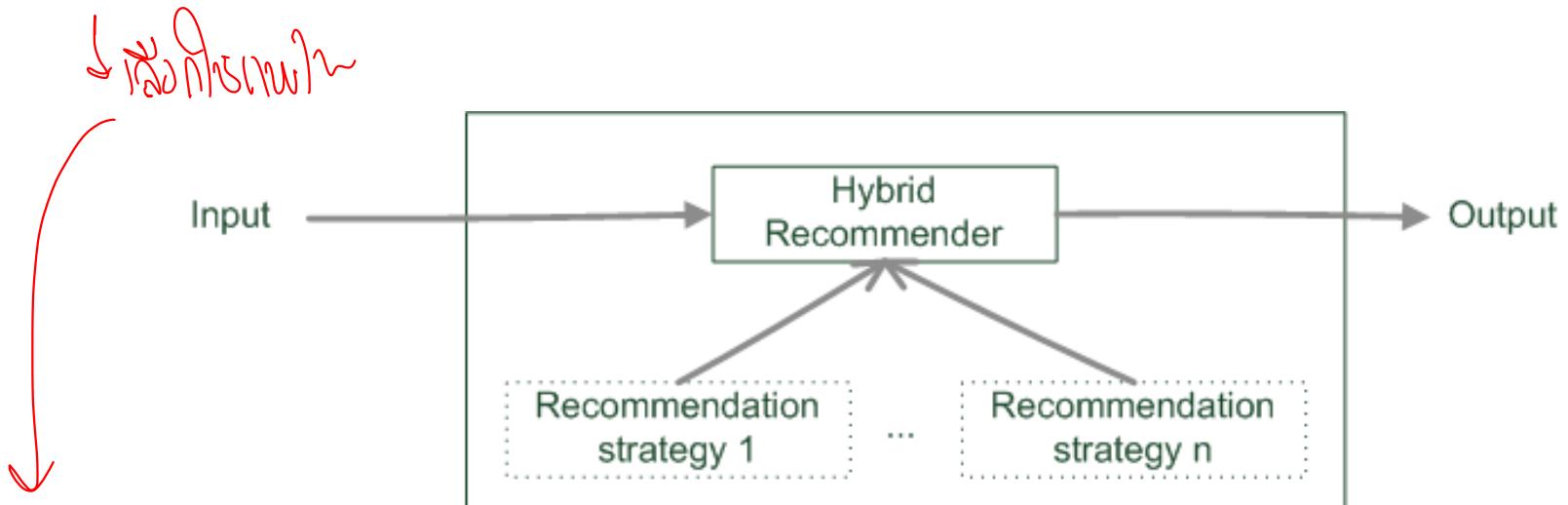
Knowledge-Based

គោលដៅទំនួរ

# Summary

---

## Hybrid Recommender Systems



- Weighted
- Switching
- Mixed

“

Machine learning allows us to build software solutions that exceed human understanding and shows us how AI can innervate every industry.

”

Steve Jurvetson

つづく

# វិធានសាស្ត្រការអនុវត្តការណ៍នៃការពារកម្មបូល់ (PCA)

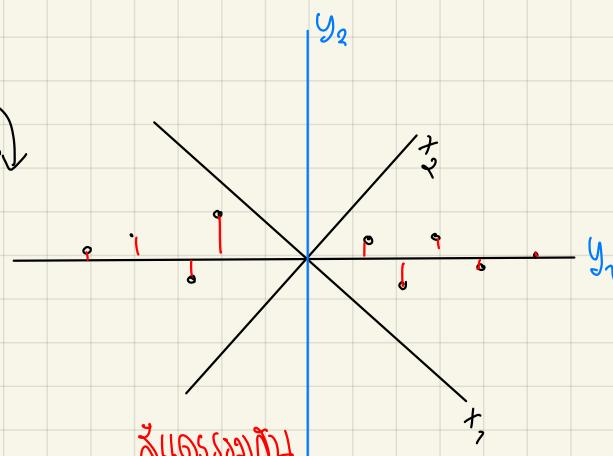
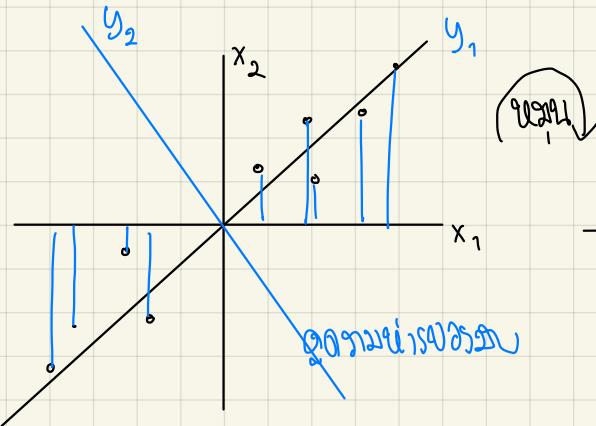
## Principal Components Analysis (PCA)

