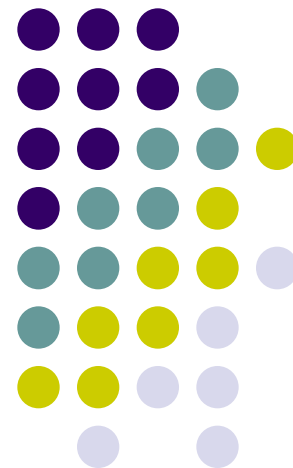




# BÀI GIẢNG KHAI PHÁ DỮ LIỆU WEB

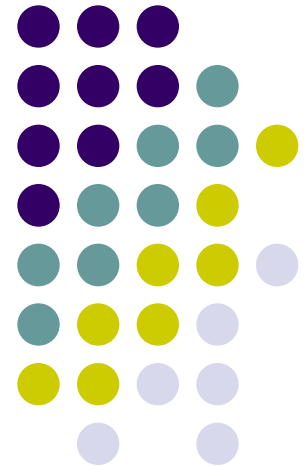
## CHƯƠNG 1. GIỚI THIỆU CHUNG

PGS. TS. HÀ QUANG THỤY  
HÀ NỘI 10-2010  
**TRƯỜNG ĐẠI HỌC CÔNG NGHỆ**  
ĐẠI HỌC QUỐC GIA HÀ NỘI



# Nội dung

1. Giới thiệu về khai phá text
2. Giới thiệu về khai phá web



# 1. Giới thiệu về khai phá text



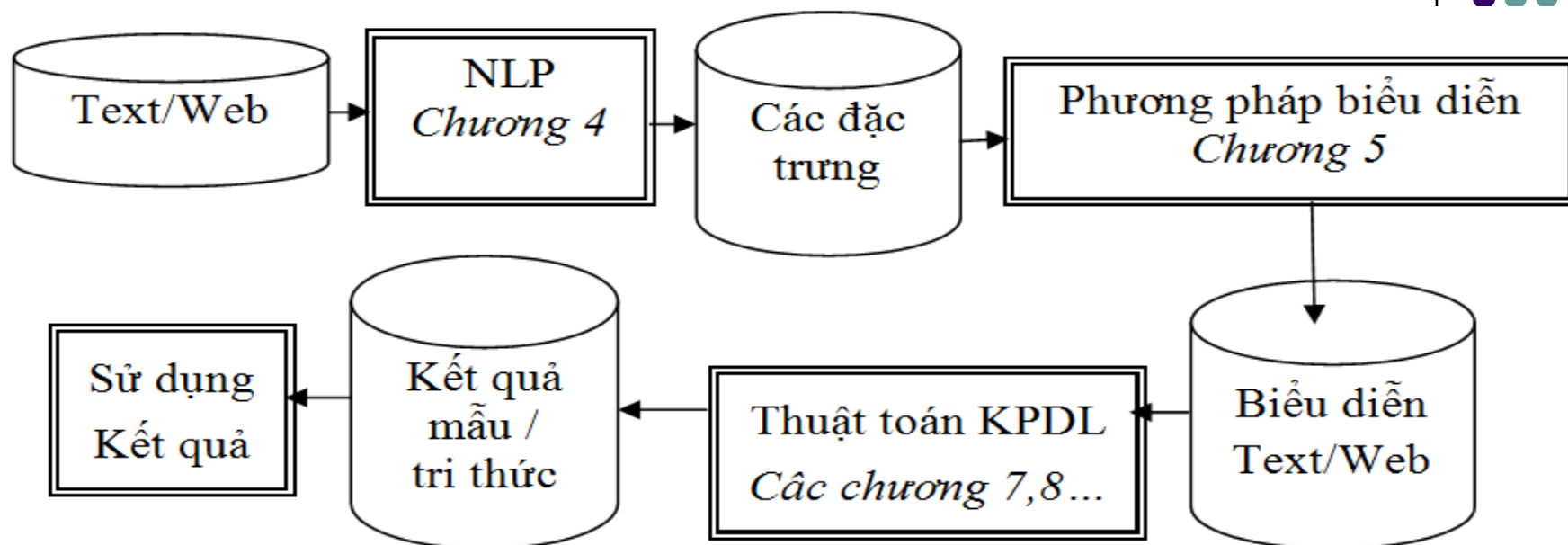
- Khái niệm
- Sự cần thiết của khai phá text
- Đặc trưng của khai phá text
- Các bài toán cơ bản trong khai phá text
- Một ví dụ về bài toán khai phá text
- Xu hướng nghiên cứu khai phá Text

# Khái niệm



- **Tiếp cận về khái niệm khai phá text**
  - Khai phá text là khai phá dữ liệu đối với loại dữ liệu text.
  - Quá trình phát hiện tri thức mới, có giá trị, tiềm ẩn trong tập hợp văn bản
  - Mang tính đa dạng về phát biểu khái niệm khai phá dữ liệu
- **Nội dung**
  - Khai phá text = Khai phá dữ liệu + Xử lý ngôn ngữ tự nhiên - XLNNTN (Natural Language Processing: NLP)
  - Các bài toán chung về khai phá dữ liệu cho dữ liệu đặc thù
  - Một số bài toán riêng điển hình cho khai phá text
- **Mối quan hệ giữa Khai phá Text và XLNNTN**
  - XLNNTN cung cấp tài nguyên, công cụ cơ sở cho khai phá Text
  - Khai phá Text mở rộng các bài toán của XLNNTN
  - Đan xen giữa Khai phá Text với XLNNTN

# Quy trình khai phá text



- **Tuân theo quy trình chung của khai phá dữ liệu**
  - Như đã trình bày trong khai phá dữ liệu
- **Quy trình tối giản**
  - Tiền xử lý
    - Công cụ của Xử lý ngôn ngữ tự nhiên
    - Mô hình cấu trúc văn bản
  - Biểu diễn văn bản
    - Phù hợp với thuật toán
  - Xử lý (khai phá) dữ liệu theo dạng biểu diễn
    - Áp dụng khai phá dữ liệu



# Sự cần thiết của khai phá text

- Text gần gũi nhất với con người
  - Là đối tượng quan trọng nhất chuyển tải thông tin của loài người
  - Phương tiện trình bày tri thức ████████ xuyên giao người khác
  - Học chữ là bài toán quan trọng của mỗi con người
- Đặc thù của ngôn ngữ tự nhiên
  - Tính đa nghĩa, đồng nghĩa của đơn vị cú pháp nhỏ nhất là từ
  - Tính cảm ngữ cảnh khi trình bày nội dung văn bản
  - Tính biến động của mỗi ngôn ngữ tự nhiên: bổ sung, thay đổi...
- Sự tăng trưởng của dữ liệu Text
  - Khả năng tạo mới
  - Khả năng lưu trữ

# Đặc trưng của khai phá text



<i>Dấu hiệu phân biệt</i>	<i>Khai phá dữ liệu</i>	<i>Khai phá Text</i>
<i>Đối tượng dữ liệu</i>	Dữ liệu số / phân loại	Văn bản
<i>Cấu trúc đối tượng</i>	CSDL quan hệ	Text dạng tự do: không cấu trúc, nửa cấu trúc
<i>Mục tiêu</i>	Dự báo, đoán nhận	Tìm kiếm thông tin liên quan, hiểu ngữ nghĩa, phân lớp / phân bố
<i>Phương pháp</i>	Học máy: DT, MBR, ...	Chỉ số, xử lý mạng nơron, ngôn ngữ, kiến trúc
<i>Kích cỡ thị trường</i>	Trăm nghìn phân tích viên từ công ty lớn và vừa	Hàng triệu người dùng từ hãng và cá nhân
<i>Tình trạng</i>	Quảng bá từ năm 1994	Mới quảng bá từ năm 2000

Sergei Ananyan (2001). Text Mining: Applications and Technologies, *Megaputer Intelligence Inc.* (truy nhập ngày 13/9/2003)



# Một số bài toán điển hình trong TM



- **Biểu diễn Text**
  - Là một trong những bài toán quan trọng nhất trong khai phá Text
  - Nghịch lý về “hiệu quả như nhau” trong tìm kiếm Text
  - Tìm biểu diễn phù hợp nhất cho bài toán khai phá text
  - Một lớp hướng mô hình biểu diễn Text: Mô hình sinh Text
  - Nội dung của chương 2.
- **Tìm kiếm/thu hồi Text (Text Search/Retrieval)**
  - Cho một tập văn bản và một yêu cầu tìm kiếm của người dùng (dạng văn bản / khác).
  - Mục đích: Tìm tập văn bản trong CSDL đáp ứng yêu cầu người dùng
  - Đã tồn tại một CSDL Text: Tìm kiếm full-text trong CSDL này
  - Tìm kiếm trên Internet. Máy tìm kiếm: Nội dung chương 5.

# Một số bài toán điển hình trong TM (2)



## ● Phân lớp văn bản

- Tương ứng học có giám sát (học có thầy)
- Cho trước tập lớp và tập ví dụ
- Mục tiêu : một mô hình phân lớp thực hiện ánh xạ mỗi văn bản vào lớp
- Ví dụ:

## ● Phân cụm văn bản

- Tương ứng học không giám sát
- Cho trước tập văn bản
- Mục tiêu : tập cụm văn bản và tóm tắt cụm.
- Ví dụ:

## ● Phân đoạn văn bản

- Phân cụm và phân lớp
- Ví dụ:

# Một số bài toán điển hình trong TM (3)

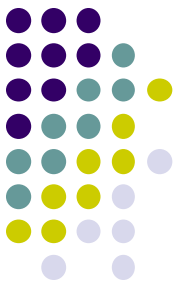


## ● Phân tích ngữ nghĩa

- Hiểu văn bản (xem DUC: Document Understanding Conferences và TAC: Text Analysis Conferences)
- Ngữ nghĩa của các thành phần trong văn bản
- Phát hiện quan hệ thực thể trong văn bản
- Taxonomy, ontology, web ngữ nghĩa (semantic Web)
- Roxana Girju [Gij08] liệt kê một số danh sách quan hệ ngữ nghĩa, trong đó có danh sách 22 quan hệ do chính tác giả tổng hợp:
  - HYPERNYMY (IS-A)      PART-WHOLE (MERONYMY)      CAUSE      POSSESSION
  - KINSHIP      MAKE/PRODUCE      INSTRUMENT      TEMPORAL
  - LOCATION/SPACE      PURPOSE      SOURCE/FROM      EXPERIENCER
  - TOPIC      MANNER      MEANS      GENT
  - THEME      PROPERTY      BENEFICIARY      MEASURE
  - TYPE      DEPICTIONDEPICTED.

[Gir08] Roxana Girju (2008). Semantic Relation Extraction and its Applications, [ESSLLI 2008](#): Invited Tutorial, Hamburg, Germany, August 2008

# Một số bài toán điển hình trong TM (4)



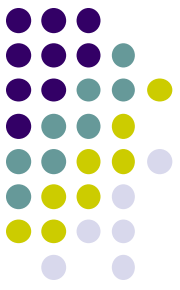
## ● Trích chọn đặc trưng

- Phát hiện/lưu trữ từ khóa (term), đặc trưng (feature), cụm từ mang nghĩa
- Đặc trưng chưa định trước: xác định đồng thời với phân tích nội dung
- Phân biệt trích chọn đặc trưng (feature extraction) với chọn lựa đặc trưng (feature selection)
- Phân tích văn bản để phát hiện tần số xuất hiện

## ● Tóm tắt văn bản

- Document Abstract/Summarization
- Xây dựng một văn bản thu gọn hơn (tỷ lệ/số lượng từ/câu) song vẫn giữ được ngữ nghĩa
- Abstract (rút trích câu) /Summarization (xây dựng câu)
- Xây dựng tự động mục lục văn bản
- Tóm tắt đơn văn bản/ tóm tắt đa văn bản
- Quan hệ chặt chẽ với “hiểu văn bản”

# Một số bài toán điển hình trong TM (5)



- **Xây dựng ontology**
  - Kho ngữ liệu về một/một nhóm lĩnh vực
  - Phục vụ, nâng cao chất lượng các bài toán ngữ nghĩa
  - Tập khái niệm, lớp khái niệm, quan hệ giữa chúng
  - Biểu diễn hình học dạng đồ thị
  - Dạng đặc biệt: Taxonomy
  - Ví dụ: WordNet, TreeBank
- **Kế thừa nguyên bản (Textual Entailment)**
  - “Văn bản T kế thừa giả thiết nguyên bản H” nếu tính chân thực của H có thể được suy diễn từ T.
  - “Ý nghĩa” của T tiềm ẩn trong H: trình bày nào đó của H có thể phù hợp trình bày nào đó của T (mức độ chi tiết hay trừu tượng)
- **Dẫn đường văn bản (Text focusing)**
  - Tích hợp xử lý văn bản với cơ sở tri thức cho phép kết nối trực tiếp tri thức trong quá trình xử lý văn bản
  - Dẫn dắt các văn bản theo tri thức đã được kết nối

# Một số bài toán điển hình trong TM (6)



- Khai phá quan điểm
  - Là chủ đề thời sự hiện nay
  - Đối tượng: không là sự vật/ hiện tượng mà là tình cảm thái độ
  - Ứng dụng: tiếp thị (quan hệ khách hàng), điều tra xã hội học...
  - Một số ví dụ
- Khai phá Text trong lĩnh vực cụ thể
  - Y Sinh học: Quan hệ tương tác protein – protein, gene – bệnh
  - Các lĩnh vực khoa học khác:



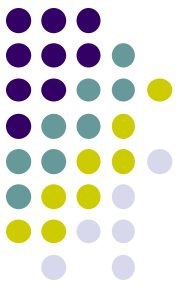
# Một số bài toán ví dụ

## ● Ví dụ 1

- **Nêu bài toán:** Nhằm mục đích quản lý, một công ty Nhật Bản muốn xây dựng một hệ thống “quản lý” các nội dung đã được máy in của công ty in ra.
- **Đặt vấn đề:**
  - Xây dựng hệ thống quản lý văn bản với thuộc tính in văn bản. Do một số lý do, đây không phải là điều công ty muốn.
  - Quản lý mọi nội dung được in ra: Dữ liệu nguồn chỉ có thể là dòng dữ liệu đi qua máy in của công ty. Cần xây dựng hệ thống có các năng lực (1) lấy được dòng dữ liệu Text đi tới các máy in; (2) Tổ chức lại hệ thống các văn bản được in ra để thuận tiện cho việc quản lý.
- **Giải pháp:**
  - Thu nhận dữ liệu: Xây dựng luồng xử lý dòng dữ liệu vào máy in, một bản đưa ra máy in và một bản đưa vào thành phần xử lý tiếp theo.
  - Tổ chức hệ thống văn bản: Tiền xử lý dữ liệu; phân lớp đã cấp (trong đó có phân cụm)

Nguồn: từ một học viên công tác tại FSOFTE làm việc với Nhật Bản

# Một số bài toán ví dụ (2)



- Ví dụ 2. Bài toán của Rich Caruana & cộng sự
  - Bài toán: Cho trước một tập (khoảng 300000) công trình nghiên cứu khoa học (bài đăng tạp chí, báo cáo hội nghị, luận án Tiến sỹ) đã được công bố. Từ nội dung văn bản của mỗi công trình nghiên cứu, chúng ta nhận được tên tác giả (các tác giả), các tài liệu tham khảo, nơi công bố (tên tạp chí, hội nghị, hội thảo ...).
  - Yêu cầu: Chỉ dùng nội dung, năm XB và tên các tác giả của tài liệu, tìm ra:
    - Tìm ra diễn biến theo thời gian của các chủ đề khoa học theo một số tiêu chí như tỷ lệ các tài liệu theo các chủ đề, các chủ đề nổi bật mới, thời điểm một chủ đề cụ thể đạt đỉnh cao nhất, chủ đề nào đang tàn lụi... và theo đó, tìm ra được các chủ đề có vai trò chủ chốt.
    - Nhận biết được các tài liệu có uy thế là tài liệu giới thiệu các ý tưởng mới và có chỉ số ảnh hưởng lớn
    - Nhận biết được tác giả có uy thế là tác giả có ảnh hưởng lớn đối với sự phát triển của các chủ đề.



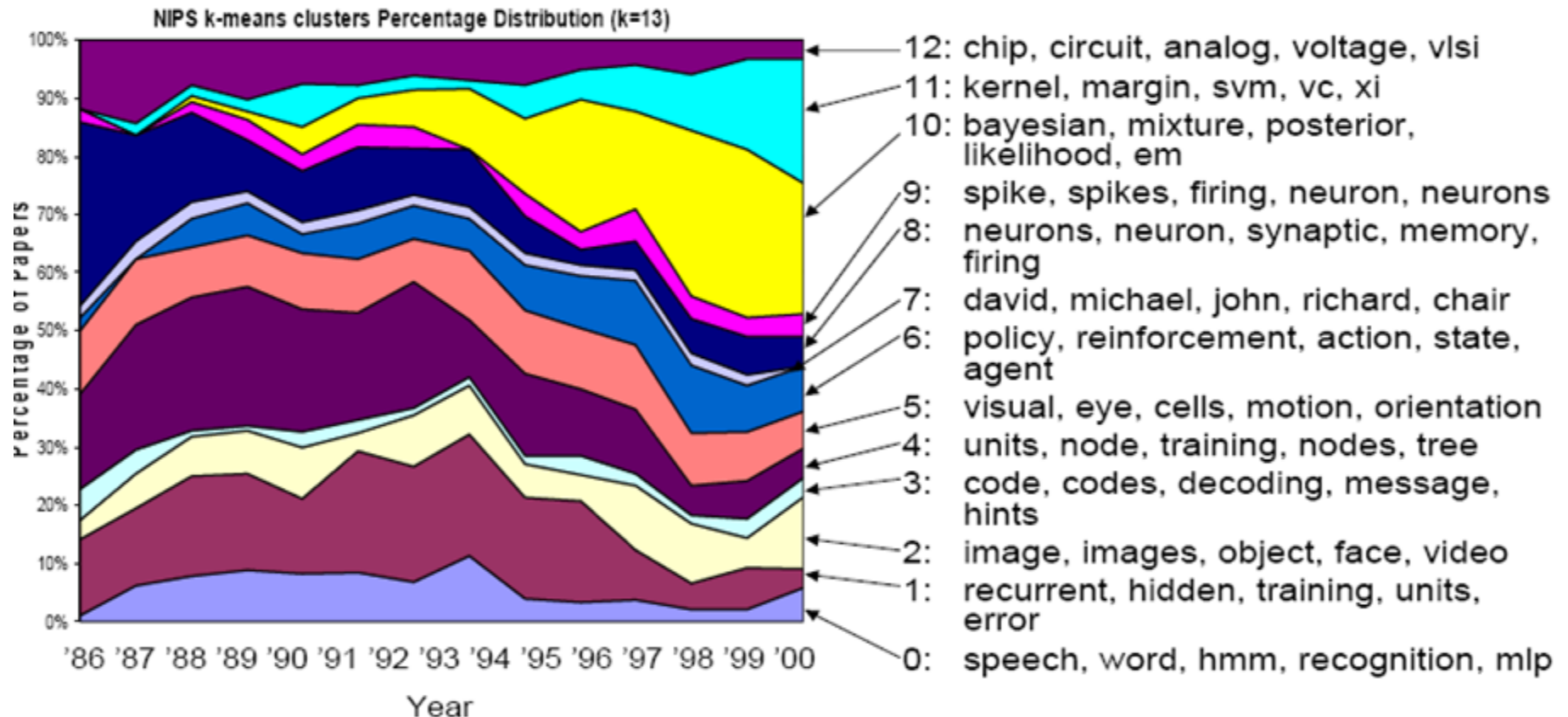
# Một số bài toán ví dụ



- Ví dụ 2. Một kết quả [CJG06]

- Phân cụm tài liệu và gán nhãn cụm (bằng các từ khóa điển hình trong cụm)
- Biểu diễn hình học theo thời gian

## Temporal Cluster Histograms: Results



Các khách hàng sử dụng sản phẩm TextAnalyst của Megaputer

# TextAnalyst

Customer base: 300+ installations

Sample customers

**Ask Jeeves** (USA)

**IMS Health** (USA)

**The Gallup Organization** (USA)

**Centers for Disease Control** (USA)

**Best Buy** (USA)

**France Telecom** (France)

**Skila.com** (USA)

**US Navy** (USA)

**Dow Chemical** (USA)

**Clontech** (USA)

**IMS HEALTH**

**Pfizer** (USA)

**TRW** (USA)

**McKinsey & Company** (USA)

**Liberty Mutual** (USA)

**Logicon** (USA)

**Net Shepherd** (Canada)

**Dept of Environmental Protection** (Australia)

**KPN Research** (Netherlands)

**Talkie.com** (USA)

**NICE Systems** (Israel)





## Ví dụ về Dự án Khai phá Text

### Text Analysis, Decision Forests Link Analysis and more in Polyanalyst 4.5

[KDnuggets](#) : [News](#) : [2004](#) : [n22](#) : [item18](#)

#### Briefs

##### **TEMIS, Leader in Text Mining in Europe, raises 3.6 million euros.**

Paris, November 9th 2004 - TEMIS, provider of corporate Text Mining solutions, has just completed a 3.6 million euro round of financing with ACE Management and Crédit Lyonnais Private Equity (CLPE).

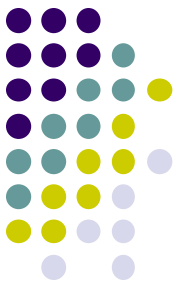
TEMIS develops and markets software solutions for Text Mining. The software transforms free text into usable data, enabling either retrieval of relevant data contained in a document or automatic document classification by topic or by recipient.

At a time when information flow in organizations is constantly increasing, TEMIS' software plays a crucial role in processing this information. It enables decisive productivity gains, in particular in the fields of Competitive Intelligence, analysis of scientific documents, and in customer relationship management.

TEMIS doubled its revenue between 2002 and 2003 and is set to achieve a similar performance for 2004. Its customers include major companies such as Novartis, IPSEN, Total, PSA Peugeot-Citroën, DaimlerChrysler, TIM-Telecom Italia Mobile, etc.

