

# Machine Learning and Its applications (and beyond)

Hisashi Kashima

Department of Mathematical Informatics  
The University of Tokyo



# What is Machine Learning?

# Machine learning as a buzzword:

## Increasing popularity of machine learning

---

- After the IT revolution, “how to exploit data” is increasingly more important than “how to store data”
  - Companies are trying to position data analytics as foundations of competitiveness
- “Machine learning” as a buzzword
  - The rise of “big data” (another buzzword)
  - Data scientist is “the sexiest job in the 21st century”

# What is machine learning? :

## Machine learning is a data analytics

---

- Originally a branch of artificial intelligence
  - Computer programs that “learns” from experience
  - Based on logical inference
- Rise of “statistical” machine learning
  - Successes in bioinformatics, natural language processing, and other business areas
  - Victory of IBM’s Watson QA system

# What can machine learning do?:

## Prediction and discovery

---

- Two categories of the use of machine learning:
  1. Prediction (supervised learning)
    - “What will happen in future data?”
    - Given past data, predict about future data
  2. Discovery (unsupervised learning)
    - “What is happening in data in hand?”
    - Given past data, find insights in them

# Formulations of machine learning problems:

## Supervised and unsupervised

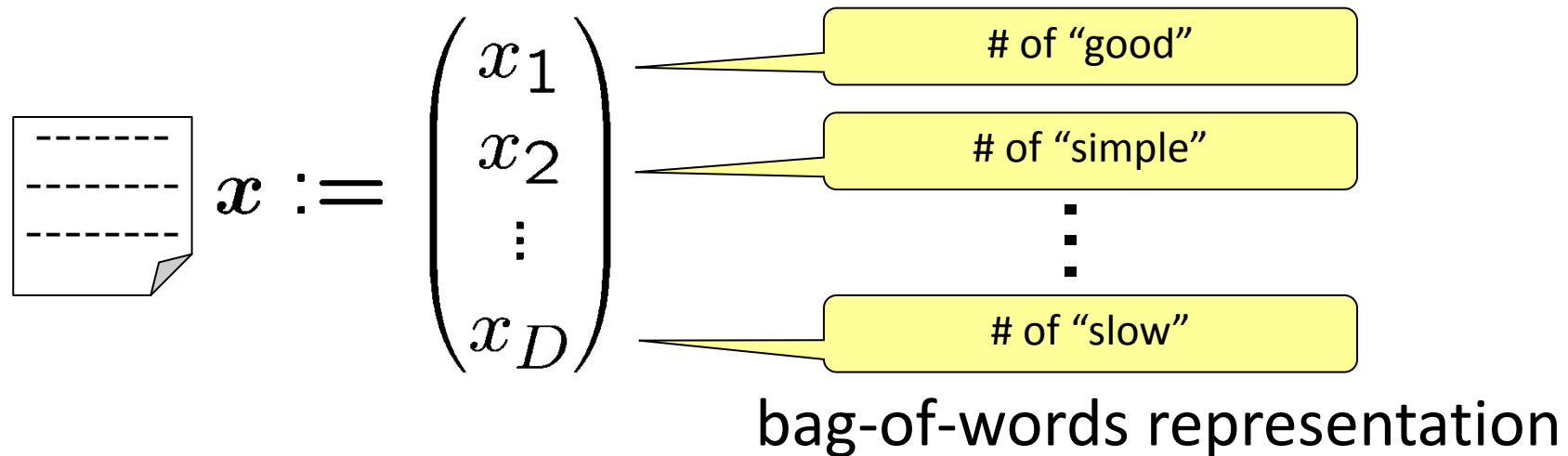
---

- Learning system as a function  $f: \mathbf{x} \rightarrow y$ 
  - $\mathbf{x} \in \mathbb{R}^D$  : input as real vector
  - $y \in \{1, 2, \dots, C\}$  : discrete output
- We want  $f \Rightarrow$  Learn from data
- Two machine learning problems:
  1. Supervised learning: have access to input and output pairs
    - $\{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$  :  $N$  in/output pairs
  2. Unsupervised learning: have access only to input data
    - $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$  :  $N$  inputs

# An application of supervised learning:

## Sentiment analysis

- Classify the polarity of a given text
  - $y \in \{+, -\}$  : Whether or not a blog post  $x$  favors a product
- $x$  is defined by using words appearing in the text



- Note: design of the feature vector is left to uses

# Various applications of machine learning: From online shopping to system monitoring

---

## ■ Marketing

- Product recommendation
- Sentiment analysis on Web
- Advertisement optimization

## ■ Finance

- Credit risk management
- Fraud detection

## ■ Bio/healthcare

- Medical diagnosis
- Gene recognition

## ■ Web

- Web search
- Spam detection
- SNS

## ■ Multimedia

- Voice recognition
- Face/object recognition

## ■ System monitoring

- fault diagnosis
- Intrusion detection



# Recommender Systems

# Recommender systems:

## Personalized information filter

- Amazon offers a list of products I am likely to buy (based on my purchase history)

amazon.co.jp    マイストア    Amazonポイント    ギフト券    タイムセール    出品サービス    ヘルプ

カテゴリーからさがす    おもちゃ    検索    こんにちは、アカウント    今すぐ体験プライム    カート    ほしい物リスト

マイストア    マイページ    お客様へのおすすめ    商品をお評価する    おすすめ商品を正確にする    プロフィール    詳しくはこちら

マイストア > おすすめ商品 > おもちゃ

これらのおすすめ商品は、[過去にお持ちの商品](#)などに基づいています。

表示: [すべて](#) | [ニューリリース情報](#) | [まもなく発売](#)    次のページ

-  **レゴ デュプロ 大きなどうぶつえん 6157**  
レゴ (2012/1/19)  
おすすめ度: ★★★★★ (2)  
在庫あり  
参考価格: ¥13,660  
価格: ¥13,000  
新品の出品: 12 ¥13,000より  
☐ 持っています    ☐ 興味がありません    ☒ ★★★★★ この商品をお評価する  
レゴ デュプロ 基礎板ミニ (赤・緑・黄) 4632を購入されたお客様におすすめします (おすすめの商品に反映させる商品の設定を変更するにはこちら)
-  **レゴ 基本セット 基礎板(青色) 620**  
レゴ (2010/2/9)  
おすすめ度: ★★★★★ (30)  
在庫あり  
参考価格: ¥1,060  
価格: ¥697  
新品の出品: 16 ¥697より  
☐ 持っています    ☐ 興味がありません    ☒ ★★★★★ この商品をお評価する  
レゴ 基本セット ブロック タイヤセット 6118などを購入されたお客様におすすめします (おすすめの商品に反映させる商品の設定を変更するにはこちら)
-  **レゴ デュプロ 基本ブロック (XL) 6176**  
レゴ (2008/1/26)  
おすすめ度: ★★★★★ (35)  
在庫あり  
参考価格: ¥3,990  
価格: ¥2,499  
新品の出品: 14 ¥2,499より  
☐ 持っています    ☐ 興味がありません    ☒ ★★★★★ この商品をお評価する  
レゴ デュプロ 基礎板ミニ (赤・緑・黄) 4632を購入されたお客様におすすめします (おすすめの商品に反映させる商品の設定を変更するにはこちら)

おもしろ商品  
おもちゃ  
おもちゃ・雑貨・手品  
お絵かき・ぬいど・シール  
きせかえ人形・ハウス  
ぬいぐるみ  
ままごと・ごっこ遊び  
アクションスポーツ玩具  
クッキング玩具  
ゲーム  
コスプレ・アクセサリ  
パズル  
パーティー小物  
ブロック  
プラモデル・模型  
メイキング玩具  
ラジコン  
ロボット・ソフビ人形  
乗用玩具・三輪車  
変身・なりきりグッズ  
美術用品  
学習・科学・工作  
工芸・民芸品  
手芸・画材  
文具・学用品  
楽器玩具  
赤ちゃん・知育玩具  
電子玩具・キッズ家電  
電車・ミニカー・乗り物

# Ubiquitous recommender systems:

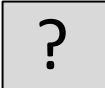
## Recommender systems are present everywhere