↪ ការបង្កើតរឹងរាល់ការពារកម្មបូល់ដែលបានបង្ហាញថា ការពារកម្មបូល់នឹងបានបង្កើតឡើងជាបន្ទាន់

↪ ផ្តល់ភាពរូបរាងនៃការពារកម្មបូល់ដែលបានបង្ហាញថាអាចបង្កើតឡើងជាបន្ទាន់

↪ ការបង្កើតរឹងរាល់ការពារកម្មបូល់នឹងបានបង្ហាញថាអាចបង្កើតឡើងជាបន្ទាន់

↪ និង complexity and storage នេះ



និងទូទាត់ការពារកម្មបូល់ ទីក្រុងក្នុង  $(x_i, y_i)$

$\downarrow$   
 $\bar{x}, \bar{y}$  បាន

plot  $(x_i, y_i)$   $\rightarrow$  រួមគាំ  $\rightarrow$  និងលក្ខណៈ  $\times$  និង  $\approx -1 - \approx 1$   
ការចូលរួមគាំ  $\rightarrow$  និងលក្ខណៈ  $\times$  និង  $\approx -1 - \approx 1$   
ដែល

↪ សម្រាប់ ① ការសម្រួលនិងបញ្ជីការណា

យានីនិង  $x, y$

ការចូលរួមគាំ

② ការកើតឡើងនូវការពារកម្មបូល់ (Covariance matrix)

= ការកើតឡើងការពារកម្មបូល់

③ ការកើតឡើងនូវការពារកម្មបូល់ (Eigenvalue and Eigenvector)

ដែល  $C$  ជាការពារកម្មបូល់

④ ការបង្កើតរឹងរាល់ការពារកម្មបូល់ដូចតែបានការពារកម្មបូល់

បានសំណងបន្ថែម

$\uparrow$ , ជាន់

① ការបង្កើតរឹងរាល់ការពារកម្មបូល់

ដែលមិនមែនការបង្ហាញ

$(x, y) \rightarrow$  ឱ្យតាមរឹងរាល់  $x, y$  ការពារកម្មបូល់  $(x - \bar{x}, y - \bar{y})$  ដែលការរឹងរាល់នេះមិនមែន

② ការកើតឡើង  $C \rightarrow C = \frac{1}{M-1} A A^T = \frac{1}{M-1} \sum_{i=1}^n a_i a_i^T$ ; ពន្លឺទីនាមុនក្នុង

ខ្លួន ៩x១០ - ១០x២ = ២x២ ដើម្បីស្វែនការពារកម្មបូល់

③ ឯក Eigenvalue  $\rightarrow$  ការបង្កើតរឹងរាល់ការពារកម្មបូល់

$C y_i = \lambda_i y_i$   $\rightarrow$  ការបង្កើតរឹងរាល់ការពារកម្មបូល់  
ដើម្បីស្វែនការពារកម្មបូល់

$$(C - \lambda I) y_i = 0$$

ដើម្បីស្វែនការពារកម្មបូល់

↪ សំណងលាស់ដើម្បីបង្ហាញការពារកម្មបូល់

$$\begin{bmatrix} C_{11} - \lambda & C_{12} \\ C_{21} & C_{22} - \lambda \end{bmatrix} y_1 = 0$$

$$(axd) - (bxc) = 0 ; \det$$

$$\lambda \text{ ជាបន្ទាន់ } \text{ eigenvalue}$$

၃ ສາမັກດຽວທີ່ປະຫຼອດເຈົ້າມາຢືນຢັນເລີ່ມຕົ້ນ → ອຳຈະເຫັນວ່າ