For more information, visit <http://www.temis-group.com/>

# Nghiên cứu về khai phá Text



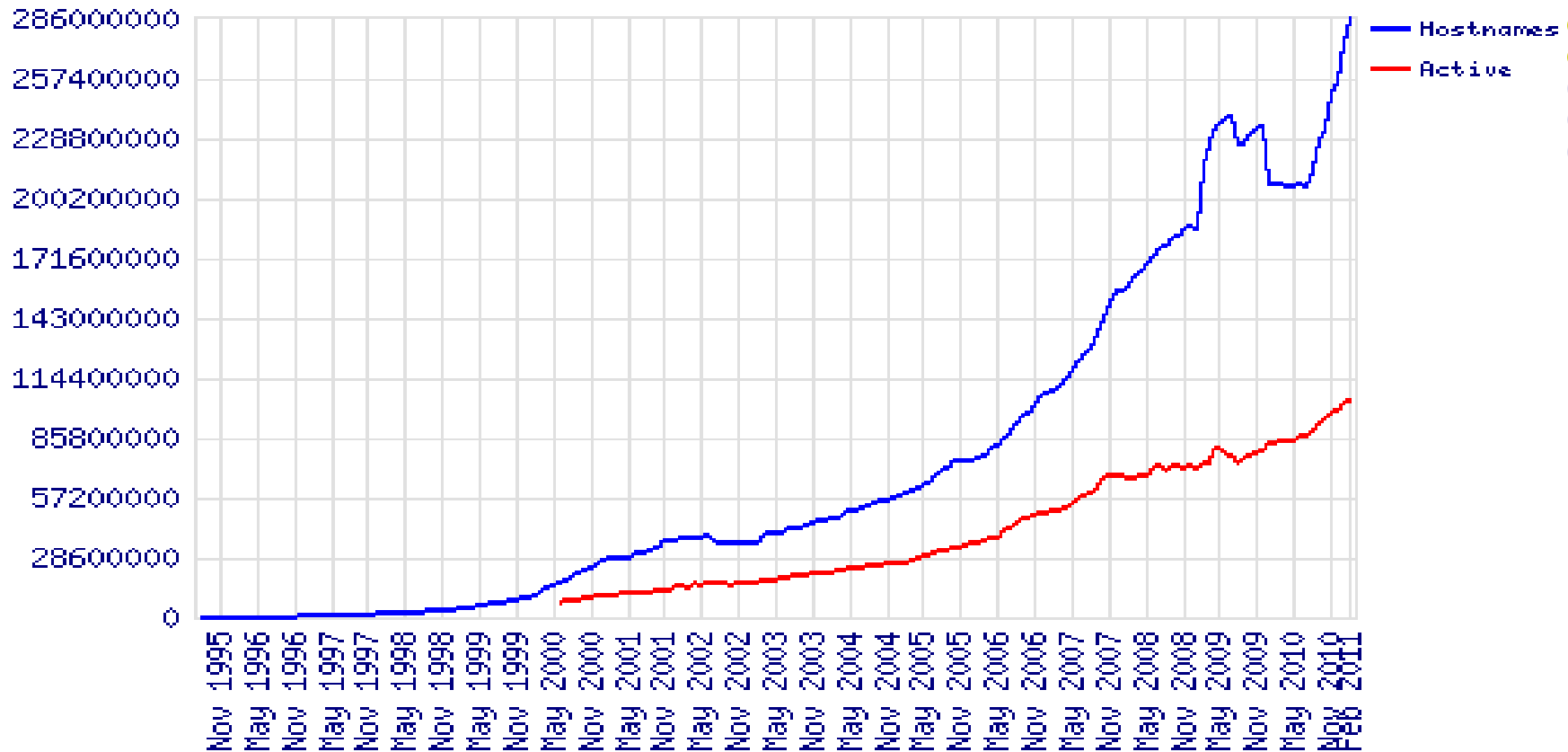
- Theo thống kê từ Google Scholar về số bài viết:
  - Với cụm từ “Text Mining”:
    - Ở tiêu đề: 2.800 bài (khoảng)
    - Ở mọi nơi: 33.000 bài (khoảng)
  - Với cụm từ “Text Analysis”:
    - Ở tiêu đề: 1.680 bài (khoảng)
    - Ở mọi nơi: 43.300 bài (khoảng)
- Nơi công bố tài liệu về Khai phá Text
  - Thường đi kèm với XLNNTN.
  - The ACL Anthology Network Corpus: <http://aclweb.org/anthology-new/>. ACL: “The Association for Computational Linguistics is THE international scientific and professional society for people working on problems involving natural language and computation”.
  - DUC (Document Understanding Conferences: <http://duc.nist.gov/> : 2001-2007) và TAC (Text Analysis Conferences: <http://www.nist.gov/tac/about/index.html>: 2008-nay)
  - Mọi hội nghị, tạp chí khoa học liên quan
  - Kdnuggets: <http://www.kdnuggets.com/>



## 2. Sự cần thiết của khai phá Web

- Web cũng rất gần gũi với con người
  - Tạo ra môi trường của xã hội ảo
  - Một phần quan trọng chuyển tải thông tin của loài người từ Web
  - Phương tiện chuyển giao tri thức
- Đặc thù của khai phá Text và Web
  - Web có bán cấu trúc
  - Kết nối không gian thời gian
  - Mở rộng giao lưu: diễn đàn, blog...
- Sự tăng trưởng của dữ liệu Web
  - Tương tự như dữ liệu Text
  - Dữ liệu đa phương tiện

Total Sites Across All Domains  
August 1995 - February 2011



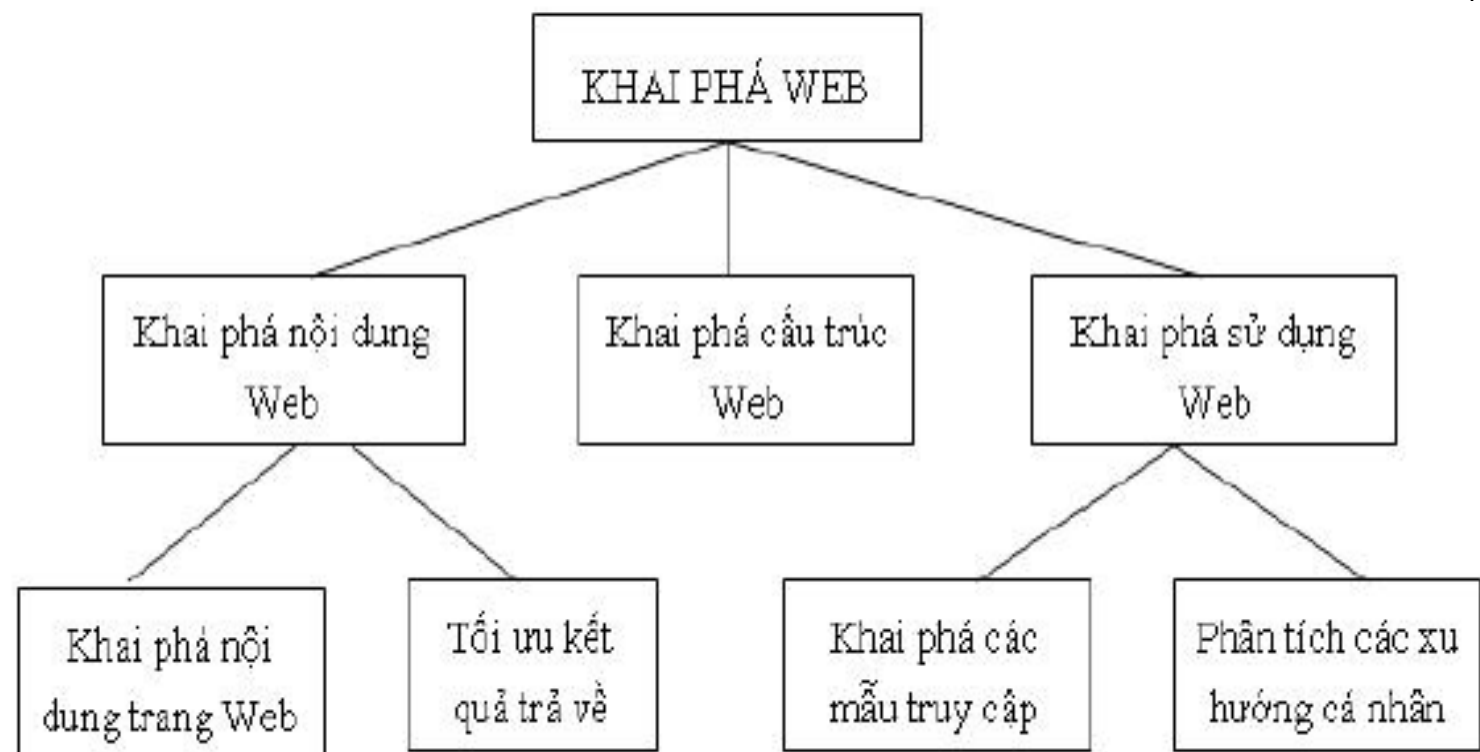
- **Hình minh họa sự tăng trưởng của Web**

- <http://news.netcraft.com/archives/category/web-server-survey/> (02/2011)

- **Khái niệm**

- Khai phá Web = Khai phá Text + WWW
- Trích chọn mẫu mới, hữu ích, hiểu được, tiềm ẩn trong Web





*Các loại khai phá Web*

- Khai phá nội dung Web
  - Khai phá nội dung trang web
  - Tối ưu hiệu quả trả về (hạng, phân cụm)
- Khai phá cấu trúc web: Độ liên quan + cách tổ chức và liên kết
- Khai phá sử dụng web
  - Phân tích các mẫu truy cập (General Access Pattern Tracking)
  - Phân tích các xu hướng cá nhân (Customized Usage tracking)

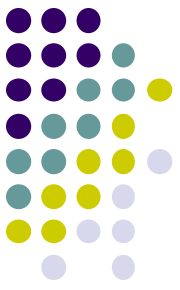
# Các chủ đề của khai phá Web



- Tìm kiếm và thu hồi: Thu hồi và tính hạng
- Phân tích đồ thị Web và Khai phá cấu trúc Web
- Phân cụm Web và Phân lớp Web
- Trích rút thông tin, Quảng cáo và tối ưu hóa Web
- Lọc cộng tác và lọc nội dung
- Phân tích web log và Khai phá sử dụng web
- Mạng xã hội trên Web
- Web ngữ nghĩa
- Khai phá quan điểm trên Web
- Các vấn đề về hệ thống Web

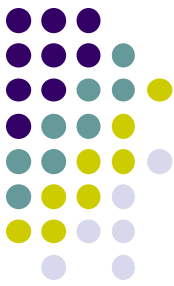


# Một số đặc điểm của khai phá Web



- Web quá lớn để tổ chức thành kho dữ liệu
  - Tăng kích cỡ DW chậm hơn nhiều tốc độ phát triển Web
- Độ phức tạp của trang Web là rất lớn
  - Các kiểu tổ chức
  - Các kiểu dữ liệu
- Web: nguồn tài nguyên thông tin có độ thay đổi cao
  - Tăng nhiều và mất nhiều
- Web phục vụ một cộng đồng người rộng lớn và đa dạng
  - Phản ánh toàn bộ thế giới
- Chỉ phần rất nhỏ thông tin trên Web là thực sự hữu ích
  - Đối với toàn bộ và từng cá nhân
- Khai phá Web có lợi thế: bán cấu trúc, giàu thông tin (thẻ, liên kết, file log)

# Nghiên cứu về khai phá Web

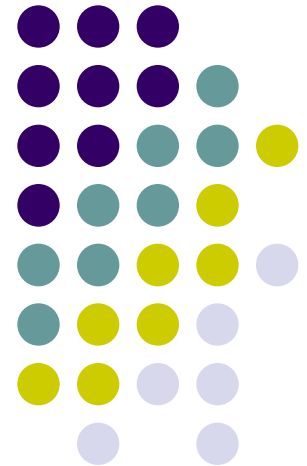


- Theo thống kê từ Google Scholar về số bài viết:
  - Với cụm từ “Web Mining”:
    - Ở tiêu đề: 2.680 bài (khoảng)
    - Ở mọi nơi: 20.000 bài (khoảng)
  - Với cụm từ “Text Analysis”:
    - Ở tiêu đề: 240 bài (khoảng)
    - Ở mọi nơi: 4.300 bài (khoảng)
  - Với cụm từ “Search Engine”:
    - Ở tiêu đề: 6.260 bài (khoảng)
    - Ở mọi nơi: 414.000 bài (khoảng)
  - Với cụm từ “Image Search”:
    - Ở tiêu đề: 890 bài (khoảng)
    - Ở mọi nơi: 15.800 bài (khoảng)
- Nơi công bố tài liệu về Khai phá Web
  - Đi kèm với XLNNTN và khai phá Text
  - Kdnuggets: <http://www.kdnuggets.com/>
  - Mọi hội nghị, tạp chí khoa học liên quan

# BÀI GIẢNG KHAI PHÁ DỮ LIỆU WEB

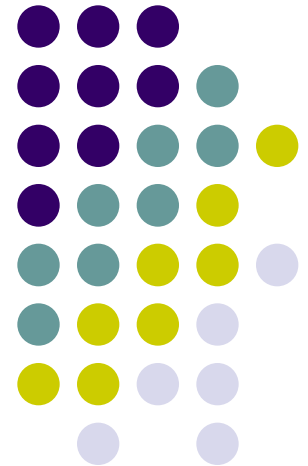
## CHƯƠNG 2. KHAI PHÁ SỬ DỤNG WEB VÀ KHAI PHÁ CẤU TRÚC WEB

PGS. TS. HÀ QUANG THỤY  
HÀ NỘI 10-2010  
**TRƯỜNG ĐẠI HỌC CÔNG NGHỆ**  
ĐẠI HỌC QUỐC GIA HÀ NỘI



# Nội dung

1. Khai phá sử dụng Web
2. Khai phá cấu trúc web





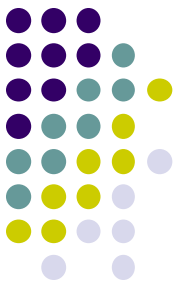
# 1. Khai phá sử dụng Web

- Giới thiệu chung
- Phân tích mẫu truy nhập Web
  - Mang tính thói quen có tính cộng đồng
  - Khai phá mẫu truy nhập theo luật kết hợp
- Khai phá xu hướng sử dụng
  - Cá nhân hóa
  - Các hệ tư vấn

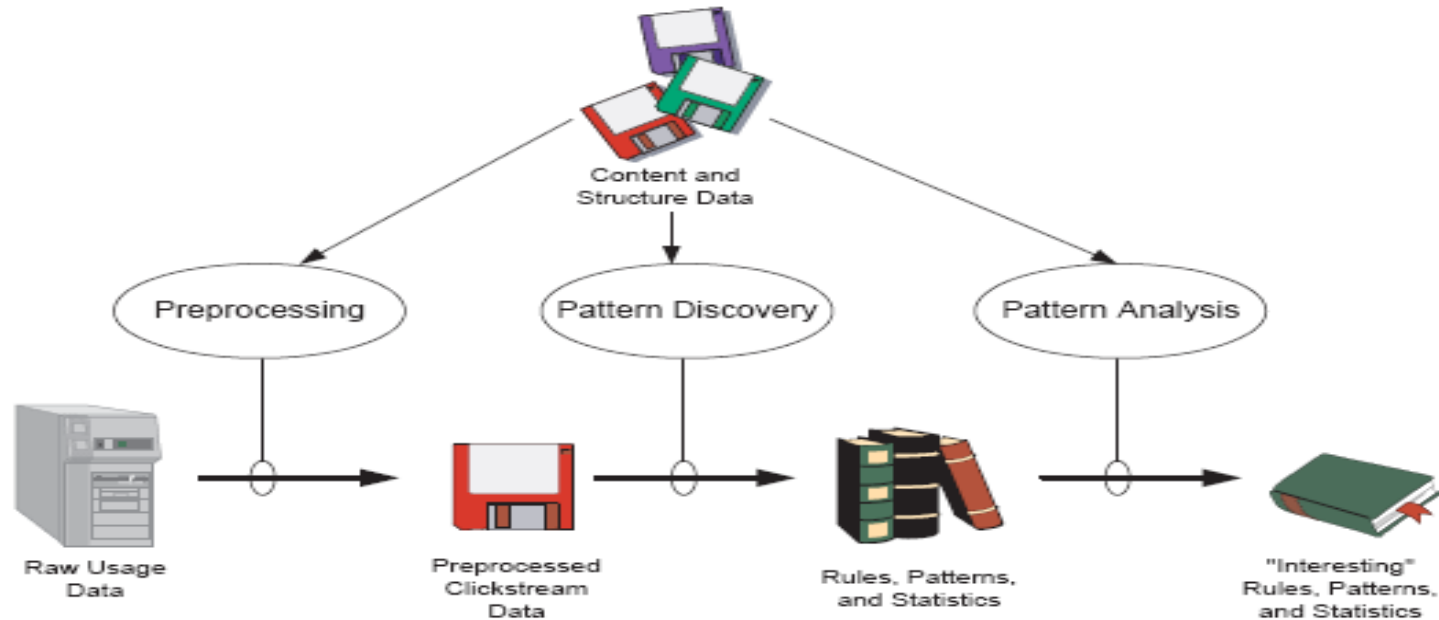
# 1.a. Giới thiệu chung



- **Nguồn dữ liệu**
  - Các logfile (máy chủ, máy khách, máy trung gian)
  - CSDL khách hàng
- **Mô hình dữ liệu**
  - Thực thể: người sử dụng, khung nhìn trang web, file trang Web, trình duyệt, phục vụ web, phục vụ nội dung, phiên người sử dụng, phiên phục vụ, dãy các sự kiện liên quan (episode).
- **Tiền xử lý dữ liệu**
  - Loại: cấu trúc, nội dung
  - Bài toán: xử lý văn bản, rút gọn đặc trưng, mô hình dữ liệu.
- **Phát hiện mẫu**
  - Mẫu quan hệ: thống kê, luật kết hợp, luật chuỗi, phân cụm, phân lớp, mô hình phụ thuộc
  - Đại chúng và cá nhân hóa



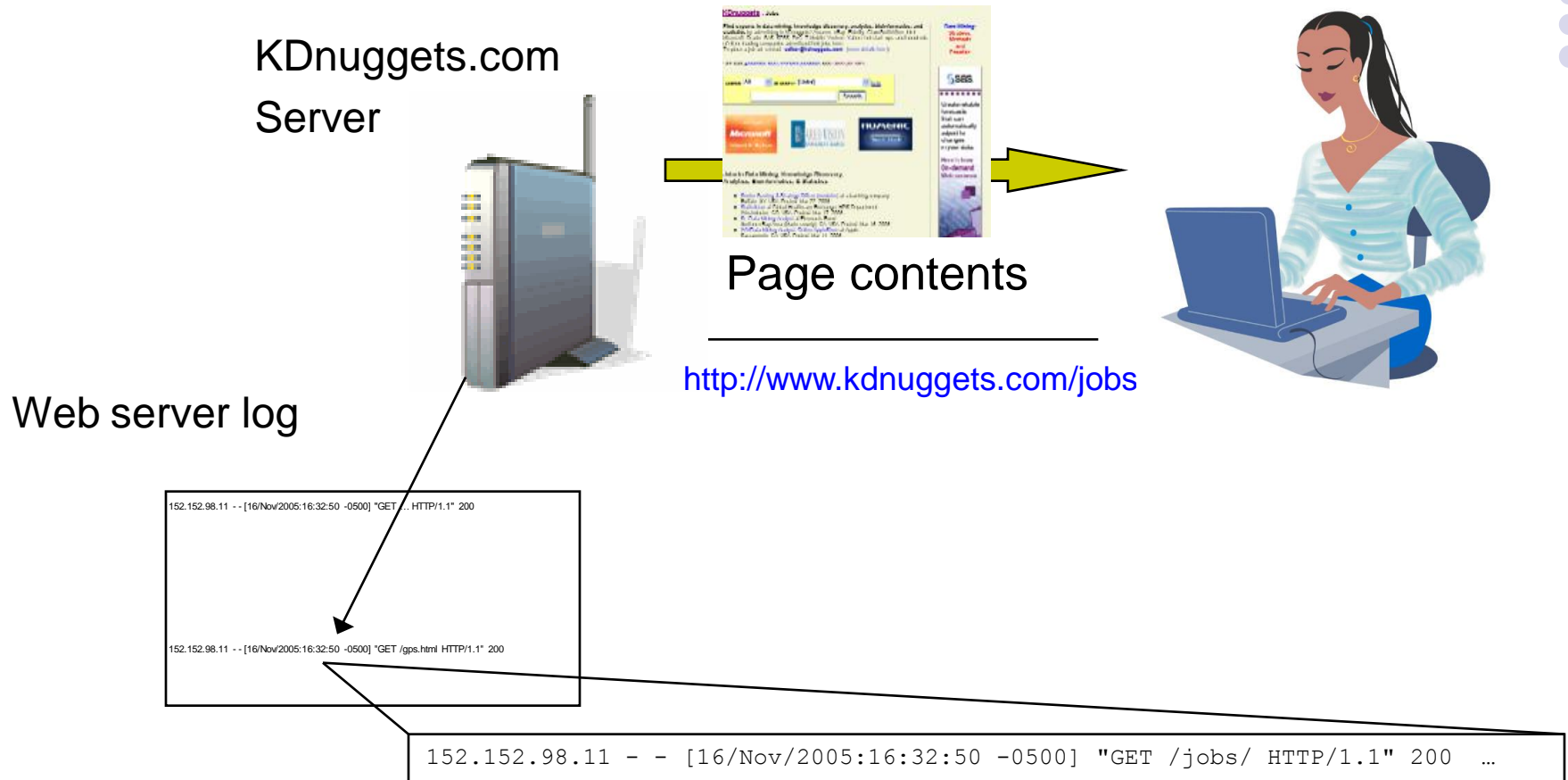
# 1.a. Một quy trình khai phá sử dụng Web



## Quá trình khai phá sử dụng Web [Coo00]

- Input: Dữ liệu sử dụng Web
- Output: Các luật, mẫu, thống kê hấp dẫn
- Các bước chủ yếu:
  - Tiền xử lý dữ liệu
  - Khám phá mẫu
  - Phân tích mẫu

# Sơ đồ ghi dữ liệu vào logfile



- **Thông tin truy nhập người dùng**

- Server tổ chức ghi nhận vào logfile
- Hỗ trợ quản lý điều hành
- Tài nguyên Khai phá dữ liệu, nâng cao hiệu năng hệ thống



# Một dòng ví dụ trong weblog



152.152.98.11 -- [16/Nov/2005:16:32:50 -0500] "GET /jobs/ HTTP/1.1" 200 15140  
"http://www.google.com/search?q=salary+for+data+mining&hl=en&lr=&start=10&sa=N" "Mozilla/4.0  
(compatible; MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4322)"

**152.152.98.11** Địa chỉ của hotname

-- Tên và login của người dùng từ xa: thường là "--

**[16/Nov/2005:16:32:50 -0500]** Ngày và giờ truy nhập.

*Giờ GMT: (+/-)HH00 US UST: -500*

**"GET /jobs/ HTTP/1.1"** Phương thức lấy thông tin, URL liên quan tới tên miền; giao thức

**200** **Trạng thái** 200 – OK (hầu hết, đạt được) | 206 – truy nhập bộ phận – chuyển hướng vĩnh viễn (truy nhập tới/ tiến trình định hướng lại /tiến trình/ )| 302 – định hướng tạm thời| 304 – không thay đổi | 404 – không thấy|...