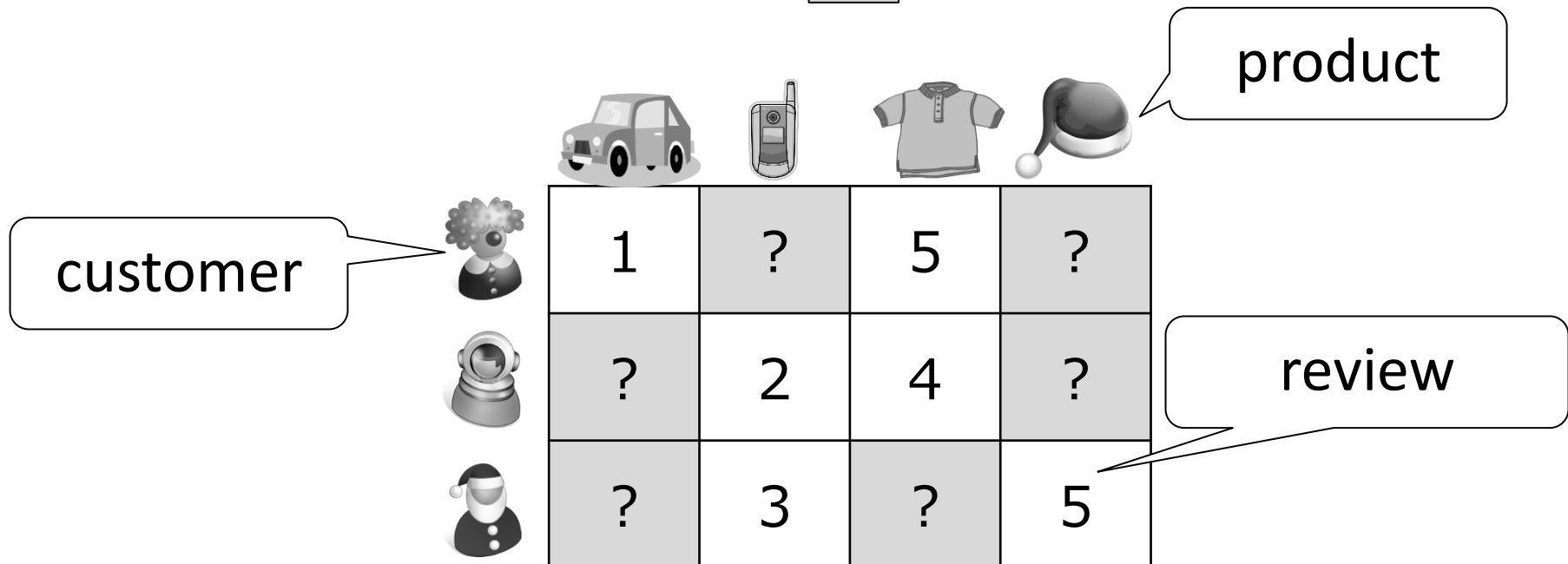
- A major battlefield of machine learning algorithms
  - Netflix challenge (with \$100 million prize)
- Recommender systems are present everywhere:
  - Product recommendation in online shopping stores
  - Friend recommendation on SNSs
  - Information recommendation (news, music, ...)
  - ...



# A formulation of recommendation problem:

## Matrix completion

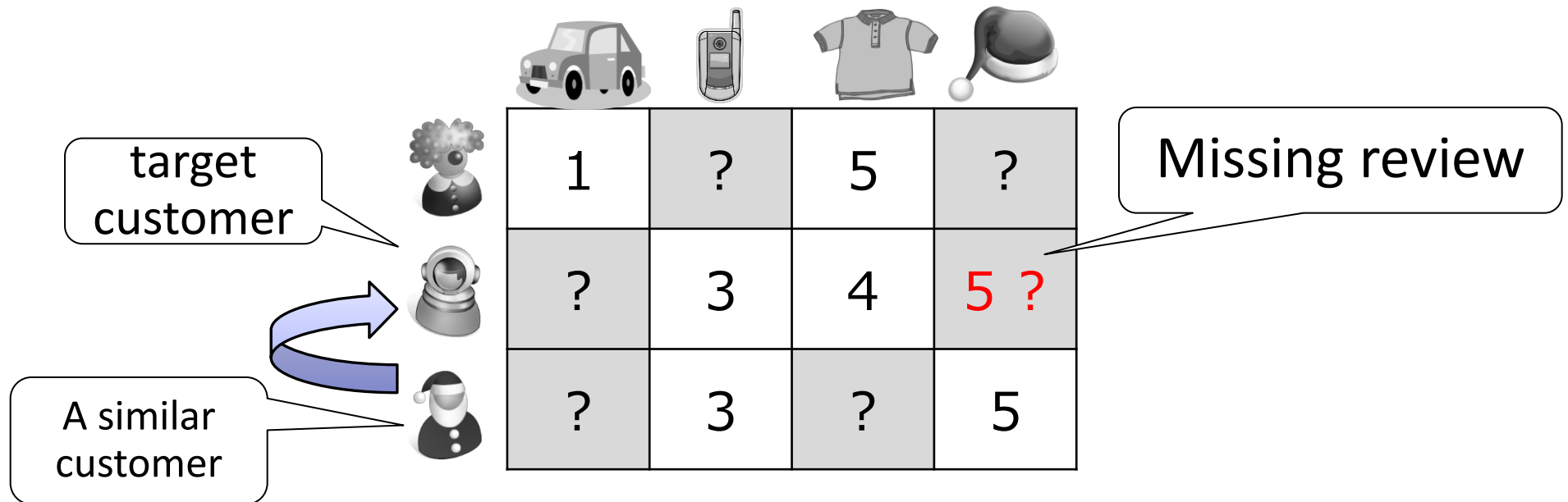
- A matrix with rows (customers) and columns (products)
  - Each element = review score
- Given observed parts of the matrix, predict the unknown parts (  )



# Basic idea of recommendation algorithms:

## “Find people like you”

- GroupLens: an earliest algorithm (for news recommendation)
  - Inherited by MovieLens (for Movie recommendation)
- Find people similar to the target customer, and predict missing reviews with theirs



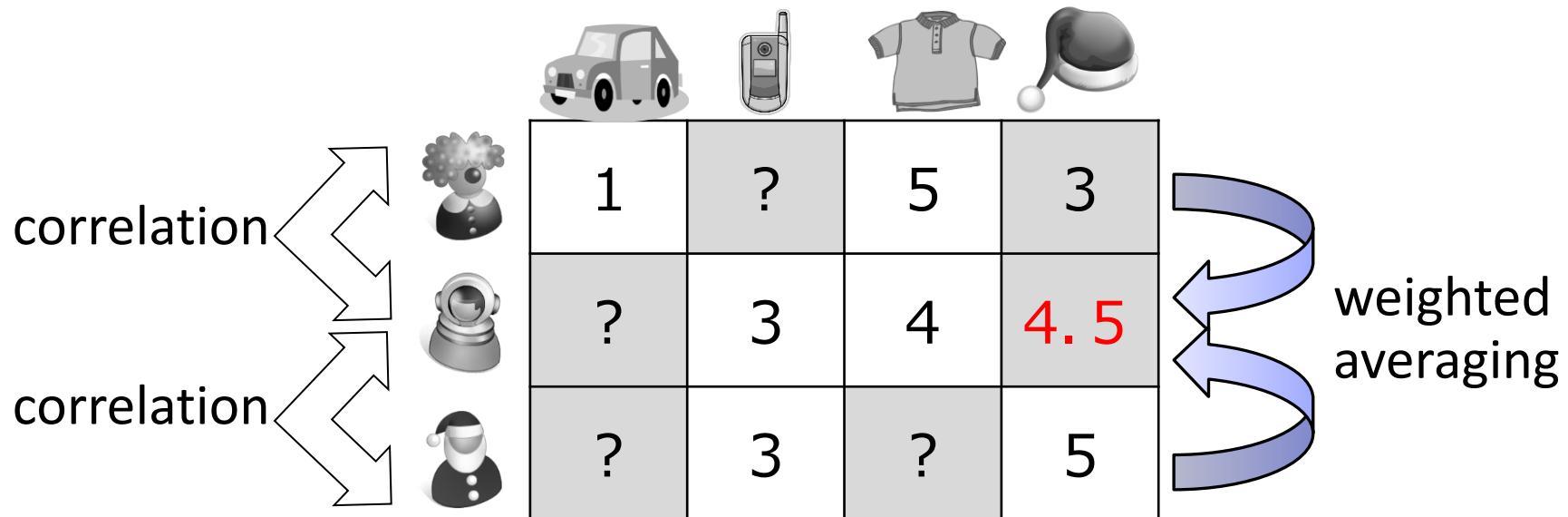
# GroupLens:

## Weighted prediction with correlations among customers

- Define customer similarity by correlation ( of observed parts )
- Make prediction by weighted averaging with **correlations** :

$$y_{i,j} = y_i + \sum_{k \neq i} \rho_{i,k} ( y_{k,j} - y_k ) / \sum_{k \neq i} \rho_{i,k}$$

correlation



# Low-rank assumption for matrix completion:

## GroupLens implicitly assumes low-rank matrices

---

- Assumption of GroupLens algorithm:  
Each row is represented by a linear combination of the other rows (i.e. linearly dependent)  
 $\Rightarrow$  The matrix is not full-rank ( $\doteq$  low-rank)
- Low-rank assumption helps matrix completion

# Low-rank matrix factorization:

## Projection onto low-dimensional latent space

- Low-rank matrix: product of two (thin) matrices

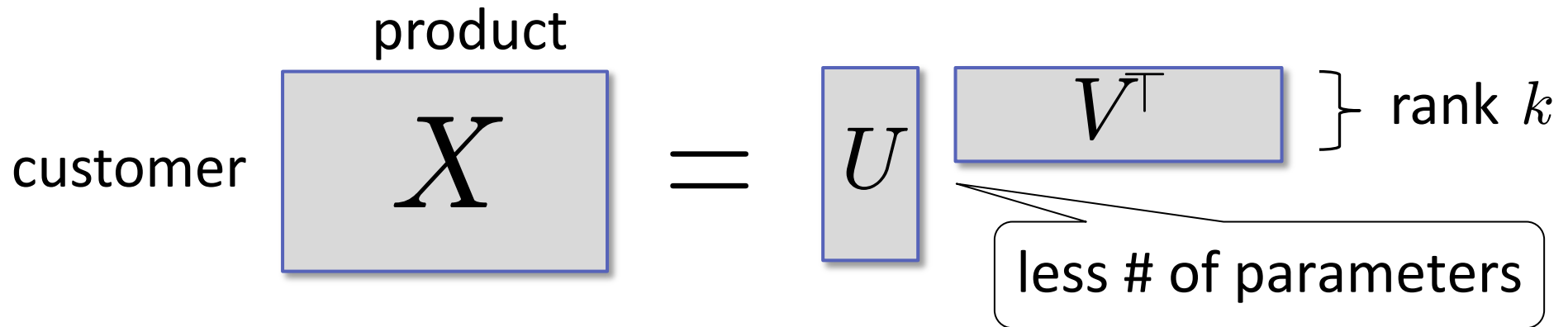
customer

product

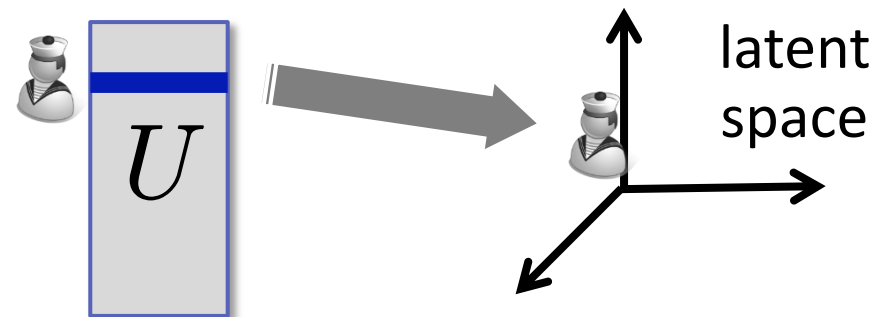
$$X = U V^T$$

} rank  $k$

less # of parameters



- Each row of  $U$  and  $V$  is an embedding of each customer (or product) onto low-dimensional latent space



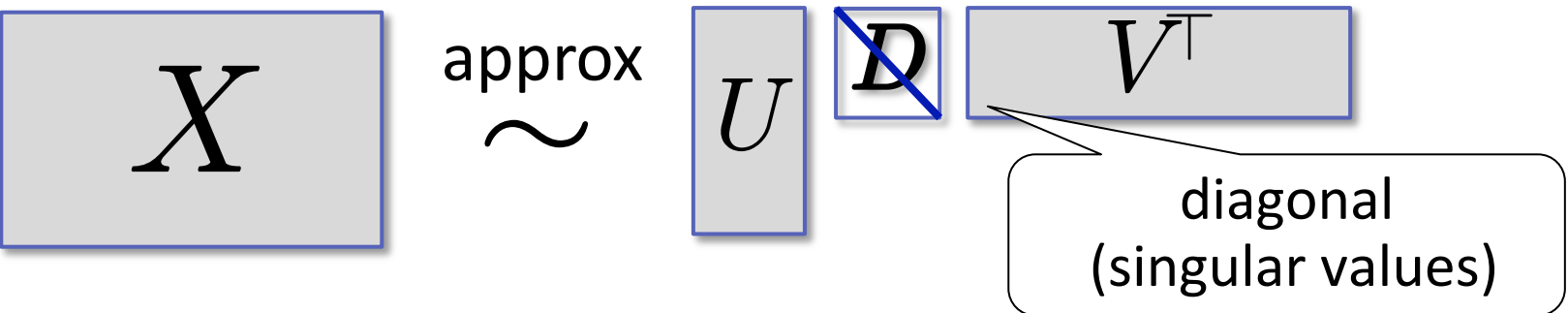


# Example of low-rank matrix decomposition: Singular value decomposition

- Find a best low-rank approximation of a given matrix

$$\underset{Y}{\text{minimize}} \quad \|X - Y\|_F^2 \quad \text{s.t.} \quad \text{rank}(Y) \leq k$$

- Singular value decomposition (SVD)

—   $X \approx U \cancel{D} V^\top$

diagonal  
(singular values)

wrt constraint:  $U^\top U = I \quad V^\top V = I$

- The largest  $k$  eigenvalues of  $X^\top X$  best approximate

# Strategies for matrices with missing values:

## EM algorithm, gradient descent, and trace norm

---

- SVD is not applicable to matrices with missing values
- For completion problem:
  - Direct application of SVD to a (somehow) filled matrix
  - Iterative applications: iterations of completion and decomposition
- For large scale data:  
Gradient descent using only observed parts
- Convex formulation: Trace norm constraint