**15140** Dung lượng tải về máy khách | "-- nếu trạng thái 304

**"http://www.google.com/search?q=salary+for+data+mining&hl=en &lr=&start=10&sa=N"** URL của người thăm (ở đây là từ Google)

**"Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4322)"** đại lý của người dùng

# Một ví dụ về log files

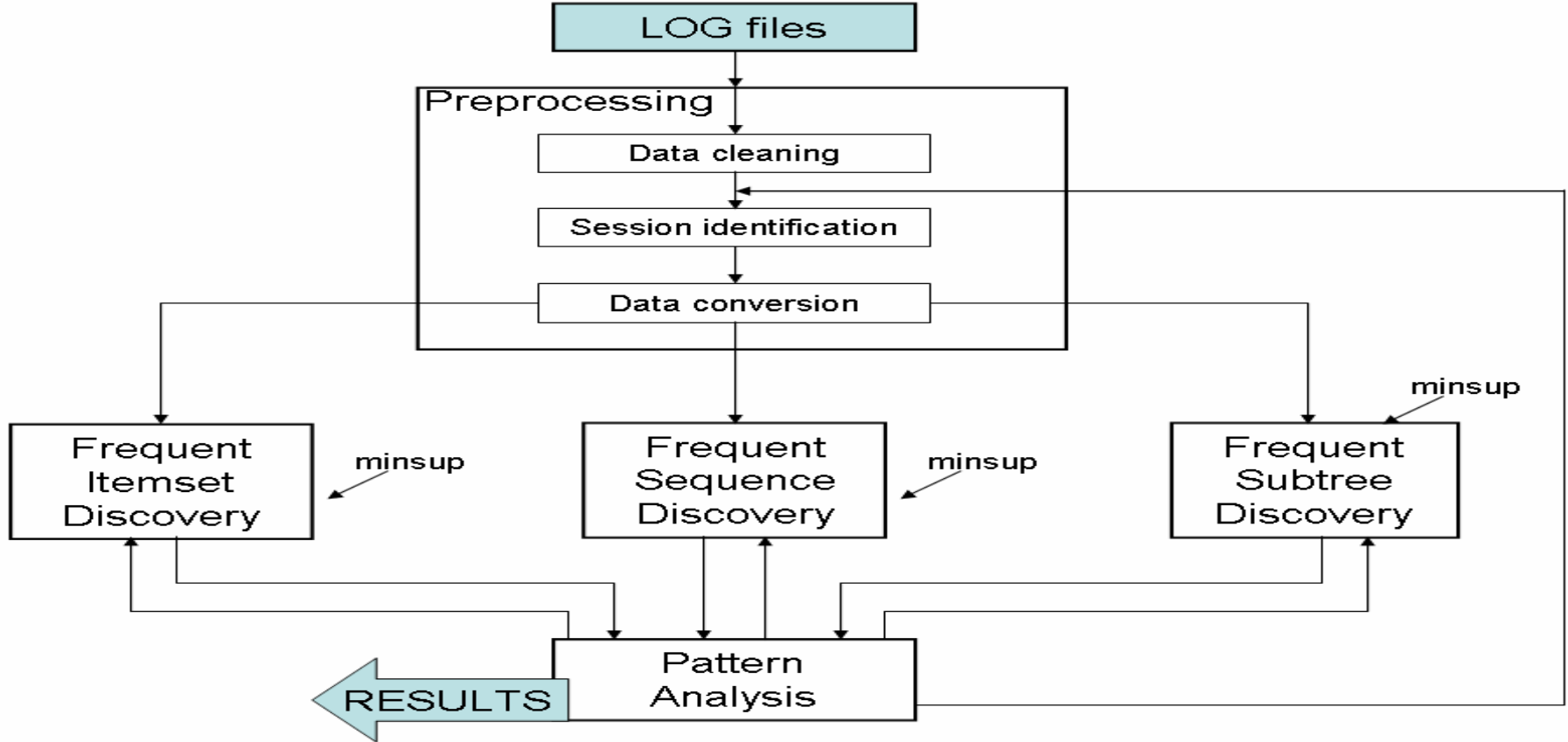


AnonID	Query	QueryTime	Rank	ClickURL
479	family guy movie references	2006-03-03 22:37:46	1	<a href="http://www.familyguyfiles.com">http://www.familyguyfiles.com</a>
479	top grossing movies of all time	2006-03-03 22:42:42	1	<a href="http://movieweb.com">http://movieweb.com</a>
479	top grossing movies of all time	2006-03-03 22:42:42	2	<a href="http://www.imdb.com">http://www.imdb.com</a>
479	car decals	2006-03-03 23:20:12	4	<a href="http://www.decaljunky.com">http://www.decaljunky.com</a>
479	car decals	2006-03-03 23:20:12	1	<a href="http://www.modernimage.net">http://www.modernimage.net</a>
479	car decals	2006-03-03 23:20:12	5	<a href="http://www.webdecal.com">http://www.webdecal.com</a>
479	car window decals	2006-03-03 23:24:05	9	<a href="http://www.customautotrim.com">http://www.customautotrim.com</a>
479	car window sponsor decals	2006-03-03 23:27:17	3	<a href="http://www.streetglo.net">http://www.streetglo.net</a>
479	bose	2006-03-03 23:30:11	1	<a href="http://www.bose.com">http://www.bose.com</a>
479	bose car decal	2006-03-03 23:31:48	1	<a href="http://stickers.signprint.co.uk">http://stickers.signprint.co.uk</a>
479	bose car decal	2006-03-03 23:31:48	1	<a href="http://stickers.signprint.co.uk">http://stickers.signprint.co.uk</a>
479	bose car decal	2006-03-03 23:31:48	7	<a href="http://www.motorcitydecals.com">http://www.motorcitydecals.com</a>
479	chicago the mix	2006-03-04 22:11:31	1	<a href="http://www.wtmx.com">http://www.wtmx.com</a>
479	chicago the drive	2006-03-04 22:14:51	2	<a href="http://www.wdrv.com">http://www.wdrv.com</a>

```
q          = cars
URL        = www.google.com/search?q=cars
IP         = 72.14.253.103
Cookie     = PREF=ID=03b1d4f329293203:LD=en:NR=10...
Browser    = Firefox/2.0.0.4;Windows NT 5.1
Time       = 25 Mar 2007 10:15:32
```

Một phần query log của AOL (trên) và Cấu trúc log của Google (dưới)

# 1.b. Phân tích mẫu truy nhập



## ● Phân tích mẫu từ logfile

- Tìm tập mục phổ biến, dãy phổ biến, cây con phổ biến
- Phân tích mẫu phổ biến tìm được

[IV06] Renáta Iváncsy, István Vajk (2006). Frequent Pattern Mining in Web Log Data, *Acta Polytechnica Hungarica*, 3(1):77-90.



# 1.b. Ví dụ về mẫu phổ biến sử dụng Web

opinion & misc & travel	→ on-air	90.26%
news & misc & business & bbs	→ frontpage	90.24%
living & business & sports & bbs	→ frontpage	90.00%
news & misc & business & sports	→ frontpage	89.68%
news & tech & living & business & sports	→ frontpage	89.00%
news & living & business & bbs	→ frontpage	88.01%
frontpage & tech & living & business & sports	→ news	87.87%
frontpage & opinion & living & sports	→ news	87.81%
frontpage & tech & opinion & living	→ news	87.60%
frontpage & tech & on-air & business & sports	→ news	87.59%
news & misc & sports & bbs	→ frontpage	87.56%
news & tech & on-air & business & sports	→ frontpage	87.43%
news & living & business & sports	→ frontpage	87.18%
news & business & sports & bbs	→ frontpage	86.70%
misc & living & travel	→ on-air	86.56%
tech & living & sports & bbs	→ frontpage	86.52%
tech & business & sports & bbs	→ frontpage	86.40%
news & misc & living & business	→ frontpage	86.22%
on-air & business & sports & bbs	→ frontpage	86.22%
news & tech & misc & bbs	→ frontpage	86.18%
on-air & misc & business & sports	→ frontpage	86.16%
tech & misc & travel	→ on-air	86.09%
tech & living & business & sports	→ frontpage	86.08%
news & living & sports & bbs	→ frontpage	85.99%
misc & business & sports	→ frontpage	85.79%
frontpage & tech & opinion & sports	→ news	85.78%
news & opinion & living & sports	→ frontpage	85.69%
misc & business & travel	→ on-air	85.66%
news & tech & misc & business	→ frontpage	85.63%
misc & business & bbs	→ frontpage	85.57%
tech & living & sports & bbs	→ news	85.49%
local & misc & business & sports	→ frontpage	85.43%
news & opinion & business & bbs	→ frontpage	85.32%
news & misc & living & sports	→ frontpage	85.19%
news & on-air & business & sports	→ frontpage	85.01%

(a)

misc → local	2.07%
frontpage → frontpage → sports	2.02%
local → frontpage	1.83%
on-air → misc → on-air	1.72%
on-air → frontpage	1.69%
on-air → news	1.51%
news → frontpage → news	1.49%
local → news	1.46%
frontpage → frontpage → business	1.35%
news → sports	1.33%
news → bbs	1.23%
health → local	1.16%
misc → frontpage → frontpage	1.16%
on-air → local	1.15%
misc → on-air → misc	1.15%
frontpage → frontpage → living	1.14%
local → frontpage → frontpage	1.13%
health → misc	1.12%
misc → on-air → on-air	1.10%
local → misc → local	1.09%
misc → news	1.06%
news → living	1.06%
on-air → misc → on-air → misc	1.00%

(b)

# 1.b. Ví dụ về mẫu kết hợp

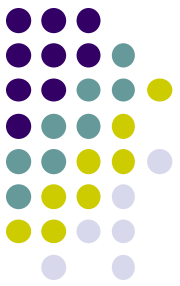


## ❖ Một số ví dụ về “luật kết hợp” (*associate rule*)

- “98% khách hàng mà mua tạp chí thể thao thì đều mua các tạp chí về ô tô” ⇒ sự **kết hợp** giữa “tạp chí thể thao” với “tạp chí về ô tô”
- “60% khách hàng mà mua bia tại siêu thị thì đều mua bím trẻ em” ⇒ sự **kết hợp** giữa “bia” với “bím trẻ em”
- “Có tới 70% người truy nhập Web vào địa chỉ *Url1* thì cũng vào địa chỉ *Url2* trong một phiên truy nhập web” ⇒ sự **kết hợp** giữa “*Url 1*” với “*Url 2*”. Khai phá dữ liệu sử dụng Web (lấy dữ liệu từ file log của các site, chẳng hạn được MS cung cấp). Các Url có gắn với nhãn “lớp” là các đặc trưng thì có luật kết hợp liên quan giữa các lớp Url này.

## ❖ Khái niệm cơ sở về luật kết hợp

# Khai phá luật kết hợp: Cơ sở



## Cơ sở dữ liệu giao dịch (transaction database)

- Tập toàn bộ các mục  $I = \{i_1, i_2, \dots, i_k\}$ : “tất cả các mặt hàng”.
- *Giao dịch*: danh sách các mặt hàng (mục: item) trong một phiếu mua hàng của khách hàng. Giao dịch  $T$  là một tập mục.

Một giao dịch  $T$  là một tập con của  $I$ :  $T \subseteq I$ . Mỗi giao dịch  $T$  có một định danh là  $T_{ID}$ .

- $A$  là một tập mục  $A \subseteq I$  và  $T$  là một giao dịch: Gọi  $T$  chứa  $A$  nếu  $A \subseteq T$ .

# Khai phá luật kết hợp: cơ sở



## Luật kết hợp

- Gọi  $A \rightarrow B$  là một “luật kết hợp” nếu  $A \neq B$  và  $A \cap B = \emptyset$
- Luật kết hợp  $A \rightarrow B$  có độ hỗ trợ (support)  $s$  trong CSDL giao dịch  $D$  nếu trong  $D$  có  $s\%$  các giao dịch  $T$  chứa  $AB$ : chính là xác suất  $P(AB)$ .

Tập mục  $A$  có  $P(A) > 0$  (với  $s$  cho trước) được gọi là tập phổ biến (*frequent set*).

- Luật kết hợp  $A \rightarrow B$  có độ tin cậy (confidence)  $c$  trong CSDL  $D$  nếu như trong  $D$  có  $c\%$  các giao dịch  $T$  chứa  $A$  thì cũng chứa  $B$ : chính là xác suất  $P(B|A)$ .

$$\text{Support}(A \rightarrow B) = P(AB) : 1 - P(A \rightarrow \emptyset)$$

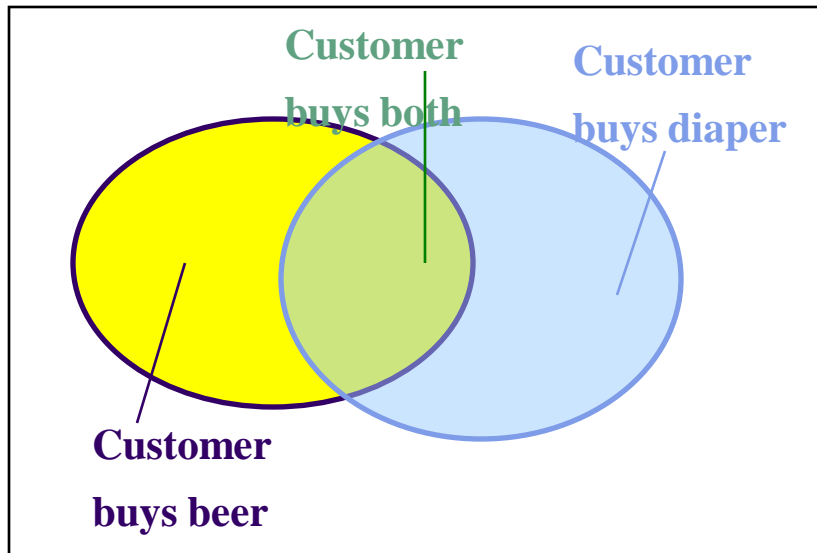
$$\text{Confidence}(A \rightarrow B) = P(B|A) : 1 - P(\emptyset | A)$$

- Luật  $A \rightarrow B$  được gọi là đảm bảo độ hỗ trợ  $s$  trong  $D$  nếu  $s(A \rightarrow B) \geq s$ . Luật  $A \rightarrow B$  được gọi là đảm bảo độ tin cậy  $c$  trong  $D$  nếu  $c(A \rightarrow B) \geq c$ . Tập mạnh.



# Ví dụ: Mẫu phổ biến và luật kết hợp

Transaction-id	Items bought
10	A, B, C
20	A, C
30	A, D
40	B, E, F



- Tập mục  $I = \{i_1, \dots, i_k\}$ . CSDL giao dịch  $D = \{d \dots\}$
- $A, B \dots A \dots = \dots A \rightarrow B$  là luật kết hợp
- Bài toán tìm luật kết hợp.

Cho trước độ hỗ trợ tối thiểu  $s > 0$ , độ tin cậy tối thiểu  $c > 0$ . Hãy tìm mọi luật kết hợp mạnh  $X \rightarrow Y$ .

Giả sử  $min\_support = 50\%$ ,  $min\_conf$

$= 50\%$ :

$A \rightarrow C$  (50%, 66.7%)

$C \rightarrow A$  (50%, 100%)

- Hãy trình bày các nhận xét về khái niệm luật kết hợp với khái niệm phụ thuộc hàm.
- Các tính chất Armstrong ở đây.





# Một ví dụ tìm luật kết hợp

Transaction-id	Items bought
10	A, B, C
20	A, C
30	A, D
40	B, E, F

Min. support 50%  
Min. confidence 50%

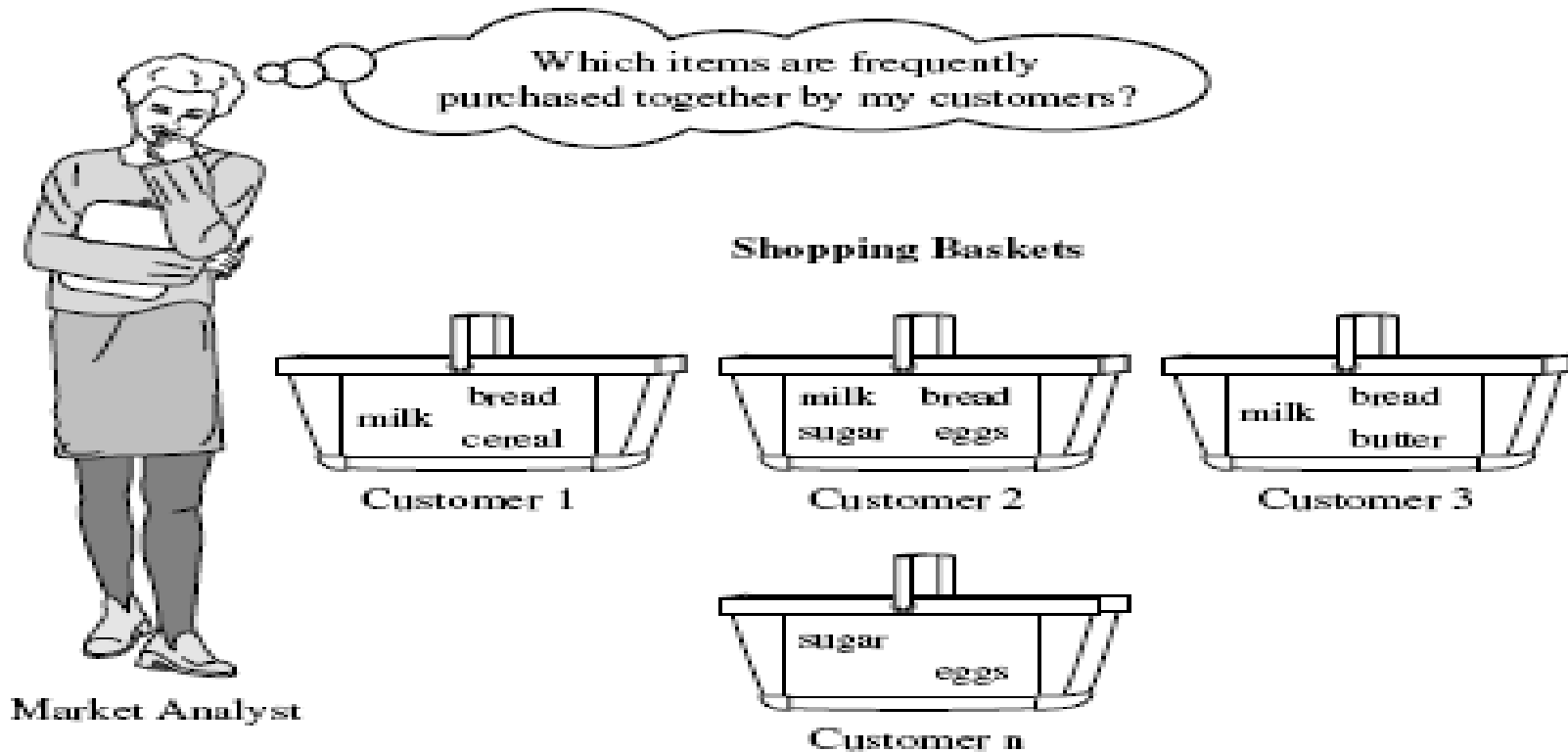
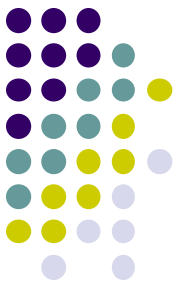
Frequent pattern	Support
{A}	75%
{B}	50%
{C}	50%
{A, C}	50%

For rule  $A \Rightarrow C$ :

$$\text{support} = \text{support}(\{A, C\}) = 50\%$$

$$\text{confidence} = \text{support}(\{A, C\}) / \text{support}(\{A\}) = 66.6\%$$

# Khai niệm khai phá kết hợp



*computer*  $\Rightarrow$  *antivirus\_software* [support = 2%, confidence = 60%]

# Khai phá luật kết hợp



- **Khai phá luật kết hợp:**
  - Tìm tất cả mẫu phổ biến, kết hợp, tương quan, hoặc cấu trúc nhan-quả trong tập các mục hoặc đối tượng trong CSDL quan hệ hoặc các kho chứa thông tin khác.
  - **Mẫu phổ biến (Frequent pattern):** là mẫu (tập mục, dãy mục...) mà xuất hiện phổ biến trong 1 CSDL [AIS93]
- **Động lực: tìm mẫu chính quy (regularities pattern) trong DL**
  - Các mặt hàng nào được mua cùng nhau? — Bia và bỉm (diapers)?!
  - Mặt hàng nào sẽ được mua sau khi mua một PC ?
  - Kiểu DNA nào nhạy cảm với thuốc mới này?
  - Có khả năng tự động phân lớp Web hay không ?

# Mẫu phổ biến và khai phá luật kết hợp là một bài toán bản chất của khai phá DL



- **Nền tảng của nhiều bài toán KPDL bản chất**
  - Kết hợp, tương quan, nhân quả
  - Mẫu tuần tự, kết hợp thời gian hoặc vòng, chu kỳ bộ phận, kết hợp không gian và đa phương tiện
  - Phân lớp kết hợp, phân tích cụm, khối tăng băng, tích tụ (nén dữ liệu ngữ nghĩa)
- **Ứng dụng rộng rãi**
  - Phân tích DL bóng rổ, tiếp thị chéo (cross-marketing), thiết kế catalog, phân tích chiến dịch bán hàng
  - Phân tích Web log (click stream), Phân tích chuỗi DNA v.v.

# Apriori: Một tiếp cận sinh ứng viên và kiểm tra



- Khái quát: Khai phá luật kết hợp gồm hai bước:
  - Tìm mọi tập mục phổ biến: theo min-sup
  - Sinh luật mạnh từ tập mục phổ biến
- Mọi tập con của tập mục phổ biến cũng là tập mục phổ biến
  - Nếu  $\{bia, bĩm, hạnh nhân\}$  là phổ biến thì  $\{bia, bĩm\}$  cũng vậy: Mọi giao dịch chứa  $\{bia, bĩm, hạnh nhân\}$  cũng chứa  $\{bia, bĩm\}$ .
- Nguyên lý của Apriori: Với mọi tập mục không phổ biến thì mọi tập bao không cần phải sinh ra/kiểm tra!
- Phương pháp:
  - Sinh các tập mục ứng viên dài  $(k+1)$  từ các tập mục phổ biến có độ dài  $k$  (Độ dài tập mục là số phần tử của nó),
  - Kiểm tra các tập ứng viên theo CSDL
- Các nghiên cứu hiệu năng chứng tỏ tính hiệu quả và khả năng mở rộng của thuật toán
- Agrawal & Srikant 1994, Mannila, và cộng sự 1994



# Thuật toán Apriori

❖ Trên cơ sở tính chất (nguyên lý tủa) Apriori, thuật toán hoạt động theo quy tắc quy hoạch động

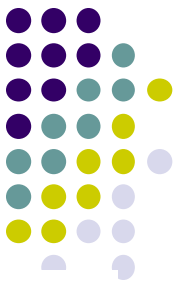
□ Từ các tập  $F_i = \{c_i \mid c_i \text{ tập phổ biến, } |c_i| = i\}$  gồm mọi tập mục phổ biến có độ dài  $i$  với  $1 \leq i \leq k$ ,

□ đi tìm tập  $F_{k+1}$  gồm mọi tập mục phổ biến có độ dài  $k+1$ .

❖ Trong thuật toán:

các tên mục  $i_1, i_2, \dots, i_n$  ( $n = |I|$ ) được sắp xếp theo một thứ tự cố định: thường được đánh chỉ số  $1, 2, \dots, n$ .

# Thuật toán Apriori



Thuật toán Apriori [WKQ08]:

Input:     - Cơ sở dữ liệu giao dịch  $D = \{t \mid t \text{ giao dịch}\}$   
          - Độ hỗ trợ tối thiểu  $\text{minsup} > 0$

Output:    - Tập hợp tất cả các tập phổ biến.

```
0:  mincount = minsup * |D|;
1.   $F_1 = \{\text{các tập phổ biến có độ dài } 1\}$ 
2.  for (k=1;  $F_k \neq \emptyset$ ; k++) do begin
3.       $C_{k+1} = \text{apriori-gen}(F_k)$ ; // sinh mọi ứng viên độ dài k+1
4.      for t  $\in D$  do begin
5.           $C_t = \{c \in C_{k+1} \mid c \subseteq t\}$ ; //mọi ứng viên chứa trong t
6.          for c  $\in C_t$  do
7.              c.count ++;
8.          end
9.           $F_{k+1} = \{c \in C_{k+1} \mid \text{c.count} \geq \text{mincount}\}$  ;
10. end
11. Answer  $\cup_k F_k$  ;
```

# Thuật toán: Thủ tục con Apriori-gen



Trong mỗi bước  $k$ , thuật toán Apriori đều phải duyệt CSDL  $D$ .

Khởi động, duyệt  $D$  để có được  $F_1$ .

Các bước  $k$  sau đó, duyệt  $D$  để tính số lượng giao dịch  $t$  thoả từng ứng viên  $c$  của  $C_{k+1}$ : mỗi giao dịch  $t$  chỉ xem xét một lần cho mọi ứng viên  $c$  thuộc  $C_{k+1}$ .

## Thủ tục con Apriori-gen sinh tập phổ biến: *tự tương*

Bước nối: Sinh các tập mục  $R_{k+1}$  là ứng viên tập phổ biến có độ dài  $k+1$  bằng cách kết hợp hai tập phổ biến  $P_k$  và  $Q_k$  có độ dài  $k$  và trùng nhau ở  $k-1$  mục đầu tiên:

$$R_{k+1} = P_k \cup Q_k = \{i_1, i_2, \dots, i_{k-1}, i_k, i_{k'}\} \text{ với}$$

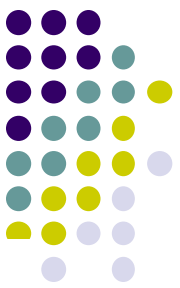
$$P_k = \{i_1, i_2, \dots, i_{k-1}, i_k\} \text{ và } Q_k = \{i_1, i_2, \dots, i_{k-1}, i_{k'}\}$$

trong đó  $i_1 \leq i_2 \leq \dots \leq i_{k-1} \leq i_k \leq i_{k'}$ .

Bước tia: Giữ lại tất cả các  $R_{k+1}$  thỏa tính chất Apriori ( $\forall X \subseteq R_{k+1}$  và  $|X|=k \Rightarrow X \in F_k$ ), nghĩa là đã loại (tia) bớt đi mọi ứng viên  $R_{k+1}$  không đáp ứng tính chất này.



# Thuật toán Apriori-gen



```
(1) for mọi tập mục phổ biến  $l_1 \in L_k$ 
(2) for mọi tập mục phổ biến  $l_2 \in L_k$ 
(3) if  $(l_1[1]=l_2[1]) \wedge (l_1[2]=l_2[2]) \wedge \dots \wedge (l_1[k-1]=l_2[k-1]) \wedge (l_1[k] < l_2[k])$ 
    then {
         $c = l_1 \cup l_2$ ; // join step: generate candidates
        //  $c = \{l_1[1], l_1[2], \dots, l_1[k-1], l_1[k], l_2[k]\}$ 
(5)    if has_infrequent_subset( $c, L_k$ ) then
(6)        delete  $c$ ; // bước thử: bỏ ứng viên không đúng
        else add  $c$  to  $C_{k+1}$ ;
(8)    }
(9) return  $C_k$ ;
```

procedure has\_infrequent\_subset( $c$ : tập ứng viên độ dài  $k+1$ ;

$L_k$ : tập các tập mục phổ biến độ dài  $k$ ); // tri thức đã có

```
(1) for mỗi tập con  $s$  độ dài  $k$  của  $c$ 
(2)    if  $s \notin L_k$  then
(3)        return TRUE;
(4) return FALSE;
```

# Một ví dụ thuật toán Apriori ( $s=0.5$ )



Database TDB

Tid	Items
10	A, C, D
20	B, C, E
30	A, B, C, E
40	B, E

$C_1$

1<sup>st</sup> scan

Itemset	sup
{A}	2
{B}	3
{C}	3
{D}	1
{E}	3

$L_1$

Itemset	sup
{A}	2
{B}	3
{C}	3
{E}	3

$C_2$

2<sup>nd</sup> scan

Itemset	sup
{A, B}	1
{A, C}	2
{A, E}	1
{B, C}	2
{B, E}	3
{C, E}	2

$C_2$

Itemset
{A, B}
{A, C}
{A, E}
{B, C}
{B, E}
{C, E}

$L_2$

Itemset	sup
{A, C}	2
{B, C}	2
{B, E}	3
{C, E}	2

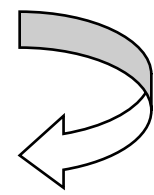
$C_3$

Itemset
{B, C, E}

3<sup>rd</sup> scan

$L_3$

Itemset	sup
{B, C, E}	2



# Chi tiết quan trọng của Apriori

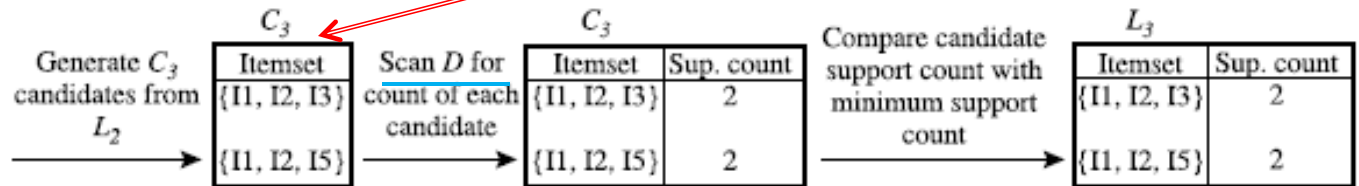
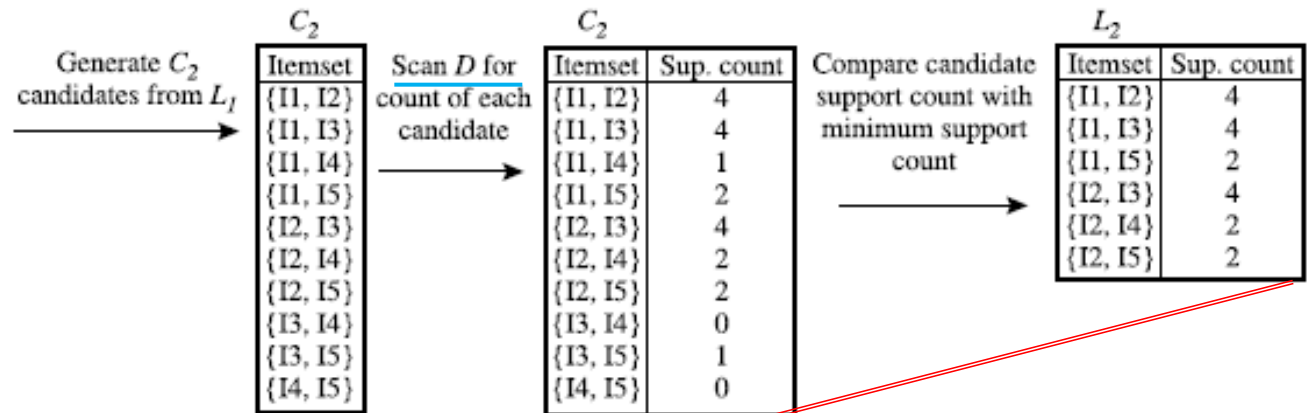
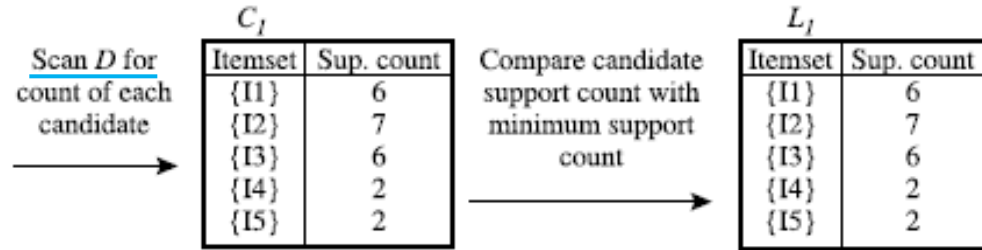


- Cách thức sinh các ứng viên:
  - Bước 1: Tụ kết nối  $L_k$
  - Step 2: Cắt tỉa
- Cách thức đếm hỗ trợ cho mỗi ứng viên.
- Ví dụ thủ tục con sinh ứng viên
  - $L_3 = \{abc, abd, acd, ace, bcd\}$
  - Tụ kết nối:  $L_3 * L_3$ 
    - $abcd$  từ  $abc$  và  $abd$
    - $acde$  từ  $acd$  và  $ace$
  - Tỉa:
    - $acde$  là bỏ đi vì  $ade$  không thuộc  $L_3$
  - $C_4 = \{abcd\}$

# Ví dụ: $D, \text{min\_sup} * |D| = 2$ ( $C_4 = \square$ )



TID	List of item_IDs
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3



# Sinh luật kết hợp



Việc sinh luật kết hợp gồm hai bước

- Với mỗi tập phổ biến  $W$  tìm được hãy sinh ra mọi tập con thực sự  $X$  khác rỗng của nó.
- Với mỗi tập phổ biến  $W$  và tập con  $X$  khác rỗng thực sự của nó: sinh luật  $X \Rightarrow W - X$  nếu  $P(W-X|X) \geq c$ .

Như ví dụ đã nêu có  $L3 = \{\{I1, I2, I3\}, \{I1, I2, I5\}\}$

Với độ tin cậy tối thiểu 70%, xét tập mục phổ biến  $\{I1, I2, I5\}$  có 3 luật như dưới đây: **Duyệt CSDL ?**

<i>TID</i>	<i>List of item_IDs</i>
T100	I1, I2, I5
T200	I2, I4
T300	I2, I3
T400	I1, I2, I4
T500	I1, I3
T600	I2, I3
T700	I1, I3
T800	I1, I2, I3, I5
T900	I1, I2, I3

$I1 \wedge I2 \Rightarrow I5,$	$confidence = 2/4 = 50\%$
$I1 \wedge I5 \Rightarrow I2,$	$confidence = 2/2 = 100\%$
$I2 \wedge I5 \Rightarrow I1,$	$confidence = 2/2 = 100\%$
$I1 \Rightarrow I2 \wedge I5,$	$confidence = 2/6 = 33\%$
$I2 \Rightarrow I1 \wedge I5,$	$confidence = 2/7 = 29\%$
$I5 \Rightarrow I1 \wedge I2,$	$confidence = 2/2 = 100\%$



# 1.b. Luật kết hợp và luật dãy sử dụng Web

opinion & misc & travel	-->	on-air	90.26%
news & misc & business & bbs	-->	frontpage	90.24%
living & business & sports & bbs	-->	frontpage	90.00%
news & misc & business & sports	-->	frontpage	89.68%
news & tech & living & business & sports	-->	frontpage	89.00%
news & living & business & bbs	-->	frontpage	88.01%
frontpage & tech & living & business & sports	-->	news	87.87%
frontpage & opinion & living & sports	-->	news	87.81%
frontpage & tech & opinion & living	-->	news	87.60%
frontpage & tech & on-air & business & sports	-->	news	87.59%
news & misc & sports & bbs	-->	frontpage	87.55%
news & tech & on-air & business & sports	-->	frontpage	87.43%
news & living & business & sports	-->	frontpage	87.18%
news & business & sports & bbs	-->	frontpage	86.70%
misc & living & travel	-->	on-air	86.55%
tech & living & sports & bbs	-->	frontpage	86.52%
tech & business & sports & bbs	-->	frontpage	86.40%
news & misc & living & business	-->	frontpage	86.22%
on-air & business & sports & bbs	-->	frontpage	86.22%
news & tech & misc & bbs	-->	frontpage	86.18%
on-air & misc & business & sports	-->	frontpage	86.16%
tech & misc & travel	-->	on-air	86.09%
tech & living & business & sports	-->	frontpage	86.08%
news & living & sports & bbs	-->	frontpage	86.99%
misc & business & sports	-->	frontpage	85.79%
frontpage & tech & opinion & sports	-->	news	85.78%
news & opinion & living & sports	-->	frontpage	85.69%
misc & business & travel	-->	on-air	85.65%
news & tech & misc & business	-->	frontpage	85.63%
misc & business & bbs	-->	frontpage	85.57%
tech & living & sports & bbs	-->	news	85.49%
local & misc & business & sports	-->	frontpage	85.43%
news & opinion & business & bbs	-->	frontpage	85.32%
news & misc & living & sports	-->	frontpage	85.19%
news & on-air & business & sports	-->	frontpage	85.01%

(a)

misc → local	2.07%
frontpage → frontpage → sports	2.02%
local → frontpage	1.83%
on-air → misc → on-air	1.72%
on-air → frontpage	1.69%
on-air → news	1.51%
news → frontpage → news	1.49%
local → news	1.46%
frontpage → frontpage → business	1.35%
news → sports	1.33%
news → bbs	1.23%
health → local	1.16%
misc → frontpage → frontpage	1.16%
on-air → local	1.15%
misc → on-air → misc	1.15%
frontpage → frontpage → living	1.14%
local → frontpage → frontpage	1.13%
health → misc	1.12%
misc → on-air → on-air	1.10%
local → misc → local	1.09%
misc → news	1.06%
news → living	1.06%
on-air → misc → on-air → misc	1.00%

(b)

## • Các loại mẫu điển hình: xu hướng chung của mọi người

- Luật kết hợp
- Luật dãy
- Cây con phổ biến

# 1.c. Nghiên cứu về luật kết hợp



- **Thống kê từ Google Scholar về số bài viết:**
  - Với cụm từ “Association Rule”:
    - Ở tiêu đề: 2.060 bài (khoảng)  
1.000 bài (2006 – nay)
    - Ở mọi nơi: 27.400 bài (khoảng)
  - Với cụm từ “Apriori Algorithm”:
    - Ở tiêu đề: 350 bài (khoảng)  
219 bài (2006 – nay)
    - Ở mọi nơi: 8.820 bài (khoảng)
  - Với cụm từ “Sequential Pattern”:
    - Ở tiêu đề: 590 bài (khoảng)  
270 bài (2006 – nay)
    - Ở mọi nơi: 15.700 bài (khoảng)

# 1.c. Khai phá xu hướng cá nhân



- **Giới thiệu**

- “Cá nhân hóa”: Thông tin cá nhân và tư vấn cá nhân hóa
- Thông tin cá nhân: CSDL quản lý; Máy khách
- Ngữ cảnh làm việc của cá nhân

- **Một số hình thức**

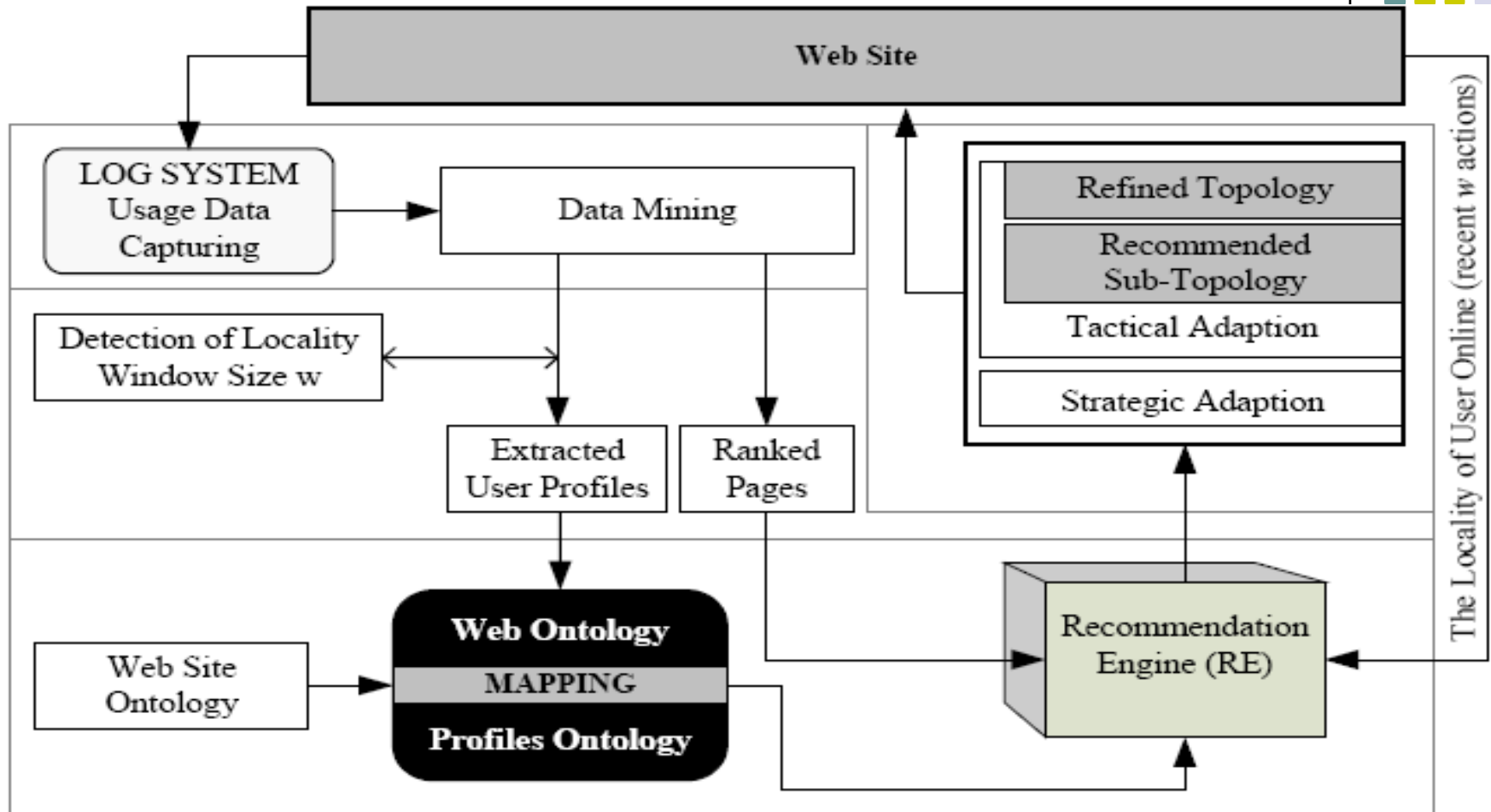
- Khai phá xu hướng cá nhân từ thông tin máy khách
- Hệ tư vấn

- **Hệ tư vấn**

- Recommendation Systems
- Lọc cộng tác, lọc nội dung, lọc kết hợp
- Hội thảo dành riêng: các năm 2007, 2009, 2010
- <http://recsys.acm.org/2010/>