# Predicting more complex relations:

## Multinomial relations

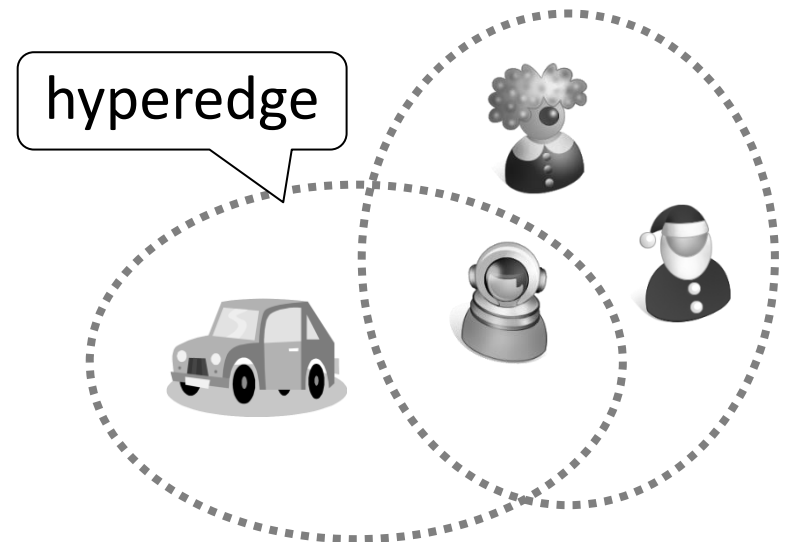
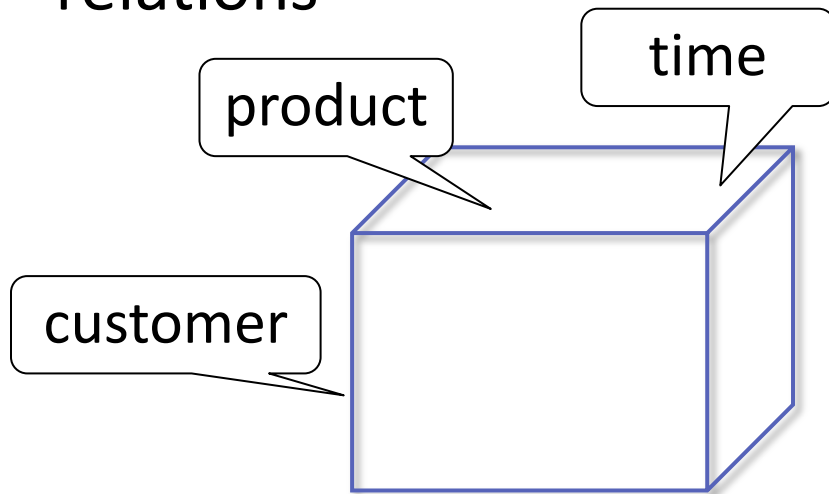
---

- Matrices can represent only one kind of relations
  - Various kinds of relations (actions):  
Review scores, purchases, browsing product information, ...
  - Correlations among actions might help
- Multinomial relations:
  - (customer, product, action)-relation:  
(Alice, iPad, buy) represents “Alice bought an iPad.”
  - (customer, product, time)-relation:  
(John, iPad, July 12<sup>th</sup>) represents “John bought an iPad on July 12th.”

# Multi-dimensional array:

## Representation of multinomial relations

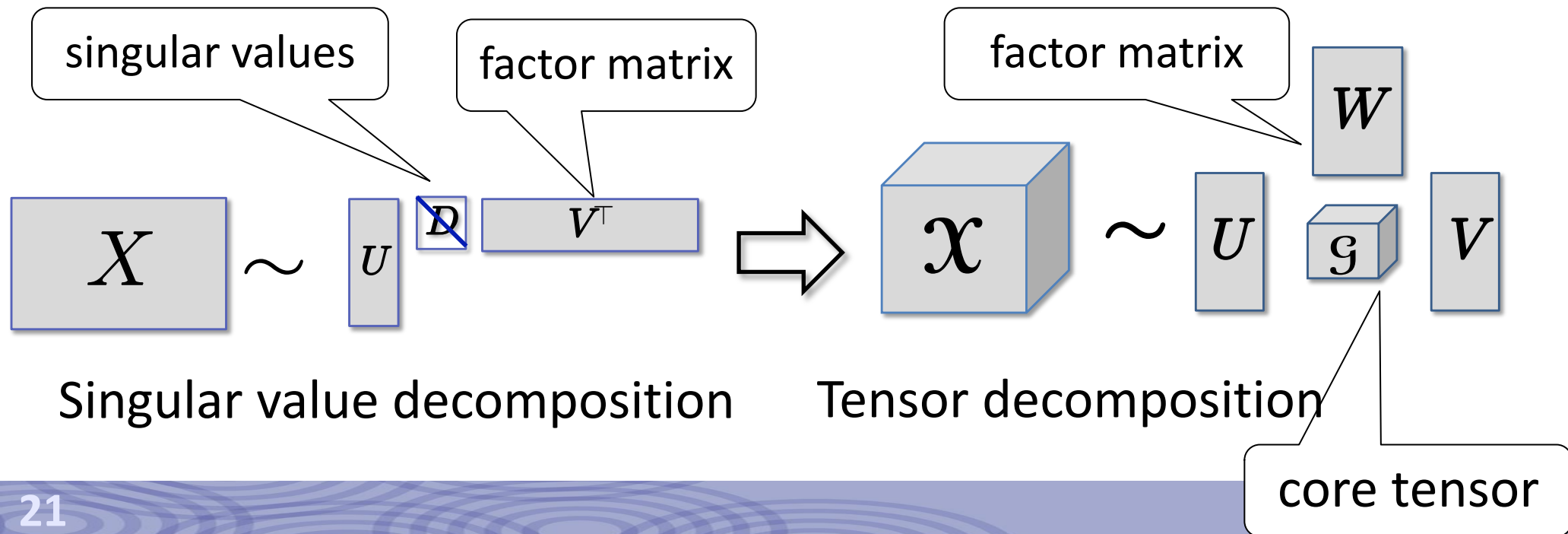
- Multidimensional array: Representation of complex relations among multiple object
  - Types of relations (actions, time, conditions, ...)
  - Relations among more than two objects
- Hypergraph: allows variable number of objects involved in relations



# Tensor decomposition:

## Generalization of low-rank matrix decomposition

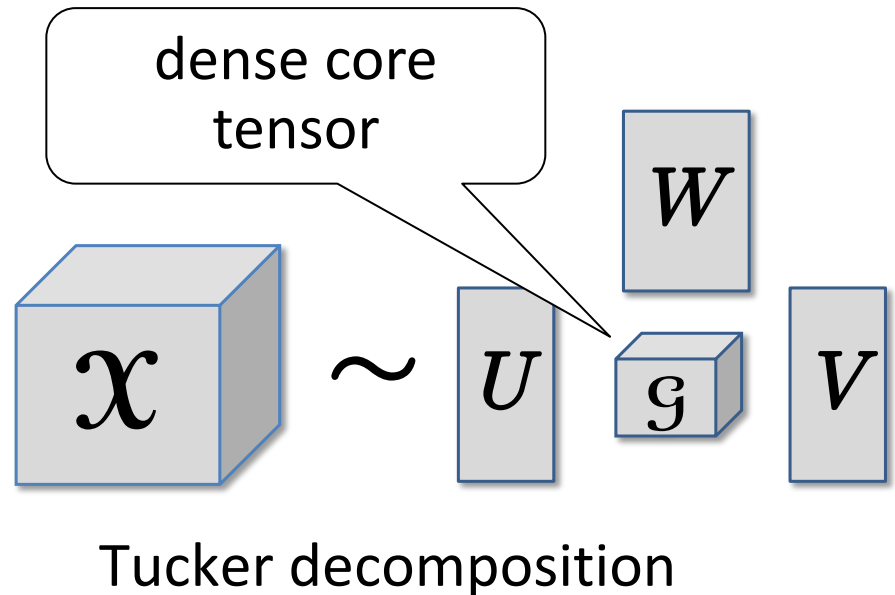
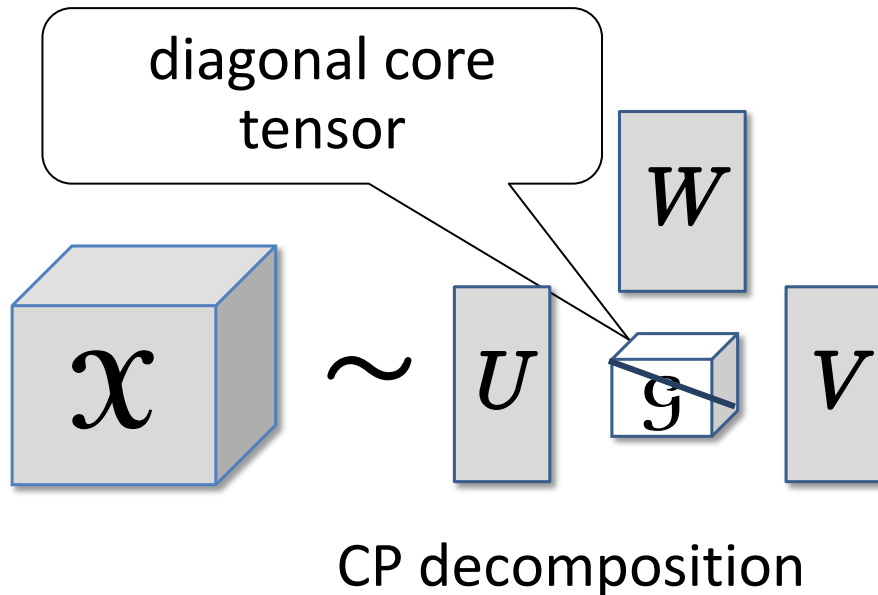
- Generalization of matrix decomposition to multidimensional arrays
  - A small core tensor and multiple factor matrices
- Increasingly popular in machine learning/data mining



# Tensor decompositions:

## CP decomposition and Tucker decomposition

- CP decomposition: A natural extension of SVD (with a diagonal core)
- Tucker decomposition: A more compact model (with a dense core)



# Applications of tensor decomposition:

## Tag recommendation, social network analysis, ...

---

- Personalized tag recommendation ( $\text{user} \times \text{webpage} \times \text{tag}$ )
  - predicts tags a user gives a webpage
- Social network analysis ( $\text{user} \times \text{user} \times \text{time}$ )
  - analyzes time-variant relationships
- Web link analysis  
( $\text{webpage} \times \text{webpage} \times \text{anchor text}$ )
- Image analysis ( $\text{image} \times \text{person} \times \text{angle} \times \text{light} \times \dots$ )

# Applications of tensor decomposition:

## Tag recommendation, social network analysis, ...

- Personalized tag recommendation (user  $\times$  webpage  $\times$  tag)

- predicts tags a user gives a webpage

- Social

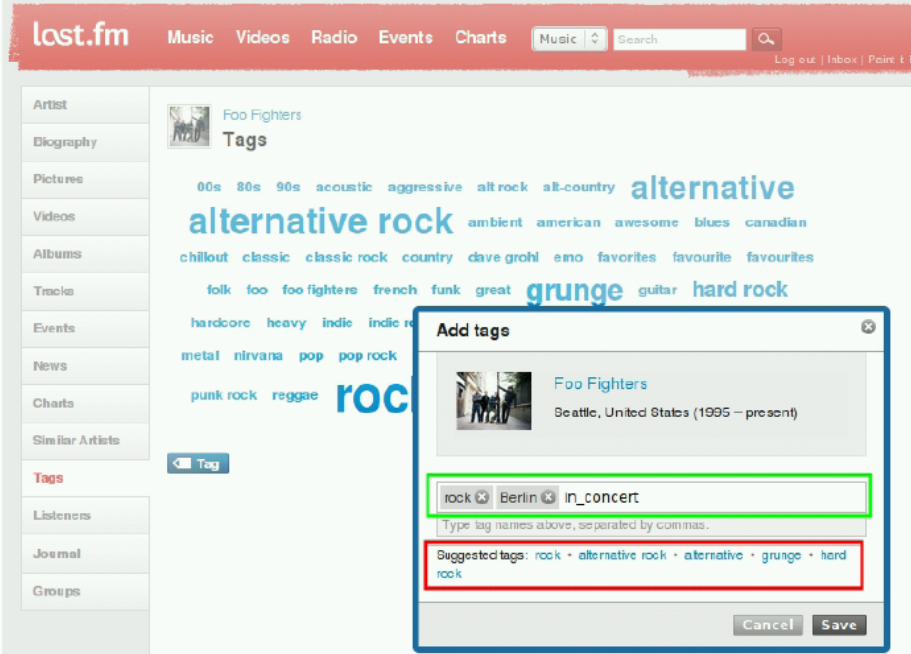
- anal

- Web I  
(web

- Image

Tag Recommendation   Optimization Criterion   Factor Models   Evaluation   Conclusion

### Personalized Tag Recommendation



The screenshot shows the lost.fm website with the 'Tags' section for the artist Foo Fighters. A modal window titled 'Add tags' is open, showing the artist's name and a list of suggested tags: rock, alternative rock, alternative, grunge, and hard rock. The user has entered 'rock', 'Berlin', and 'in\_concert' in the input field. The modal also includes a 'Cancel' button and a 'Save' button.

**Task:** Recommend a user a (personalized) list of tags for a specific item.

Steffen Rendle, Lars Schmidt-Thieme   2 / 14   ISMILL, University of Hildesheim

(...)

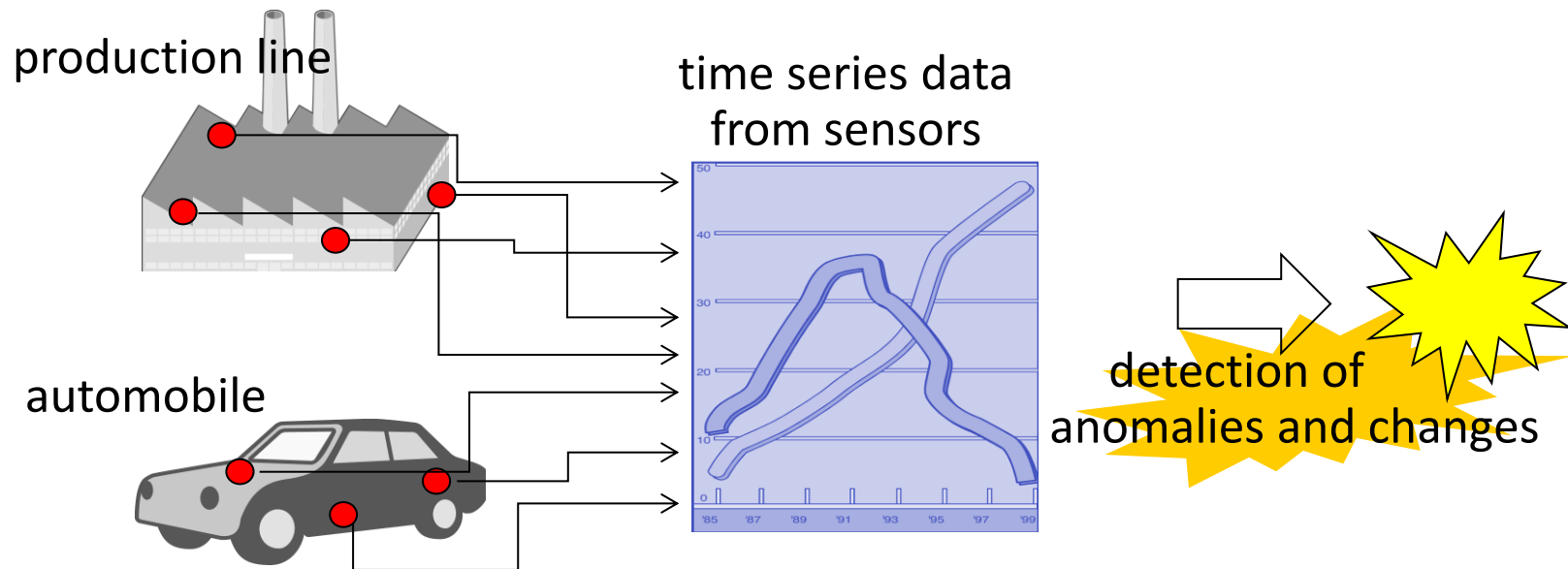


# Anomaly Detection

# Anomaly detection:

## Early warning for system failures reduces costs

- A failure of a large system can cause a huge loss
  - production line in factory
  - Infection of computer virus/intrusion to computer systems
- Early detection of failures from data collected from sensors