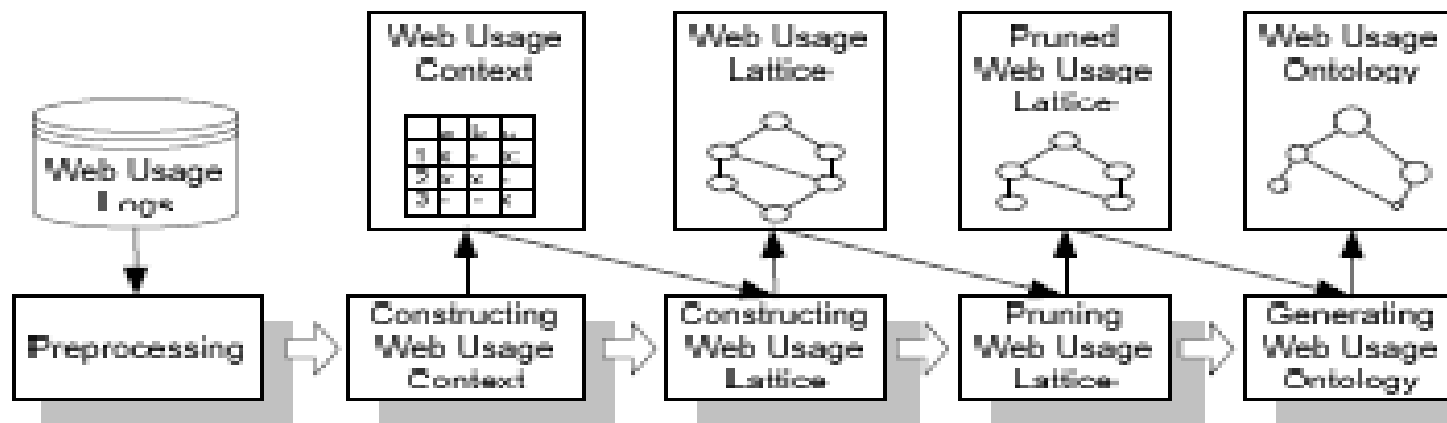
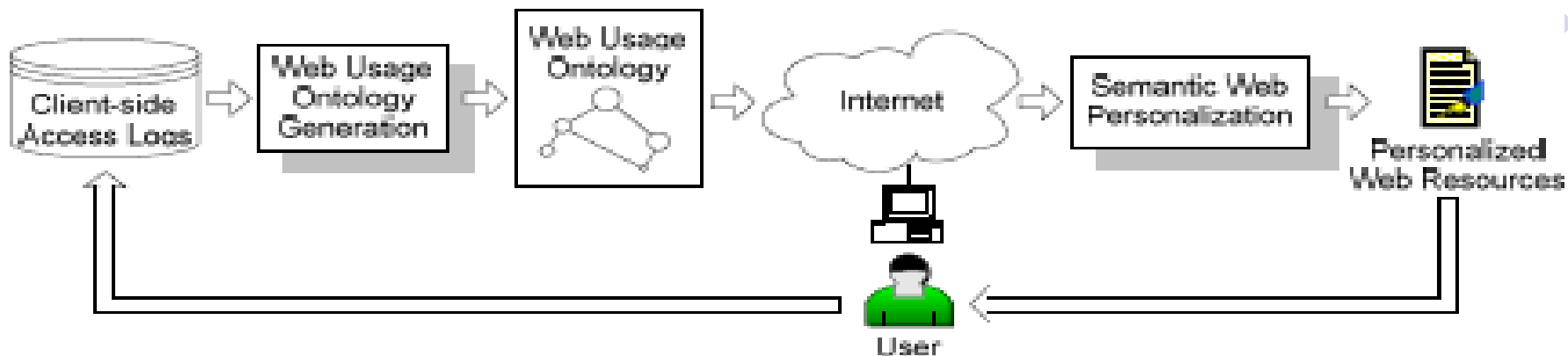


# 1.c. Sinh tự vận dựa theo tiêu sử người dùng



[RK07] Tarmo Robal, Ahto Kalja (2007). Applying User Profile Ontology for Mining Web Site Adaptation Recommendations, *ADBIS Research Communications 2007*

# 1.c. Khai phá sử dụng Web

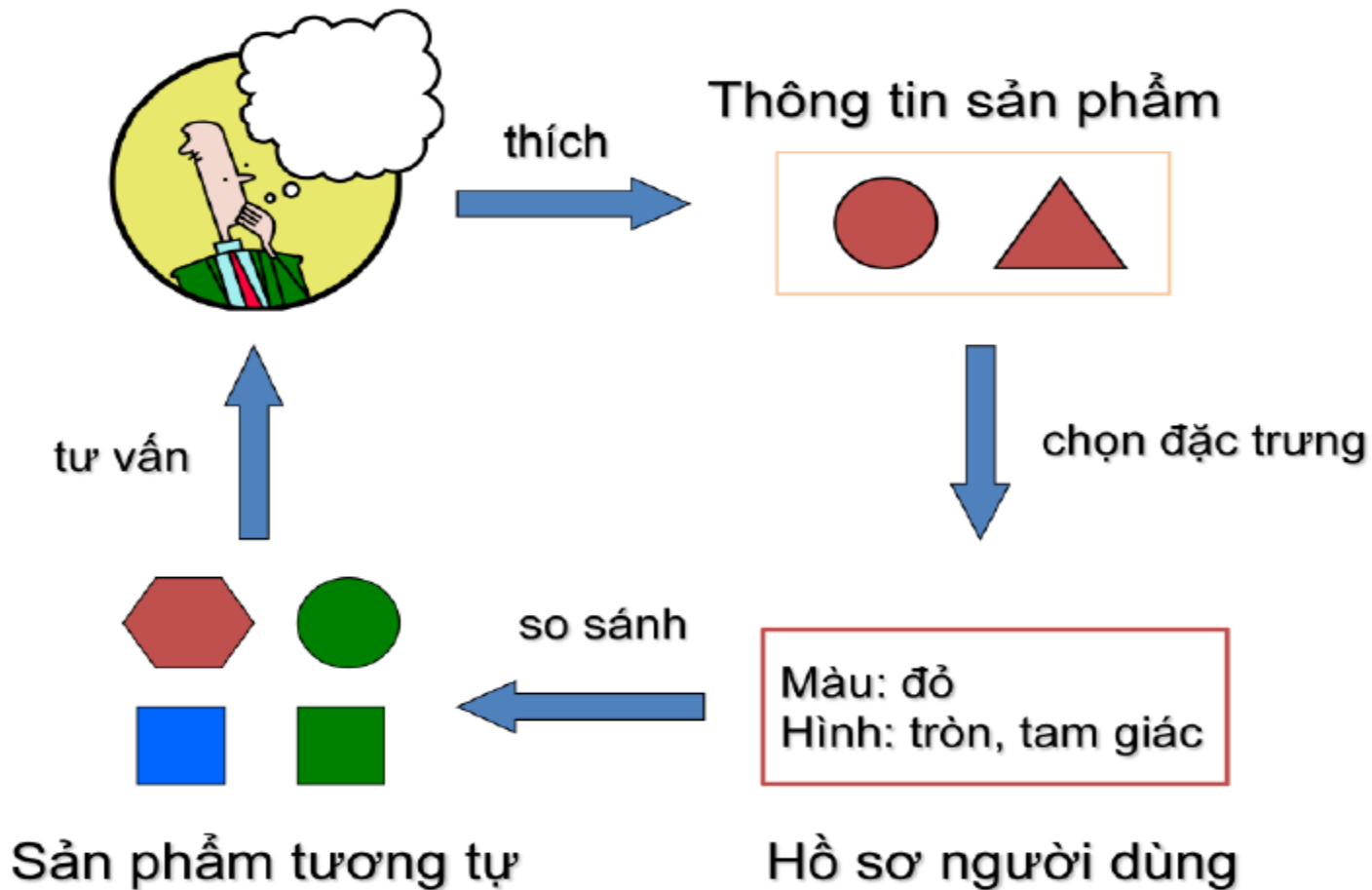


## Hệ thống khai phá sử dụng Web tự vấn hướng cá nhân

- Kiến trúc hệ thống (trên)
- và sinh ontology sử dụng Web (dưới)

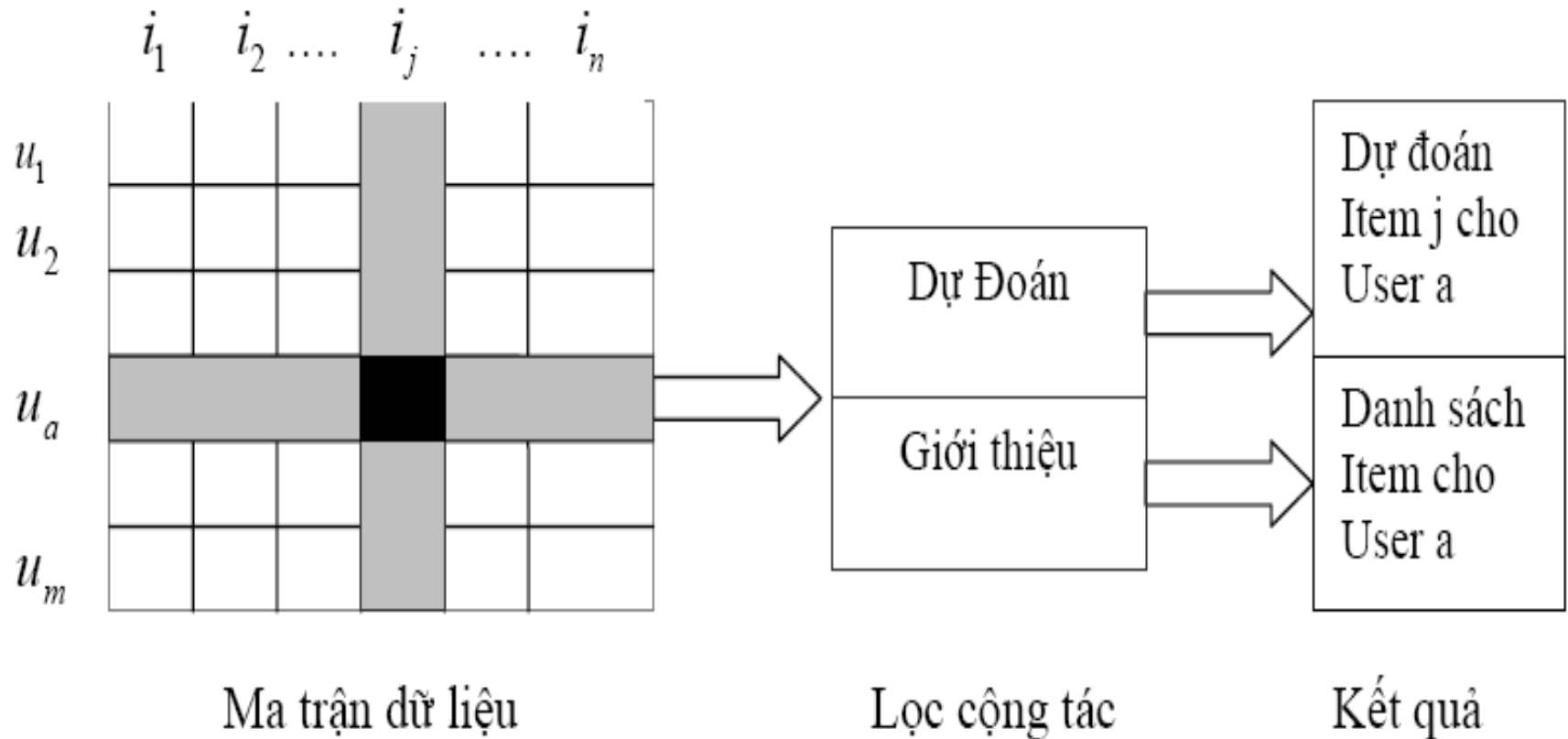
Baoyao Zhou, Siu Cheung Hui, Alvis C. M. Fong (2005). Web Usage Mining for Semantic Web Personalization, *Workshop on Personalization on the Semantic Web*, 66–72, Edinburgh, UK, 2005.

# 1.c. Hệ thống tư vấn: lọc nội dung



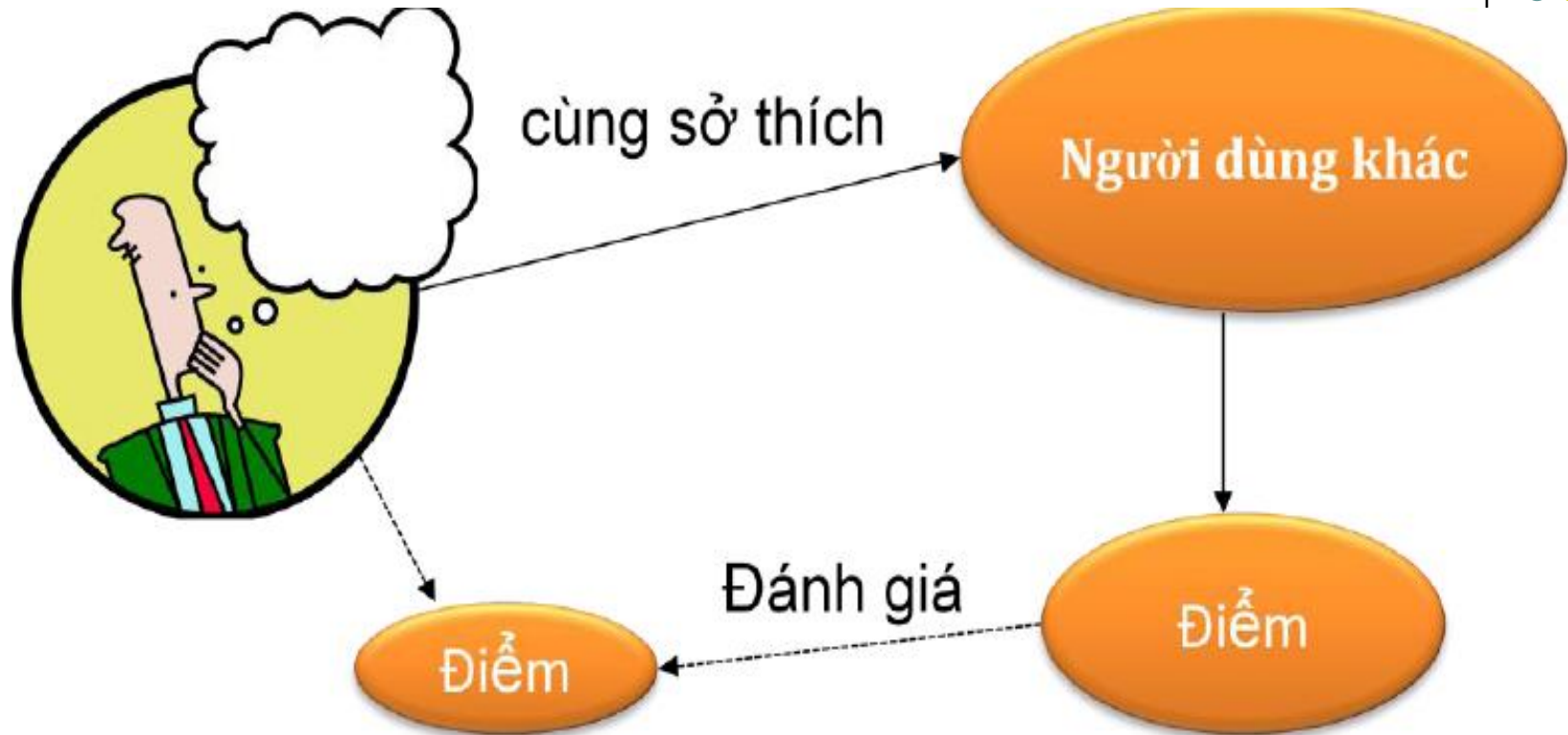
Lấy nội dung thuộc tính các sản phẩm người dùng đã ưa thích để dự đoán sản phẩm ưa thích tiếp theo

# 1.c. Hệ thống tự vấn: lọc cộng tác

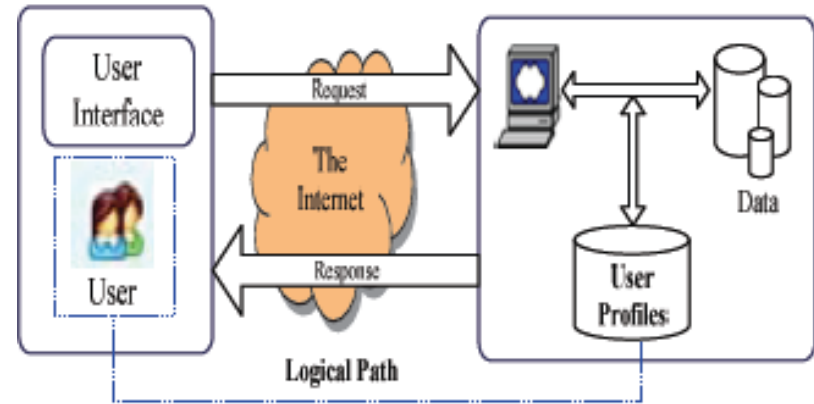
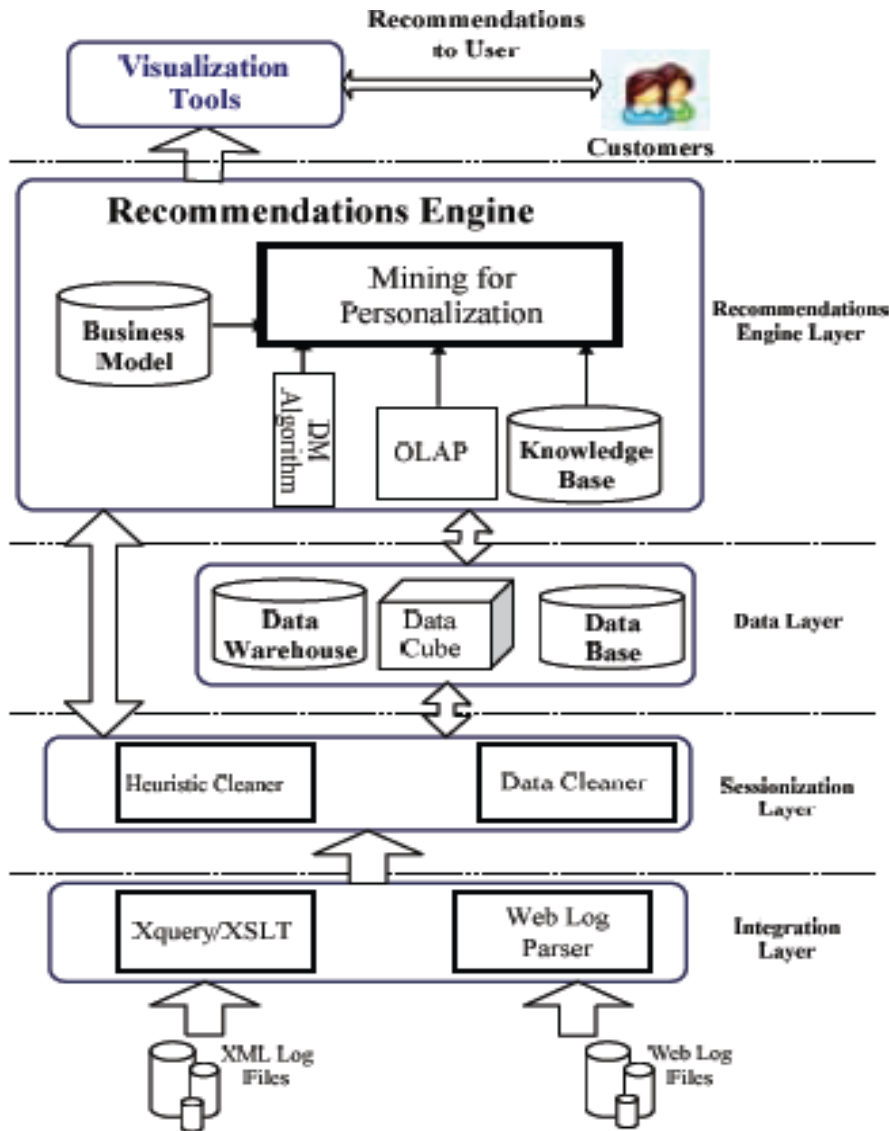


Quan hệ người dùng – sản phẩm: nhóm người dùng “tương tự nhau” và khi có một người dùng trong “thích” thì các người khác cũng “thích” tương tự

# 1.c. Hệ thống tư vấn: lọc cộng tác



# 1.c. Hệ thống tư vấn: lọc cộng tác



Jinhua Sun, Yanqi Xie (2009). A Web Data Mining Framework for E-commerce Recommender Systems, *Computational Intelligence and Software Engineering, 2009. CiSE 2009*.