# Difficulty in anomaly detection:

## Failures are often unknown

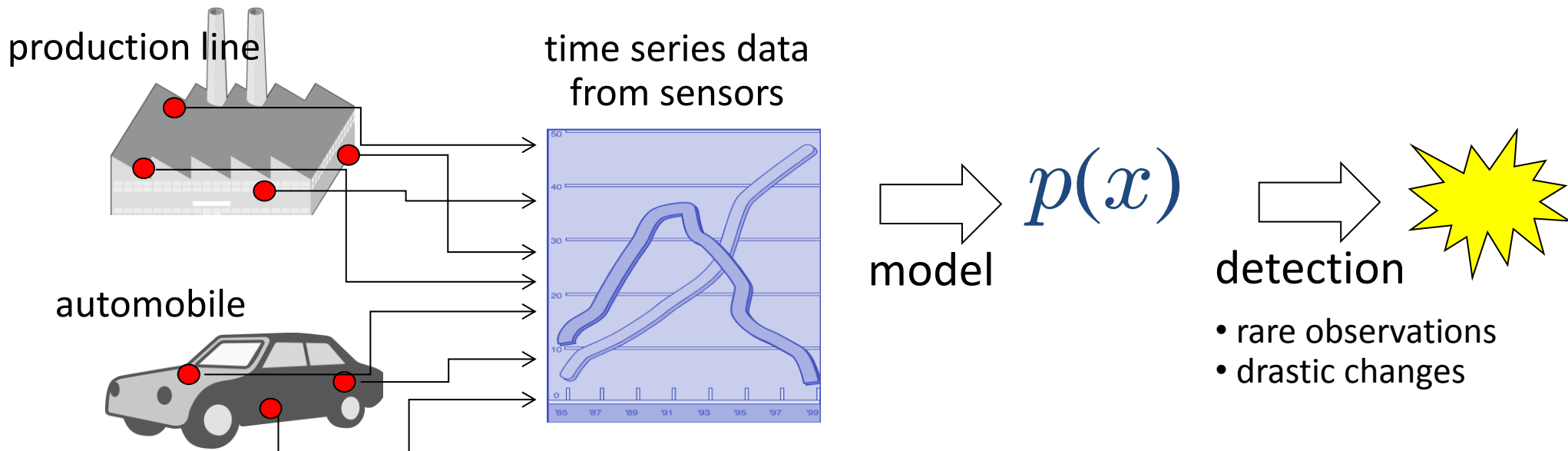
---

- Known failures are detected by using supervised learning:
  1. Construct a predictive model from past failure data
  2. Apply the model to system monitoring
- Serious failures are often new ones
  - No past data are available
- There are many cases where supervised learning is not applicable

# Change the strategy:

## Model the normal times, detect deviation from them

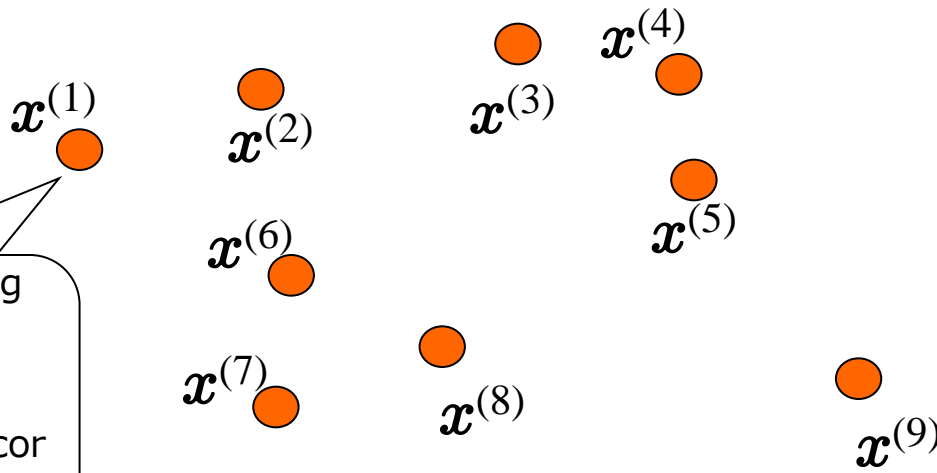
- Difficult to model anomalies → Model normal times
  - Data at normal times are abundant
- Observation of rare data is a precursor of failures



# Clustering:

## Model the normal times by grouping the data

- Divide normal time data  $\{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$  into  $K$  groups
  - Group is represented by centers  $\{\mu^{(1)}, \mu^{(2)}, \dots, \mu^{(K)}\}$
- data  $x$  is an “outlier” if it lies far from all of the centers  
= system failures, illegal operations, instrument faults

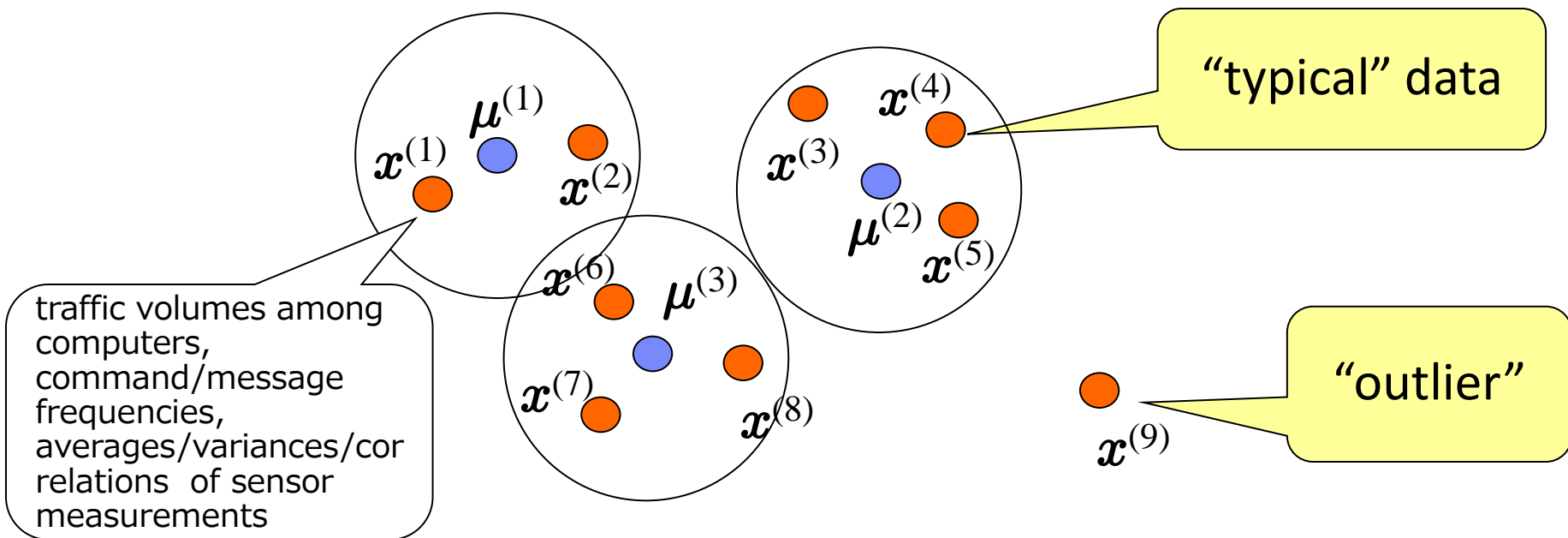


traffic volumes among  
computers,  
command/message  
frequencies,  
averages/variances/cor  
relations of sensor  
measurements

# Clustering:

## Model the normal times by grouping the data

- Divide normal time data  $\{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$  into  $K$  groups
  - Group is represented by centers  $\{\mu^{(1)}, \mu^{(2)}, \dots, \mu^{(K)}\}$
- data  $x$  is an “outlier” if it lies far from all of the centers  
= system failures, illegal operations, instrument faults



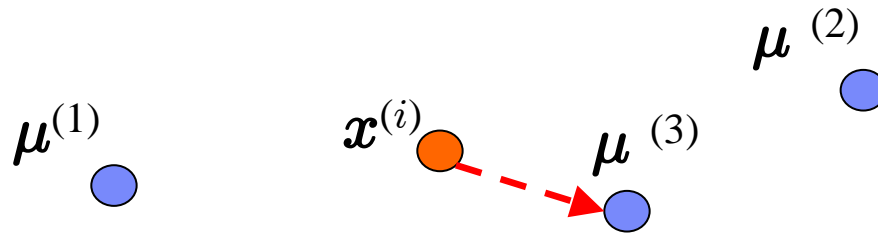
# $K$ -means algorithm:

## Iterative refinements of groups

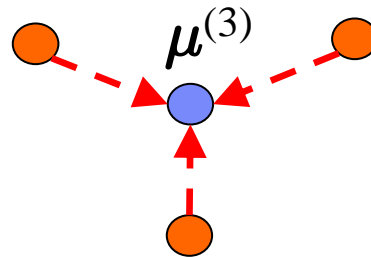
---

- Repeat until convergence:

1. Assign each data  $x^{(i)}$  to its nearest center  $\mu^{(k)}$



2. Update each center to the center of the assigned data

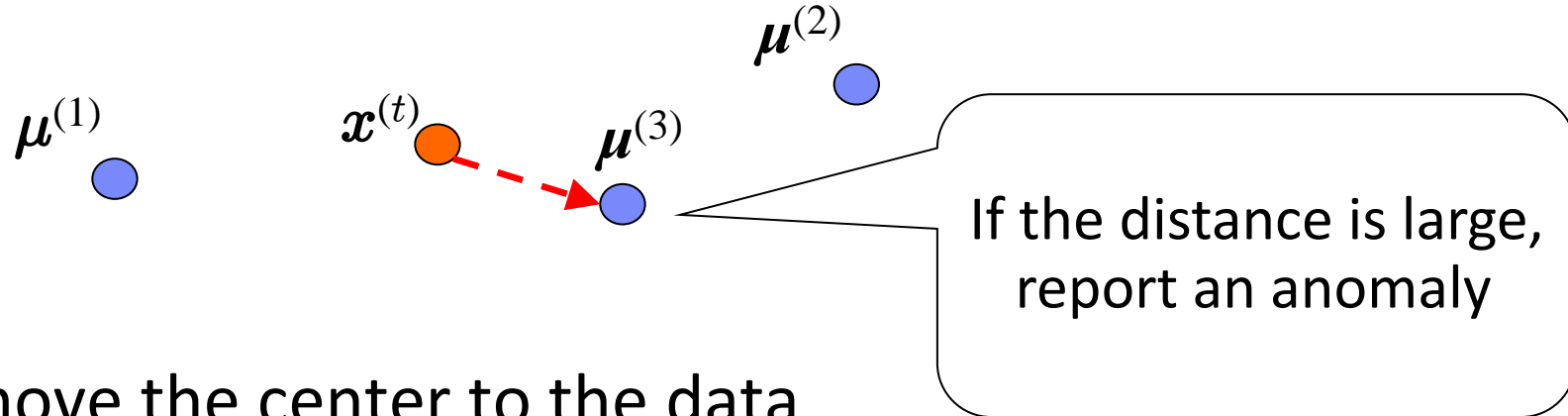


# Sequential $K$ -means:

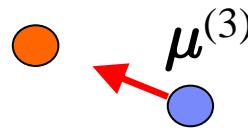
## Simultaneous estimation of clusters and outliers

- Data arrives in a streaming manner, and apply clustering and anomaly detection at the same time

1. Assign each data  $x^{(t)}$  to its nearest center  $\mu^{(k)}$



2. Slightly move the center to the data





# Limitation of anomaly detection:

## Failures are unknown

---

- Anomalies are not defined in advance  
→ Details of failures are unknown
- Report to system administrators  
→ Investigations by administrators
- Change detection: detects model changes

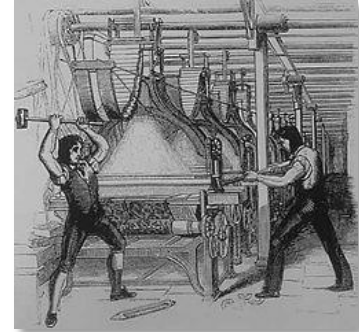
# Future Direction

# Can machines with big data go beyond humans?:

## Many instances. But, not always

---

- Computers overwhelming humans on many “intelligent” tasks:
  - e.g. the victory of IBM’s QA system WATSON trained with big data
- Concerns about replacement of humans by computers:
  - Luddite movement in the industrial revolution in 19<sup>th</sup> century
- With big data, do computer always defeat humans?
  - No, there are still many only human can do



# ReCAPTCHA :

## Authentication with a “hidden” task

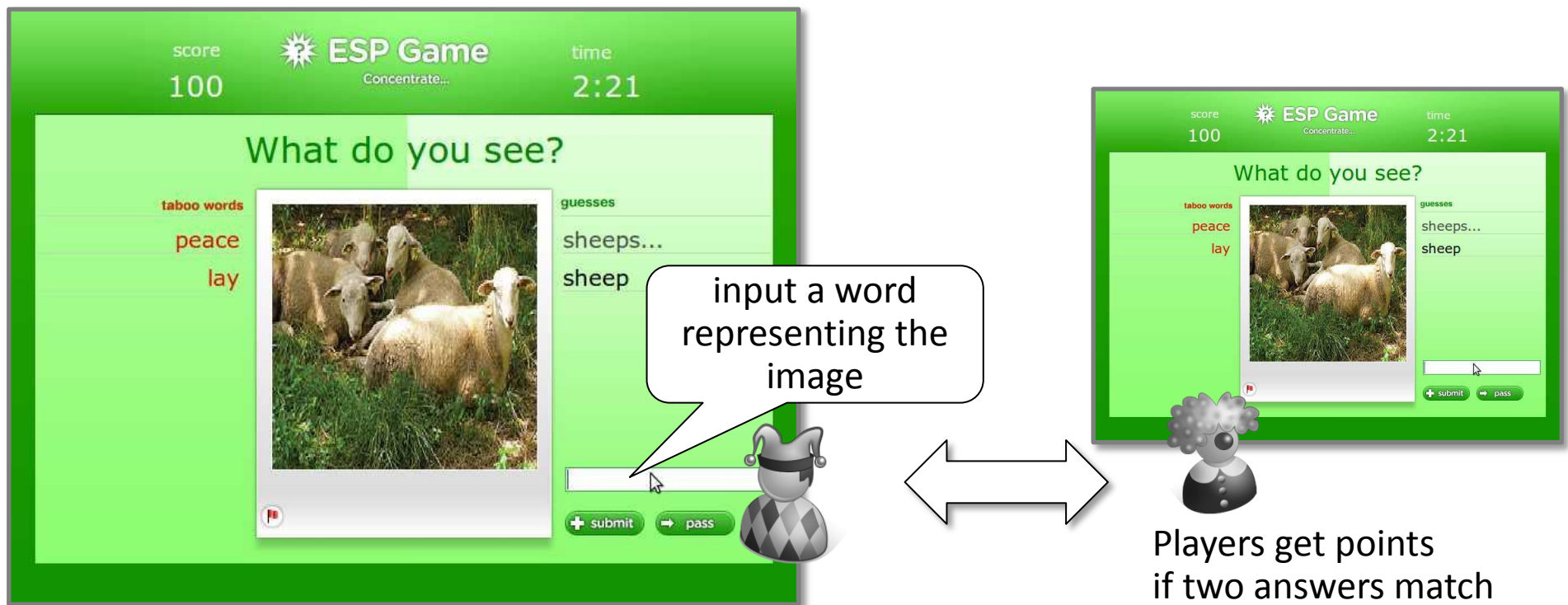
- A Turing test discriminating humans from computers
  - Authentication of an access by human to a Web site requires users to read and input given two words
  - Hard for computers, but easy for humans



# ESP game:

## A game with a “hidden” task

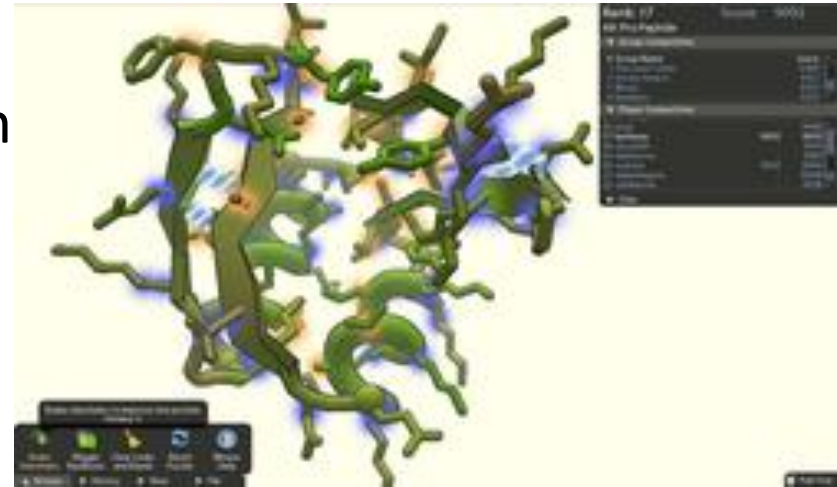
- A cooperative game played by two online players
- Players get points if two keywords given by them to a presented image coincide