# Nghiên cứu về khai thác sử dụng Web



- **Thống kê từ Google Scholar về số bài viết:**
  - Với cụm từ “Web Usage Mining”:
    - Ở tiêu đề: 860 bài (khoảng)  
280 bài (2006 – nay)
    - Ở mọi nơi: 171.000 bài (khoảng)
  - Với cụm từ “Web Log Mining”:
    - Ở tiêu đề: 340 bài (khoảng)  
140 bài (2006 – nay)
    - Ở mọi nơi: 137.000 bài (khoảng)
  - Với cụm từ “Recommendation System”:
    - Ở tiêu đề: 1.750 bài (khoảng)  
750 bài (2006 – nay)
    - Ở mọi nơi: 1.760.000 bài (khoảng)

## 2. Khai phá cấu trúc Web

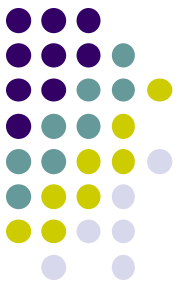


- Hai bài toán điển hình
  - Khai phá liên kết Web
  - Khai phá cấu trúc trang Web
- Khai phá liên kết Web
  - Mỗi trang Web là một đỉnh
  - Liên kết các trang Web hình thành các cung
  - Đồ thị có hướng hoặc vô hướng
  - Web phản ánh xã hội: đồ thị Web là một loại mạng xã hội
  - Hạng trang Web, một bài toán điển hình: tính “độ quan trọng” của một trang Web (một nút trên đồ thị Web)
  - Khai phá liên kết Web: Phân lớp trang web dựa theo liên kết, Phân tích cụm dựa theo liên kết, Kiểu liên kết; Độ mạnh liên kết;



## 2. Khai phá liên kết Web

- Phân lớp Web dựa theo liên kết
  - Khai thác thông tin liên kết cho phân lớp Web
- Phân cụm Web dựa theo liên kết
  - Tìm ra sự xuất hiện tự nhiên các lớp con: dữ liệu là liên kết
- Phân tích kiểu liên kết
  - Dự báo về sự tồn tại của liên kết
  - Dự báo mục đích của liên kết
- Phân tích độ mạnh liên kết
  - Độ mạnh của cung và đỉnh (hạng trang)
- Phân tích số lượng liên kết
  - Dự báo số lượng liên kết giữa các đối tượng



## 2. Khai phá cấu trúc trang Web



- **Cấu trúc trang Web**

- Trang Web được viết theo ngôn ngữ trình bày Web: chẳng hạn HTML, XML
- Trang web được tổ chức dưới dạng hình cây
- Cấu trúc trình bày nội dung trang web

- **Phân tích cấu trúc trang Web**

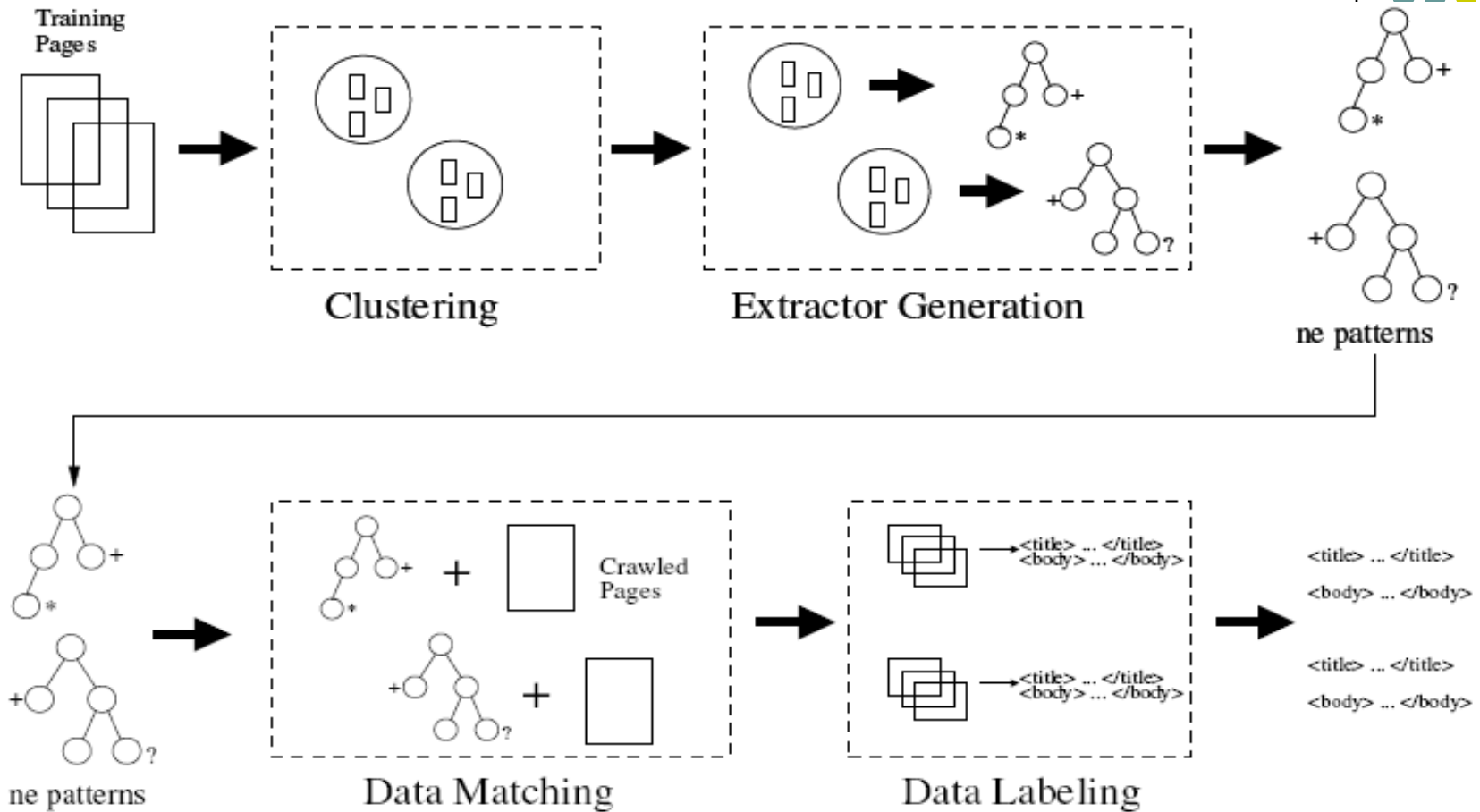
- Tìm các mẫu cấu trúc trang Web
- Kết hợp với khai phá nội dung Web

## 2. Khai phá cấu trúc trang báo điện tử



Davi de Castro Reis, Paulo B. Golgher, Altigran S. da Silva, Alberto H. F. Laender (2004). Automatic Web News Extraction Using Tree Edit Distance, *Proceedings of the Thirteenth International World Wide Web Conference*: 502-601, ACM Press, New York, NY, May 2004, ISBN 1581139128

## 2. Khai phá cấu trúc trang báo điện tử



Davi de Castro Reis, Paulo B. Golgher, Altigran S. da Silva, Alberto H. F. Laender (2004). Automatic Web News Extraction Using Tree Edit Distance, *Proceedings of the Thirteenth International World Wide Web Conference*: 502-601, ACM Press, New York, NY, May 2004, ISBN 1581139128

## 2. Áp dụng: báo điện tử Việt Nam



**NGHIÊN CỨU CÔNG NGHỆ KHAI PHÁ DỮ LIỆU VĂN BẢN, ÁP DỤNG CHO CÁC TRANG TIN TỨC TRÊN CÁC THIẾT BỊ CẦM TAY (PDAs & SMARTPHONES)**



*Sinh viên thực hiện:* **KS. Vũ Ngọc Anh**

*Giáo viên hướng dẫn:* **TS. Hà Quang Thụy**

*Lớp:* **K9T3**



## 2. Áp dụng: báo điện tử Việt Nam



**MỘT SỐ HÌNH ẢNH THỰC NGHIỆM**

Wednesday, November 10, 2010

Kênh tin tức điện tử cho PDAs & Smartphones

Trang 29

Vũ Ngọc Anh (2006). Kênh tin tức điện tử cho PDAs & Smartp, *Luận văn Thạc sỹ*, Trường ĐHCN-ĐHQGHN



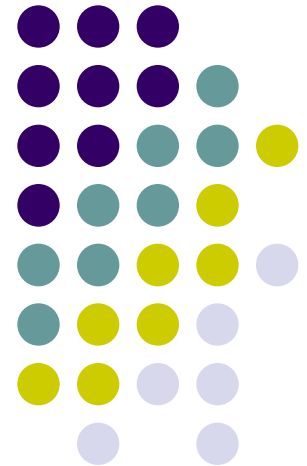
## 2. Áp dụng: báo điện tử Việt Nam

- <http://vietbao.vn/Vi-tinh-Vien-thong/12-san-pham-vao-vong-chung-khao-Tri-tue-Viet-Nam/20641855/217/> ; Thứ sáu, 08 Tháng mười hai 2006, 02:31 GMT+7
  - “4. Vienews - kênh báo điện tử trên thiết bị điện thoại di động thông minh (Vũ Ngọc Anh, Hà Duyên Hòa - Hà Nội): Sản phẩm hỗ trợ thiết bị di động cầm tay đọc báo điện tử qua môi trường Internet không dây”.
- <http://www.tapchibcv.gov.vn/vi-vn/dacsan/2006/8/17521.bcv> ; 7:58, 02/01/2007
  - 7. **Giải Ba:** Sản phẩm đoạt giải: “**Các kênh báo điện tử trên thiết bị điện thoại di động thông minh**” của Hà Duyên Hoá (Hà Nội).

# BÀI GIẢNG KHAI PHÁ DỮ LIỆU WEB

## CHƯƠNG 3. MỘT SỐ KIẾN THỨC TOÁN HỌC BỔ TRỢ CHƯƠNG 4. MỘT SỐ BÀI TOÁN XỬ LÝ NGÔN NGỮ TỰ NHIÊN NỀN TẢNG

PGS. TS. HÀ QUANG THỤY  
HÀ NỘI 10-2010  
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ  
ĐẠI HỌC QUỐC GIA HÀ NỘI



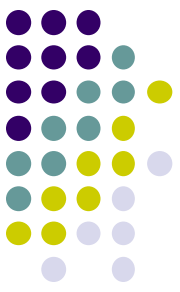


# Nội dung

1. Một số kiến thức Toán học bổ trợ
2. Một số bài toán xử lý ngôn ngữ tự nhiên nền tảng



# C3. Một số kiến thức Toán học hỗ trợ



- **Toán học Internet**

- Ra đời một lĩnh vực mới: Internet Mathematics
- Cộng đồng Toán học Internet: Internet Mathematics Community

- **Đối tượng và các chủ đề**

- Đối tượng: Mạng phức tạp trên Internet và Web: đồ thị Web, đồ thị Internet, mạng xã hội trực tuyến (Facebook, LinkedIn, và Twitter...), mạng sinh học trên Web...
- Các chủ đề thuộc khai phá và mô hình hóa web (cơ sở lý thuyết và ứng dụng thực tiễn) trong môi trường mạng phức tạp.

- **Tạp chí Internet Mathematics**

- <http://www.internetmathematics.org/> (2/2011 - xem trang sau)
- Đồng Trưởng ban biên tập:
  - Fan Chung Graham (<http://www.math.ucsd.edu/~fan/>). DBLP: 137 bài báo
  - Anthony Bonato (<http://www.math.ryerson.ca/~abonato/>). DBLP: 35 bài báo
- Công bố bài báo chất lượng cao về mạng phức

# Tạp chí Internet Mathematics



Internet Mathematics: Editorial Board

## Internet Mathematics

<http://www.internetmathematics.org/board.htm>

Statement of  
Philosophy

Subscription  
Information

Submission  
Guidelines

Articles

Editorial Board

### Editorial Board

#### Editors-in-Chief

Fan Chung Graham  
Anthony Bonato

#### Managing Editors

Xiaotie Deng  
Nelly Litvak

#### Editorial Board

Noga Alon  
Albert-László Barabási  
Elwyn Berlekamp  
Béla Bollobás  
Andrei Broder  
Jennifer Chayes  
Persi Diaconis  
Ding-Zhu Du  
Rick Durrett  
Cynthia Dwork  
Alan Frieze  
Tim Griffin  
Ronald Graham  
Monika Henzinger

Frank Kelly  
Jon Kleinberg  
Tom Leighton  
Michael Mitzenmacher  
S. Muthu Muthukrishnan  
Andrew Odlyzko  
Christos Papadimitriou  
Prabhakar Raghavan  
Peter Sarnak  
Joel Spencer  
Walter Willinger  
Peter Winkler  
Andrew Yao

To order a  
subscription, or to  
request further  
information or a  
sample issue, [send  
e-mail to us](#) or  
contact the publisher  
at:

A K Peters  
5 Commonwealth Rd.  
Suite 2C  
Natick, MA 01760-1526  
phone: 508-651-0887  
fax: 508-651-0889



Copyright ©2010  
A.K Peters, Ltd.  
All rights reserved.



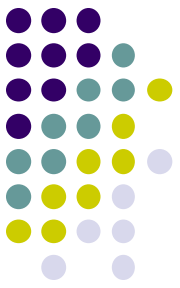
- **Ban biên tập tạp chí: Bổ sung một số chuyên gia**  
Jennifer Tour Chayes <http://research.microsoft.com/en-us/um/people/jchayes/>. “She is the co-author of over 100 scientific papers and the co-inventor of more than 25 patents”  
Rick Durrett <http://www.math.duke.edu/~rtd/> .  
Andrew Tomkins <http://www.tomkinshome.com/andrew/paperlist>. DBLP: 88 bài báo
- **Một số biên tập viên được lưu ý**  
Ronald L. Graham (<http://www.math.ucsd.edu/~ronspubs/>). DBLP:116 bài báo. Nhiều giải thưởng  
[Frank Kelly](http://www.statslab.cam.ac.uk/~frank/) (<http://www.statslab.cam.ac.uk/~frank/> )

# Một số nội dung Toán học bổ trợ



- **Mô hình đồ thị**
  - Một số kiến thức cơ sở
  - Đồ thị ngẫu nhiên
  - Mạng xã hội
- **Học máy xác suất Bayes**
  - Một số kiến thức cơ sở
  - Học máy xác suất Bayes
  - Ước lượng giá trị tham số
- **Thuật toán Viterbi**
  - Lý thuyết quyết định hỗn hợp
  - Nội dung thuật toán

# Đồ thị Web và đồ thị ngẫu nhiên



- **Đồ thị Web**

- Web có cấu trúc đồ thị
  - Đồ thị Web: nút  $\leftrightarrow$  trang Web, liên kết ngoài  $\leftrightarrow$  cung (có hướng, vô hướng).
  - Bản thân trang Web cũng có tính cấu trúc cây (đồ thị)
- Một vài bài toán đồ thị Web
  - Biểu diễn nội dung, cấu trúc
  - Tính hạng các đối tượng trong đồ thị Web: tính hạng trang, tính hạng cung..

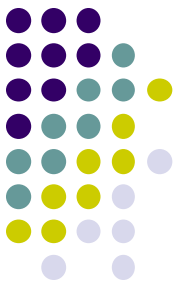
Nghiên cứu về đồ thị Web (xem trang sau)

- **Đồ thị ngẫu nhiên**

- Tính ngẫu nhiên trong khai phá Web
  - WWW có tính ngẫu nhiên: mới, chỉnh sửa, loại bỏ
  - Hoạt động con người trên Web cũng có tính ngẫu nhiên
- Là nội dung nghiên cứu thời sự

# Bibliography Webgraph Papers

Dragomir R. Radev, 03/4/2010



Toàn bộ	2007	2008	2009	To 04/10	2007-10
1542	127	61	36	13	237

- So many webgraph research papers.
- Some previous versions of “Bibliography Webgraph Papers” by Dragomir R. Radev
- 1601: <http://clair.si.umich.edu/~radev/webgraph/webgraph-bib.html>

5/2005	5/2007	5/2008	1/2009	8/2009	4/2010	11/2010
496	1212	1361	1457	1471	1542	1601

# Lý thuyết về đồ thị lớn



## Đồ thị lớn

- Số đỉnh lên tới hàng tỷ
- Biểu diễn cung chính xác không còn là quan trọng

## Cơ sở lý thuyết trong nghiên cứu đồ thị lớn

- Khả năng là lý thuyết sinh đồ thị
- Bất biến tới một số thay đổi nhỏ trong định nghĩa
- Phải có năng lực chứng minh các định lý cơ bản

[Hop07] John E. Hopcroft (2007). Future Directions in Computer Science, <http://www.cs.cornell.edu/jeh/China%202007.ppt>

# Đồ thị ngẫu nhiên: Mô hình Erdős-Rényi

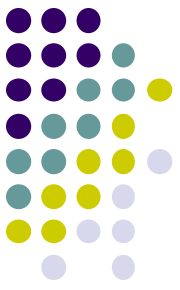


- Đồ thị ngẫu nhiên: có thể mô hình mạng thế giới thực.
- Định nghĩa: có hai định nghĩa
  - Chọn ngẫu nhiên:  $G_{n, N}$  được chọn ngẫu nhiên từ  $\mathcal{G}_{n, N} = \{\text{mọi đồ thị có } n \text{ đỉnh và } N \text{ cung}\}$  các phần tử trong  $\mathcal{G}_{n, N}$  là đồng khả năng được chọn với xác suất  $1/\binom{n}{2}/N$ ;
  - Quá trình hình thành các cung trong  $G_{n, N}$  là ngẫu nhiên: mỗi cạnh xuất hiện với xác suất  $p$ , sự xuất hiện hay vắng mặt hai cạnh là độc lập nhau.

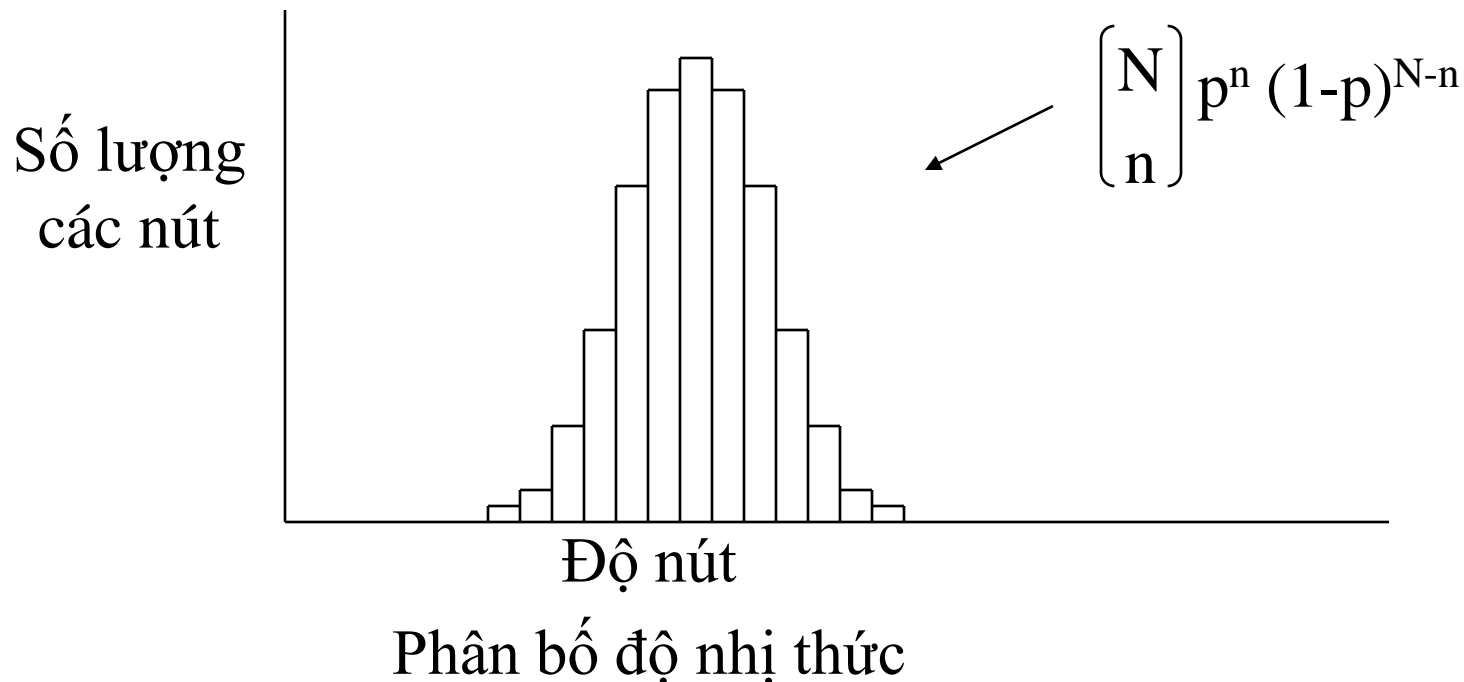
[ER61] P. Erdős, A. Rényi (1961). On the evolution of random graphs, *Théorie de L'Information*: 343-347, 1961.