# “Hidden” tasks:

## Tasks hard for computers, but easy for humans

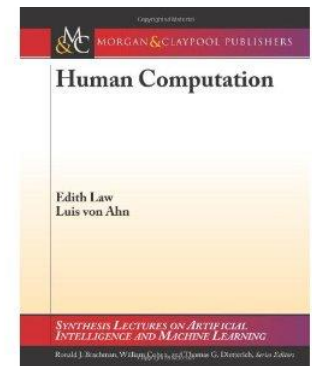
- ReCAPCHA: digitization of paper documents
  - The system does not know one of two presented words
  - Human helps digitization of paper documents
- ESP game: image tagging
  - Tags for images improves image search
  - Human helps tagging of images
  - “Game with a purpose (GWAP)”
    - Music tagging, protein folding, ...
- Both examples cleverly embed computer-unfriendly tasks into other forms



# Human computation:

## Cooperative Problem solving by machines and humans

- Human computation:
  - regards humans as computational resources
  - solves problems that computers can not solve
- Efficient cooperative problem solving
  - E.g. quicksort with humans
- Employment of human resources:
  - Gamification: Tasks implicitly embedded into games
  - Crowdsourcing (e.g. Mechanical Turk)
  - . . .

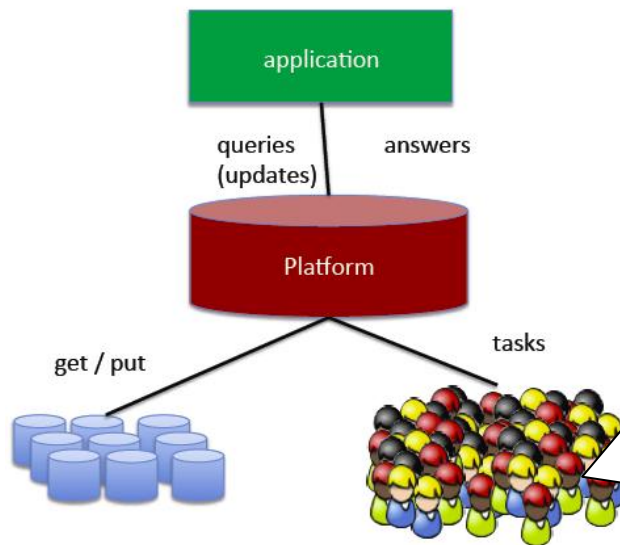


Law & Von Ahn (2012), Human Computation

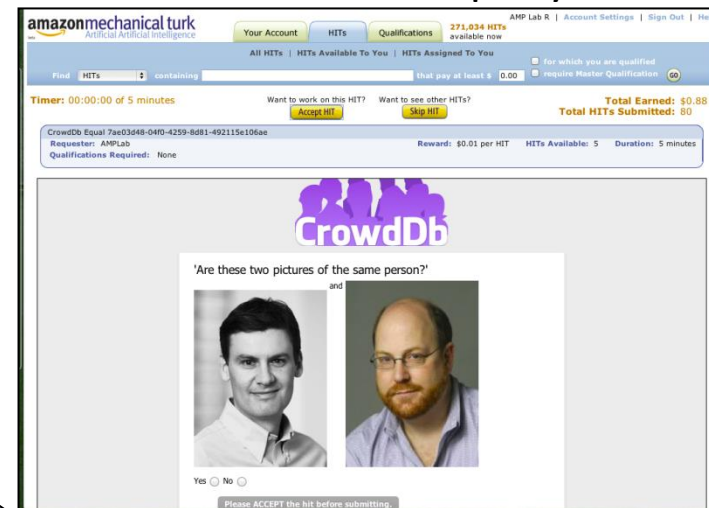


# Human computation in computer science: NLP, computer vision, HCI, web search, DB, ...

- Natural language processing: NL understanding, translation, ...
- Computer vision: image understanding, annotation, detection,...
- DB/IR: data generation/integration, evaluation, ...
- Machine learning: data collection



Call crowdsourcing services to  
evaluate a SQL query



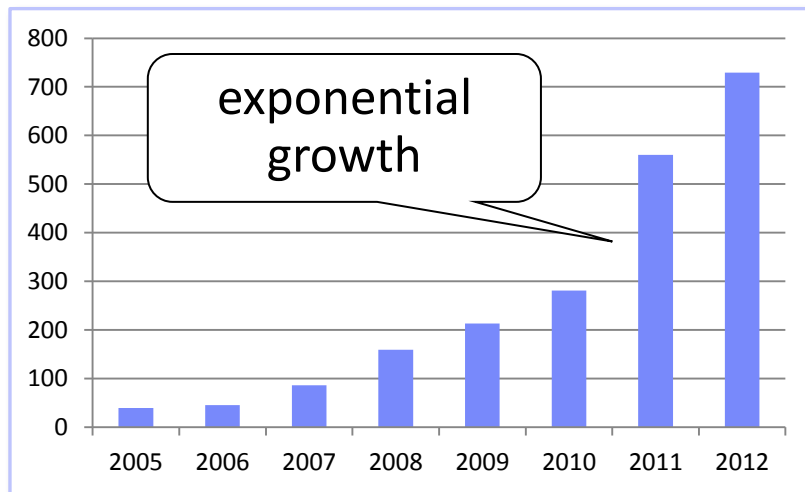


# Research trend:

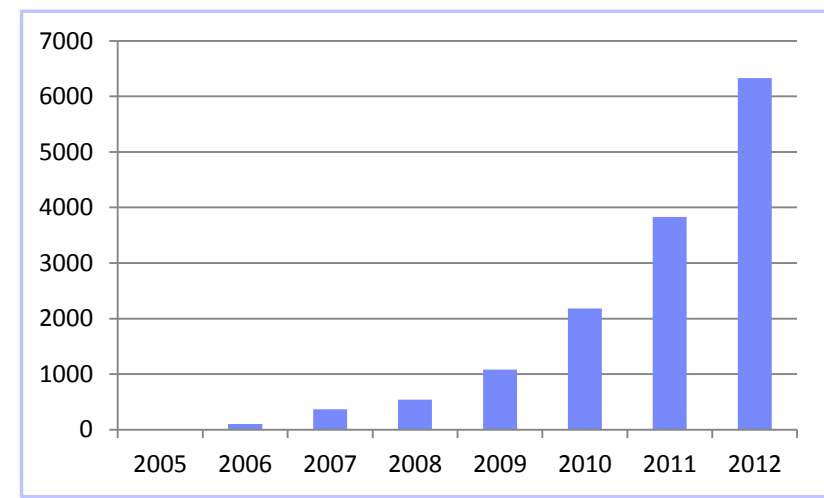
## Exponential growth of human computation/crowdsourcing

### ■ Exponential growth of related research

#papers related to “human computation”



#papers related to “crowdsourcing”



✂ According to Google Scholar

- 2005: Amazon Mechanical Turk
- 2006: “human computation” “crowdsourcing”
- 2013: The first international conference on human computation and crowdsourcing HCOMP (Human Computation & Crowdsourcing)

# Future of artificial intelligence research: Is human computation a compromise? or a new direction?

- Ask humans to do tasks computers can not do
  - Does it mean giving up “artificial intelligence”?
- This might be a new research direction
- “Freestyle” chess tournament:
  - A team of 2 amateurs and 3 programs won
  - Computers + humans overwhelms best of each
- Not only being a compromise,  
this might be a new artificial intelligence



# Summary:

## Machine learning and its applications (and beyond)

---

- Increasing popularity of machine learning:
  - supported by the rise of “big data” and “data scientists”
- Many applications:
  - Supervised machine learning: recommender systems
    - Matrix factorization/tensor decomposition
  - Unsupervised machine learning: Anomaly detection
    - Capture normal times (by clustering)
- Human computation: Collaborative problem solving by machines and humans