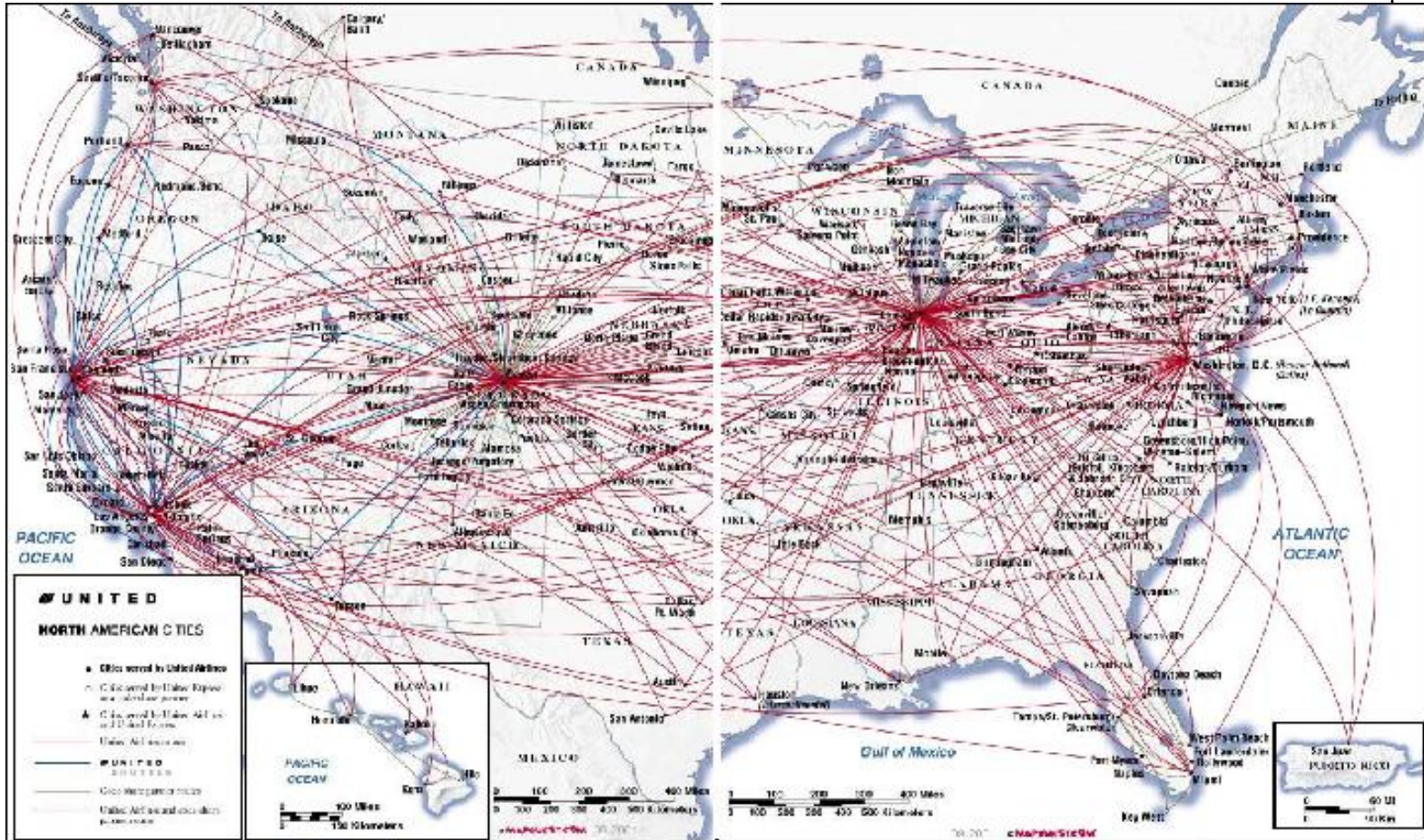
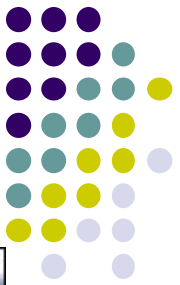
# Đồ thị ngẫu nhiên: Mô hình Erdős-Rényi



- Đặt tên: Paul Erdős và Alfréd Rényi
- Là một trong hai mô hình sinh các đồ thị ngẫu nhiên
- Chứa tập các nút mà mỗi nút trong mỗi tập đó có xác suất như nhau, độc lập với các cung khác
- $n$  nút: Mỗi bộ  $n^2$  cung tiềm năng được biểu diễn với xác suất độc lập



# Đồ thị ngẫu nhiên

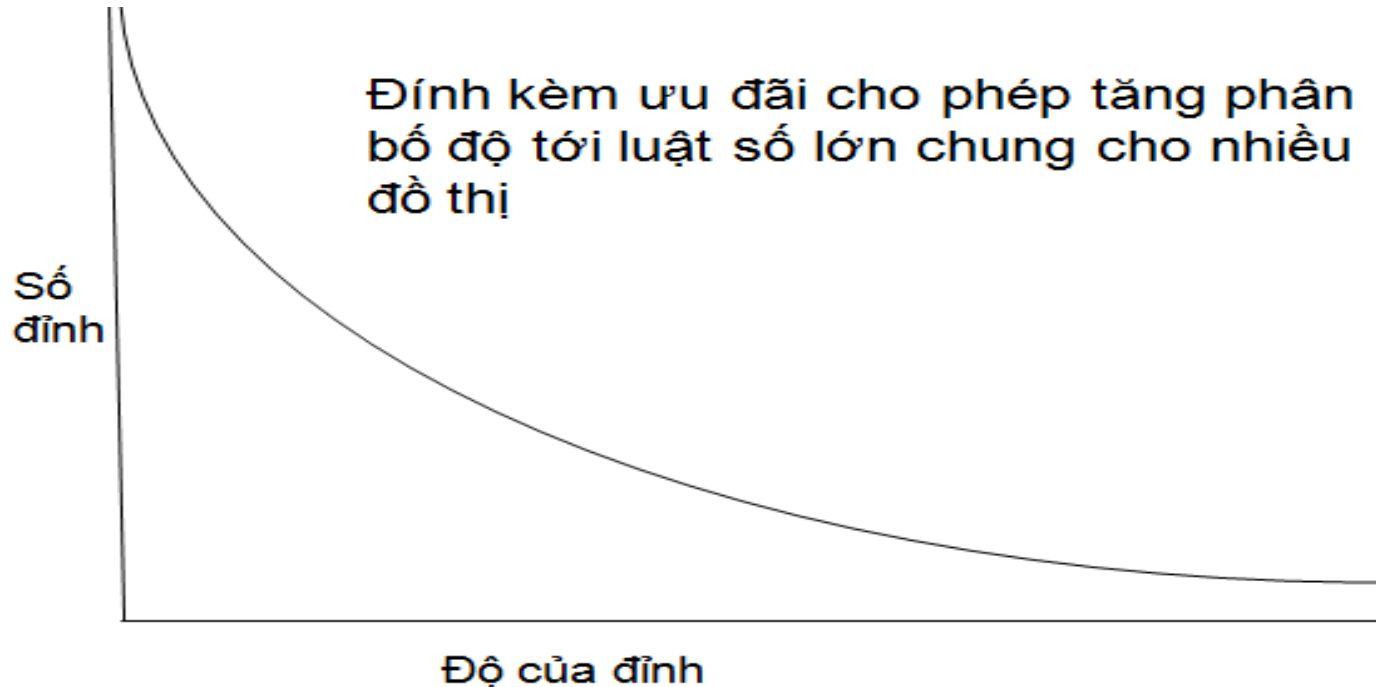


[Hop07] John E. Hopcroft (2007). Future Directions in Computer Science, <http://www.cs.cornell.edu/jeh/China%202007.ppt>

# Mô hình sinh đồ thị



- Các nút và cung được bổ sung sau mỗi đơn vị thời gian
- Quy tắc xác định nơi cung xuất hiện (nơi đặt cung mới)
  - Xác suất đồng nhất
  - Đỉnh kèm ưu đãi – đưa đến phân bố theo luật số lớn



[Hop07] John E. Hopcroft (2007). Future Directions in Computer Science, <http://www.cs.cornell.edu/jeh/China%202007.ppt>

# Mạng xã hội

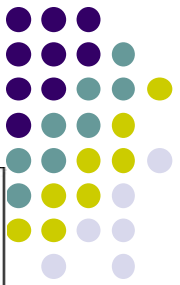


- Mạng xã hội

- Internet, Web là một xã hội ảo
  - Nhiều hoạt động (đặc biệt là hoạt động thông tin) trong thế giới thực được thi hành
  - “Thế giới phẳng”, “toàn cầu hóa” và “bản địa hóa”
- Khái niệm
  - ❖ Mạng xã hội là mạng của một nhóm người có hoạt động và các mối quan hệ gắn kết họ với nhau.
  - ❖ Mạng xã hội là một kiểu của mạng phức tạp
- Một số ví dụ mạng xã hội trên Internet
  - ❖ Diễn đàn, Blog, Mạng e-mail, mạng xã hội chuyên đề
  - ❖ Một số ví dụ khác (trang bên)
- Nghiên cứu mạng xã hội
  - ❖ Vấn đề nghiên cứu thời sự.
  - ❖ Kết hợp nhiều lĩnh vực, chẳng hạn như CNTT + Xã hội học



# Mạng xã hội: ví dụ



**Events and News**  
Duncan J. Watts's new book is out now!

**Project Information**  
In the Press  
Description  
Procedures  
Security and Privacy  
Articles/References  
Results

**Research Team**  
Duncan J. Watts  
Peter Dodds  
Roby Muhamad

**Web Development**  
Peter Hausel

Vijay (Delhi, India) worked at an engineering firm with

Sameer (Kolkata, India) whose daughter

Prema (Berkeley, USA) goes to school in California and plays soccer with

Christie (Berkeley, USA) whose best friend from high school

Alice (New York, USA)

William (New York, NY) is studying medicine with

The **SMALL WORLD** project is an online experiment to test the idea that any two people in the world can be connected via 'six degrees of separation'.

Your objective is to get a message to a "target person", somewhere in the world, by forwarding the message to a friend of yours--someone who is "closer" to the target than you are. (If you happen know the target, you can of course send it to them)

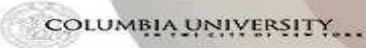
If we have asked you to participate (you would have received a message from a friend of yours), you should **continue** the chain.

If you are just visiting us, sign up to start a new chain.

home  
my small world  
chat  
FAQ  
related links

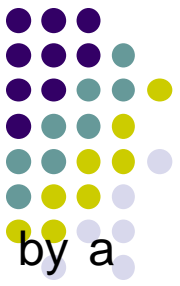
login

sign up

 COLUMBIA UNIVERSITY  
THE CITY OF NEW YORK

<http://www.uvm.edu/~pdodds/teaching/courses/2008-01UVM-295/docs/2008-01UVM-295smallworldnetworks-slides-handout.pdf>

# Social Networks: Properties



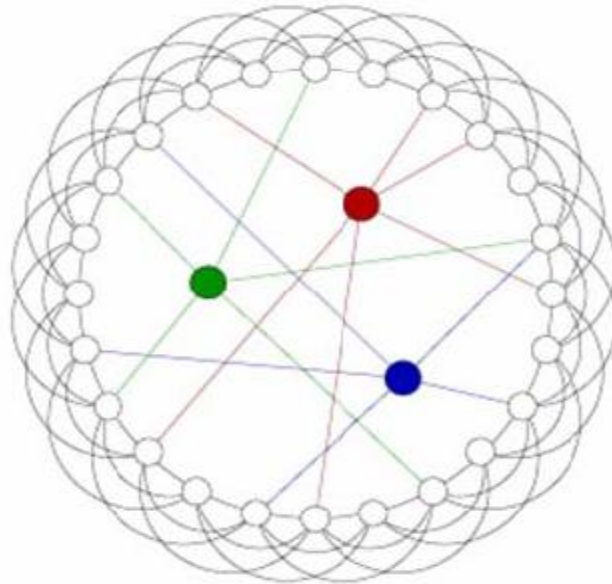
- The small-world property
    - Almost any pair of people in the world can be connected together by a short chain of intermediate acquaintances, usually about six lengths.
- [TM69] Jeffrey Travers, Stanley Milgram (1969). An Experimental Study of the Small World Problem, *Sociometry*, **32**(4): 425-443, Dec., 1969.
- Power-law degree distributions / the scale – free property
    - Social network's nodes (also edges) are distributed under the power-law degree
  - Network transitivity
    - Structure and dynamics of the network influenced by nodes with the large number of connectings (using to detect communities in a social network!)
  - **Community structure**
    - Networks are divided into communities in which the nodes in the same community closed links, and links communities liquid
    - A community in social networks as an “**interest group**” in the real world. [http://en.wikipedia.org/wiki/Interest\\_group\\_\(disambiguation\)](http://en.wikipedia.org/wiki/Interest_group_(disambiguation)) as meaning of “**nhóm lợi ích**” in Vietnamese. See also “**Advocacy group**”, “**Lobby group**”. 5P&5C marketing model: **People** ⇨ **Customer approach** (Product ⇨ Consumer desire; Price ⇨ Cost; Place ⇨ Convenience; Promotion ⇨ Communication)
    - Flexible community structure: one community structure for one case.

# Social Networks: Properties



## Small-world Networks

Almost any pair of people in the world can be connected to one another by a short chain of intermediate acquaintances, of typical length about six.



Bui, N., Lan; Tran, Q., Anh; Ha, Q., Thuy

User's authentic rating based on email networks

Lan N. Bui, Anh Q. Tran, Thuy Q. Ha (2006). User authentic Rating based on Email Networks, *ICMOCCA2006*: 144-148, Seoul, Korea & *International Journal of Natural Sciences and Technology*, 1(2): 173-180, 2006.



### The quantitative definition of the clustering coefficient

$$C_i = \frac{2E_i}{k_i(k_i - 1)} \quad (1)$$

- $C_i$  is clustering coefficient of node  $i$  in email networks.
- $k_i$  is degree of node  $i$  (or node  $i$  has  $k_i$  neighbors)
- $E_i$  is the number of links between  $k_i$  neighbors of  $i$ .

### Comments

- Can't calculate in case  $K_i = 1$ .
- If  $E_i = 0$  then  $C_i = 0$ , independent of  $k_i$





### Clustering coefficient

$$C_i = \frac{2(E_i + 1)}{k_i(k_i - 1) + 1} \quad (2)$$

### Comments

- Formula (2) can't discriminate between users who sent e-mails from users who received e-mails.
- For this reason, we consider the email network graph as directed graph.



## Proposed method

### New clustering coefficient formula

$$C_i = a \times \frac{2(E_i + 1)}{S_i(S_i - 1) + 1} + b \times R_i \quad (3)$$

### Formulas of $E_i, S_i, R_i$

$$E_i = \sum_{j=1}^{Edge} (1 + w_i \times 0.05) \quad (4)$$

$$S_i = \sum_{j=1}^{Send} (1 + w_i \times 0.05) \quad (5)$$

$$R_i = \sum_{j=1}^{Recieve} (1 + w_i \times 0.05) \quad (6)$$



# E-mail Networks



## Experimental Data

- The e-mail network studied here is constructed from log files of the VNUH e-mail server.
- Logs over of a period of one week (from March, 28th to April, 4th, 2006.)
- Consists of **19876** users (1149 in-users, 18727 out-users).
- Total of exchanged messages in this time is **88842** messages



# E-mail Networks



## Email Network Graph

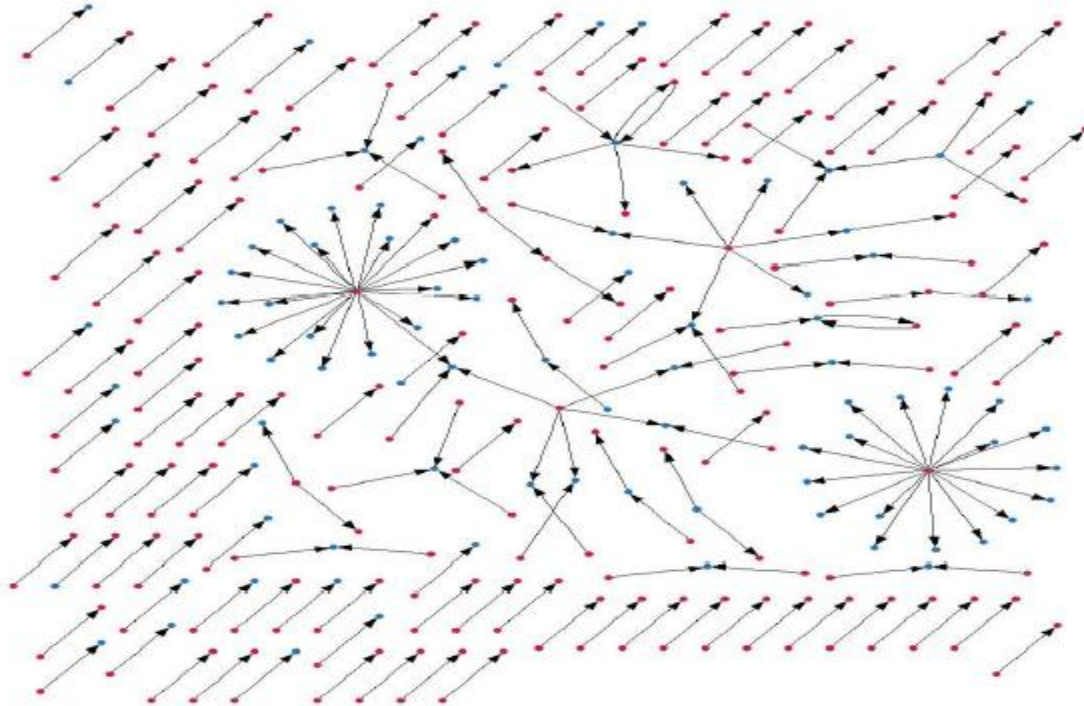


Figure: Email network graph is constructed from log files of VNUH in an hour



Bui, N., Lan; Tran, Q., Anh; Ha, Q., Thuy

User's authentic rating based on email networks

Lan N. Bui, Anh Q. Tran, Thuy Q. Ha (2006). User authentic Rating based on Email Networks, *ICMOCCA2006*: 144-148, Seoul, Korea & *International Journal of Natural Sciences and Technology*, 21(2): 173-180, 2006.

# E-mail Networks

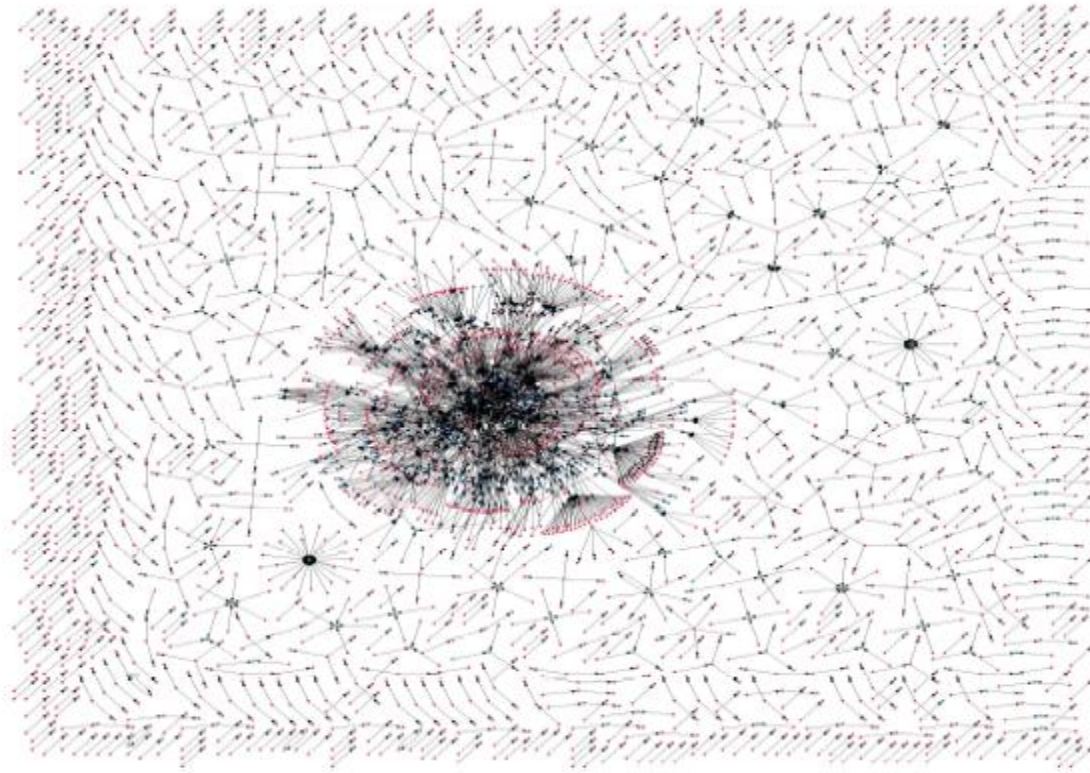


Figure: Email network graph is constructed from log files of VNUH in a week



Bui, N., Lan; Tran, Q., Anh; Ha, Q., Thuy

User's authentic rating based on email networks

Lan N. Bui, Anh Q. Tran, Thuy Q. Ha (2006). User authentic Rating based on Email Networks, *ICMOCCA2006*: 144-148, Seoul, Korea & *International Journal of Natural Sciences and Technology*, 1(2): 173-180, 2006.

# E-mail Networks



## Result

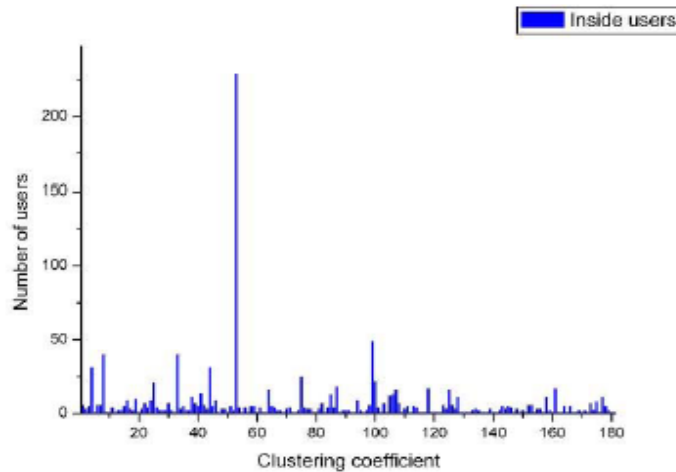


Figure: Clustering coefficient distribution diagram of in-users

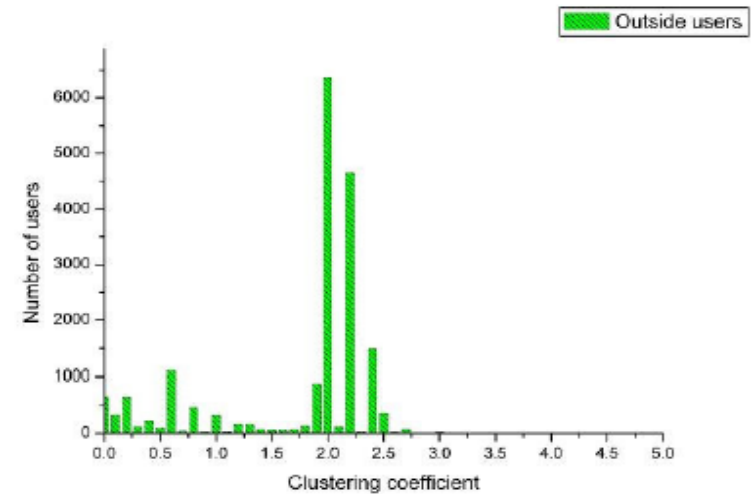
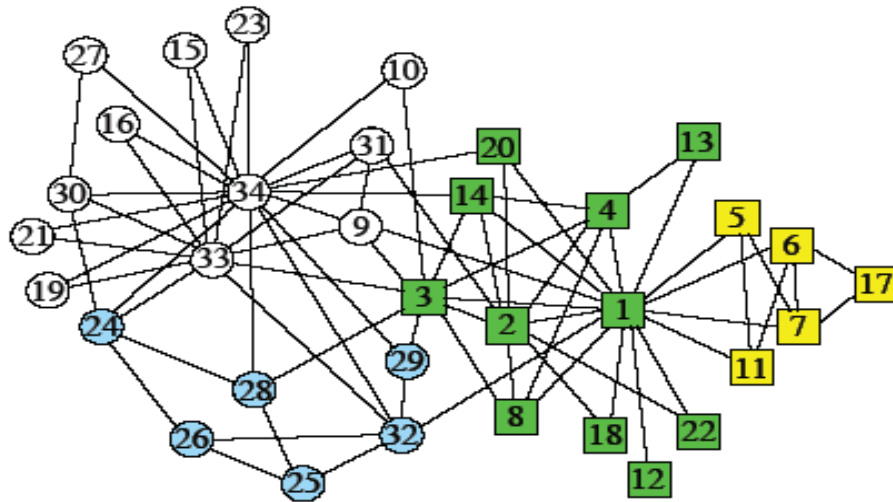


Figure: Clustering coefficient distribution diagram of out-users





# Mạng XH và cộng đồng [For10]

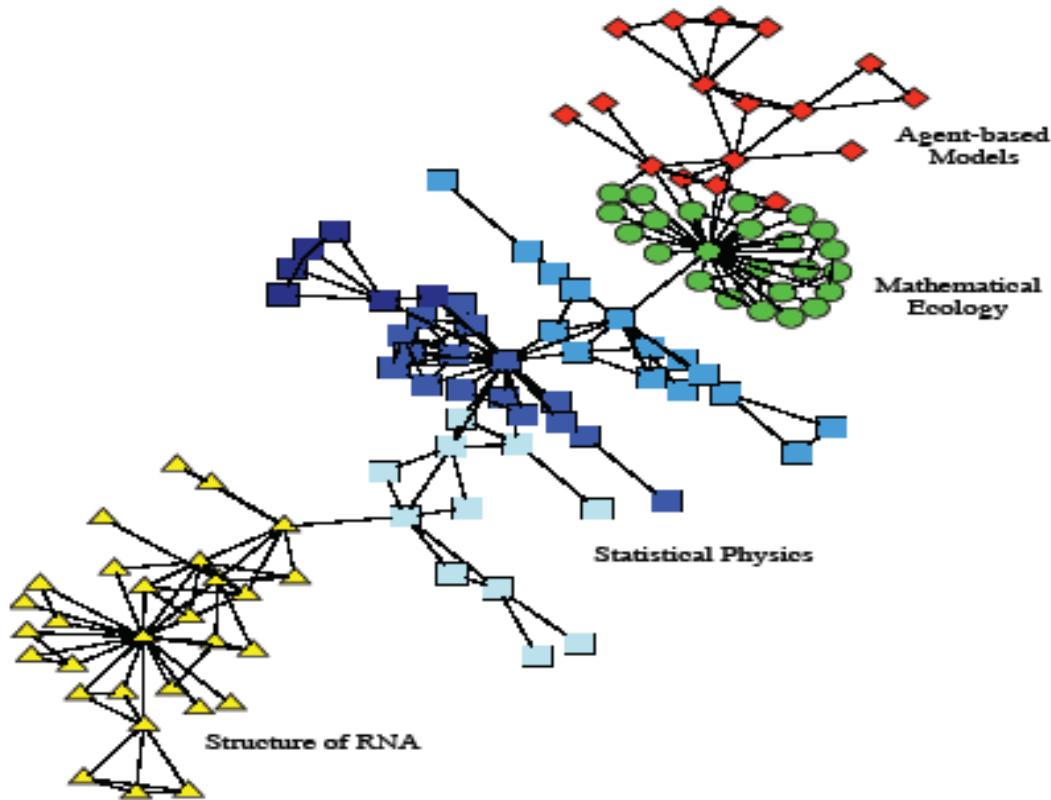


- Câu lạc bộ karate của Zachary (được quan sát trong 3 năm), một kiểm chứng chuẩn cho phát hiện cộng đồng. Các màu sắc tương ứng với phân hoạch tốt nhất tìm được bằng cách tối ưu các mô đun của Newman và Girvan.

- Đồ thị gồm 34 đỉnh thành viên của câu lạc bộ. Bên ngoài các hoạt động của câu lạc bộ. Theo quan sát, có xung đột giữa chủ tịch câu lạc bộ và người hướng dẫn dẫn đến sự phân hoạch câu lạc bộ thành hai nhóm riêng biệt, tương ứng ủng hộ người hướng dẫn và chủ tịch (chỉ dẫn hình vuông và hình tròn). Câu hỏi đặt ra là liệu từ cấu trúc mạng ban đầu có thể suy luận các thành phần của hai nhóm.
- Nhìn vào hình, có thể phân biệt hai tập hợp, một tập quanh các đỉnh 33 và 34 (34 là chủ tịch), tập còn lại quanh đỉnh 1 (người hướng dẫn).
- Cũng có một số đỉnh nằm giữa hai cấu trúc chính, chẳng hạn như 3, 9, 10; đỉnh như vậy thường không phân loại được theo phương thức phát hiện cộng đồng.

[For10] Santo Fortunato (2010), Community detection in graphs, *Technical Report, Complex Networks and Systems Lagrange Laboratory, ISI Foundation, Torino, ITALY.*

# Mạng XH và cộng đồng [For10]

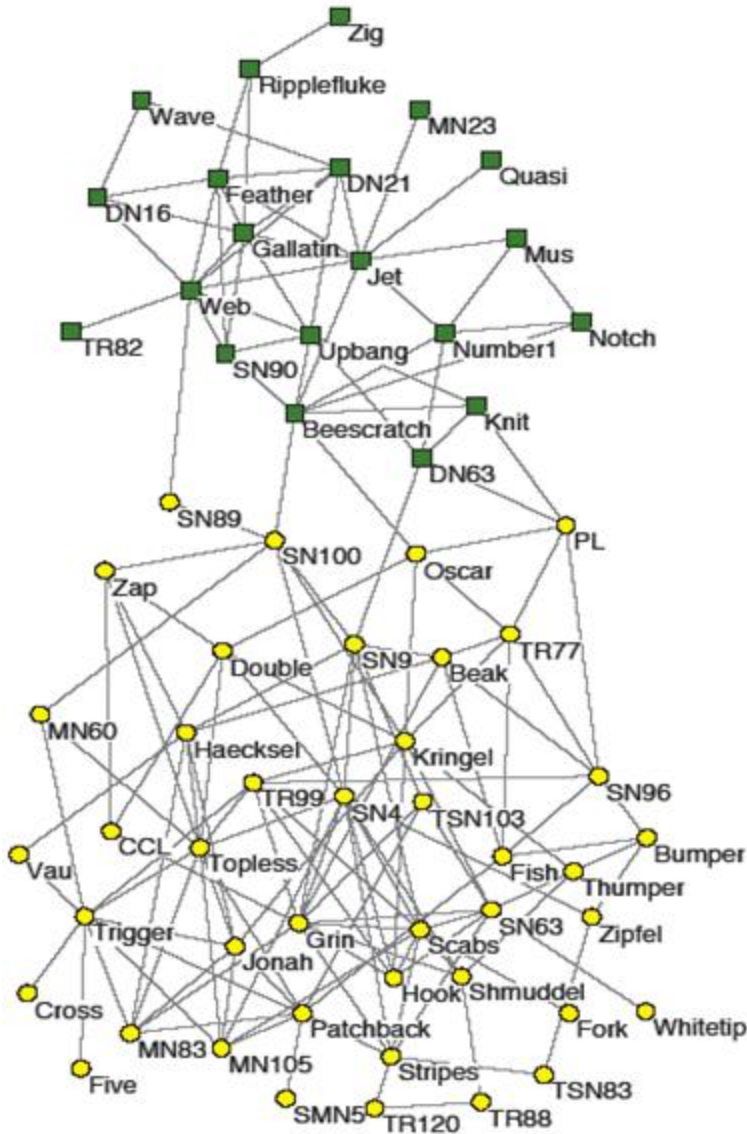


- Mạng hợp tác giữa mạng các nhà khoa học làm việc tại học viện Santa Fe (SFI). Các màu chỉ dẫn cộng đồng ở mức độ cao thu được theo thuật toán của Girvan và Newman (mục VA) và tương ứng khá chặt chẽ với các đơn vị nghiên cứu của học viện. Phân chia nhỏ hơn tương ứng với các nhóm nghiên cứu nhỏ hơn, xoay quanh các lãnh đạo dự án.

Đồ thị hiện có 118 đỉnh (các nhà khoa học đại diện cho cư dân tại SFI và cộng tác viên của họ). Các cạnh nối các nhà khoa học đã cùng công bố ít nhất một bài báo. Trực quan cho phép phân biệt được các nhóm chuyên ngành. Trong mạng này, khi quan sát nhiều nhóm, là tác giả của một bài báo thì tất cả cùng liên kết với nhau. Có chỉ một số ít các kết nối giữa hầu hết các nhóm.



# Mạng XH và cộng đồng [For10]

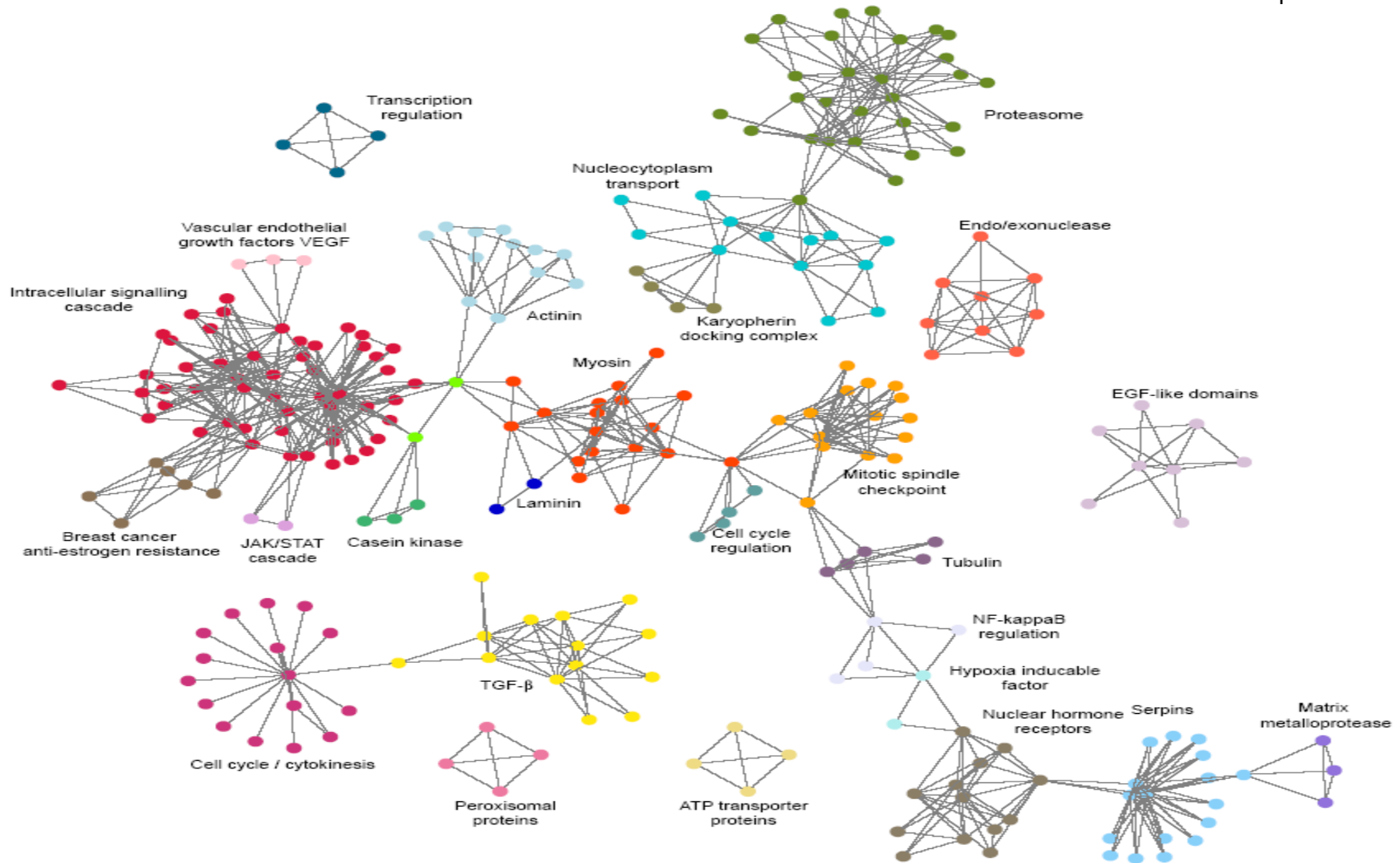


Mạng Lusseau các cá heo mũi to. Những màu nhãn cộng đồng được xác định qua tối ưu hóa một phiên bản mô đun của Newman và Girvan, theo đề xuất của Arenas và cộng sự. Phân hoạch phù hợp với các lớp sinh học của cá heo được Lusseau đề xuất.

Hiện có 62 cá heo, các cạnh nối các cá heo được nhìn thấy thường xuyên hơn so với dự kiến. Tập cá heo bị tách thành hai nhóm sau khi cá heo một trái nơi dành cho một số thời gian (hình vuông và hình tròn trong Hình vẽ). Các nhóm như vậy là khá cố kết, với một vài cụm (clique) nội bộ, và dễ dàng định danh: chỉ có sáu cạnh nối các đỉnh của nhóm khác nhau.

Do mạng này phân lớp tự nhiên cá heo Lusseau, cũng như câu lạc bộ karate của Zachary, thường được dùng để kiểm tra thực nghiệm thuật toán phát hiện cộng đồng

# Mạng XH và cộng đồng: tương tác protein - protein [For10]



# Học máy xác suất Bayes



- Một số kiến thức cơ sở về lý thuyết xác suất
  - Không gian đo được
  - Không gian xác suất
  - Sigma – trường
  - Hệ thống động
  - Quá trình ngẫu nhiên thời gian rời rạc
  - Kỳ vọng
  - Entropy
  - Trong tài liệu
- Nhắc thêm về học máy xác suất
  - ...

# Học máy xác suất Bayes



- **Mô hình tần số**

- Tiến hành thử nghiệm lặp đi lặp lại
- Tỷ số xuất hiện trên toàn bộ số lần thử

- **Mô hình xác suất**

- Xác suất có điều kiện: sự kiện  $e$  với tri thức nền  $D$  là  $P(e|D)$
- Tri thức nền là sự xuất hiện của tài liệu (trái) hoặc sự xuất hiện của tài liệu mới.
- Xác suất tiên nghiệm và xác suất hậu nghiệm.

$$P(e|D) = \frac{P(D|e)P(e)}{P(D)}$$

$$P(e|D, D_2) = \frac{P(D_2|e, D)P(e|D)}{P(D_2|D)}$$

# Ước lượng giá trị tham số



- Hai bài toán

- Lựa chọn mô hình hay dạng hàm: Dựa trên tri thức miền đã có
- Mỗi mô hình/hàm có bộ tham số tương ứng
- Cần xác định bộ tham số này

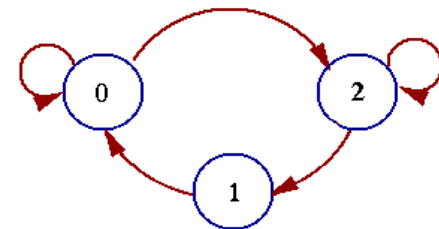
- Xác định tham số

- Thường theo ước lượng
- Ước lượng tham số mô hình
- Ước lượng tham số cho trường hợp cụ thể

# Thuật toán Viterbi

## • Thuật toán Viterbi

- Mô hình máy trạng thái hữu hạn
  - ❖ xác định tham số mô hình phù hợp tập ví dụ học
- Lý thuyết quyết định hỗn hợp
- Bài toán giải mã
  - ❖ Đã có mô hình máy trạng thái hữu hạn
  - ❖ Tìm dãy trạng thái phù hợp nhất với trường hợp cụ thể
- Nội dung thuật toán
  - ❖ Xem trong giáo trình



Input:  $Z=z_1, z_2, \dots, z_n$  // dãy quan sát đầu vào

Output: Đường đi ngắn nhất tương ứng với dãy quan sát đầu vào

Khởi tạo:

$k \leftarrow 1$  // chỉ số lặp

$S(c_1) \leftarrow c_1$

$L(c_1) \leftarrow 0$  // biên chứa tổng độ dài, khởi tạo là 0

Đệ quy:

repeat

For  $\forall$  bộ chuyển  $t_k = (c_k, c_{k+1})$

$L(c_k, c_{k+1}) \leftarrow L(c_k) + 1 [t_k = (c_k, c_{k+1})]$  theo  $\forall c_k$

Tìm  $L(c_{k+1}) = \min L(c_k, c_{k+1})$

For mỗi  $c_{k+1}$

lưu  $L(c_{k+1})$  và vết  $S(c_{k+1})$  tương ứng

$k \leftarrow k+1$

until  $k=n$

# C4. Một số bài toán xử lý tiếng Việt



- **Lĩnh vực xử lý ngôn ngữ tự nhiên**
  - Xử lý ngôn ngữ tự nhiên (tự động hóa)
  - Ra đời khoảng những năm 1950
  - Ngày càng phát triển
- **Phân loại**
  - Xử lý
    - ❖ Cơ bản
    - ❖ Ứng dụng
  - Tài nguyên
    - Cơ bản
    - Mức cao

# Bài toán tách câu



- Đây là bài toán khá đơn giản
- Khái niệm
  - Chuỗi ký tự kết thúc bằng dấu chấm, chấm hỏi, chấm than
  - Vẫn có trường hợp ngoại lệ (khoảng 10%)
    - ❖ Các dấu trên không kết thúc câu (trong nháy kép)
    - ❖ Một số dấu khác kết thúc câu
- Một số trường hợp
  - Dựa theo kinh nghiệm
  - Mô hình ME (Le Hong Phuong & Ho Tuong Vinh)
  - Xem giáo trình



# Bài toán tách từ



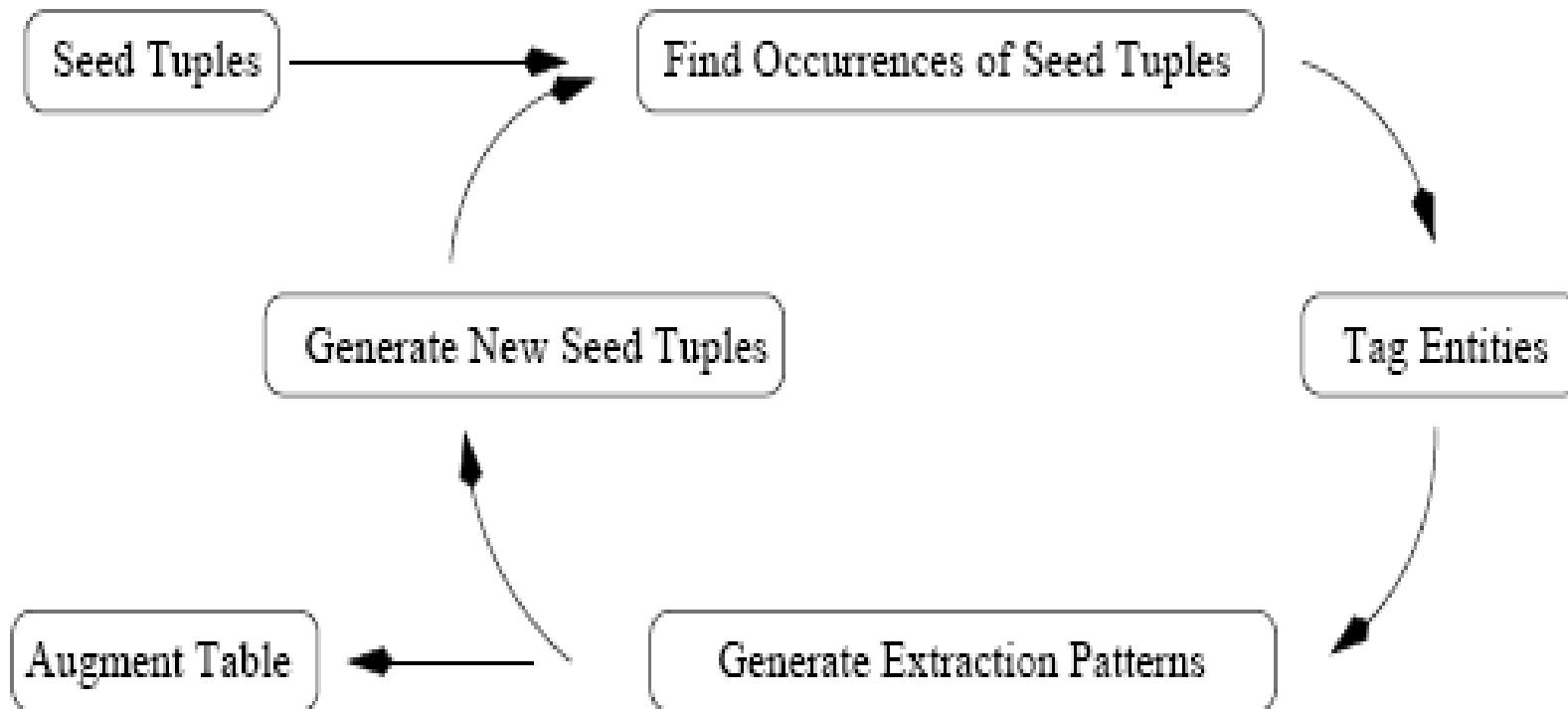
- Đây là bài toán rất cơ bản, luôn thời sự
  - Từ vẫn phát triển bổ sung, thay đổi
  - Ngăn cách hiển, nhập nhằng, mờ
  - “Ông già đi nhanh quá” | “Học sinh học sinh học” ...
- **Khái niệm**
  - Cho một câu hãy xác định các từ trong câu
  - “Phù hợp ngữ cảnh”
- **Một số phương pháp**
  - Khớp từ đa
  - Máy trạng thái hữu hạn có trọng số (WSFT)
    - ❖ Trường ngẫu nhiên có điều kiện
  - Xem giáo trình

# Phát hiện quan hệ ngữ nghĩa



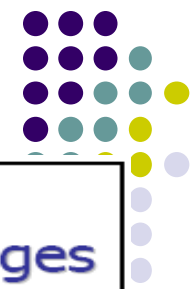
- Là bài toán cơ bản
  - Quan hệ ngữ nghĩa giữa các đối tượng ngữ pháp
  - Một số quan hệ ngữ nghĩa: theo cách tiếp cận
- Khái niệm
  - Cho một tập các văn bản
  - Tìm ra các đối tượng ngữ pháp và các quan hệ giữa chúng
- Một số phương pháp
  - DIPRE
  - SNOWBALL
  - Xem giáo trình

# Phương pháp Snowball

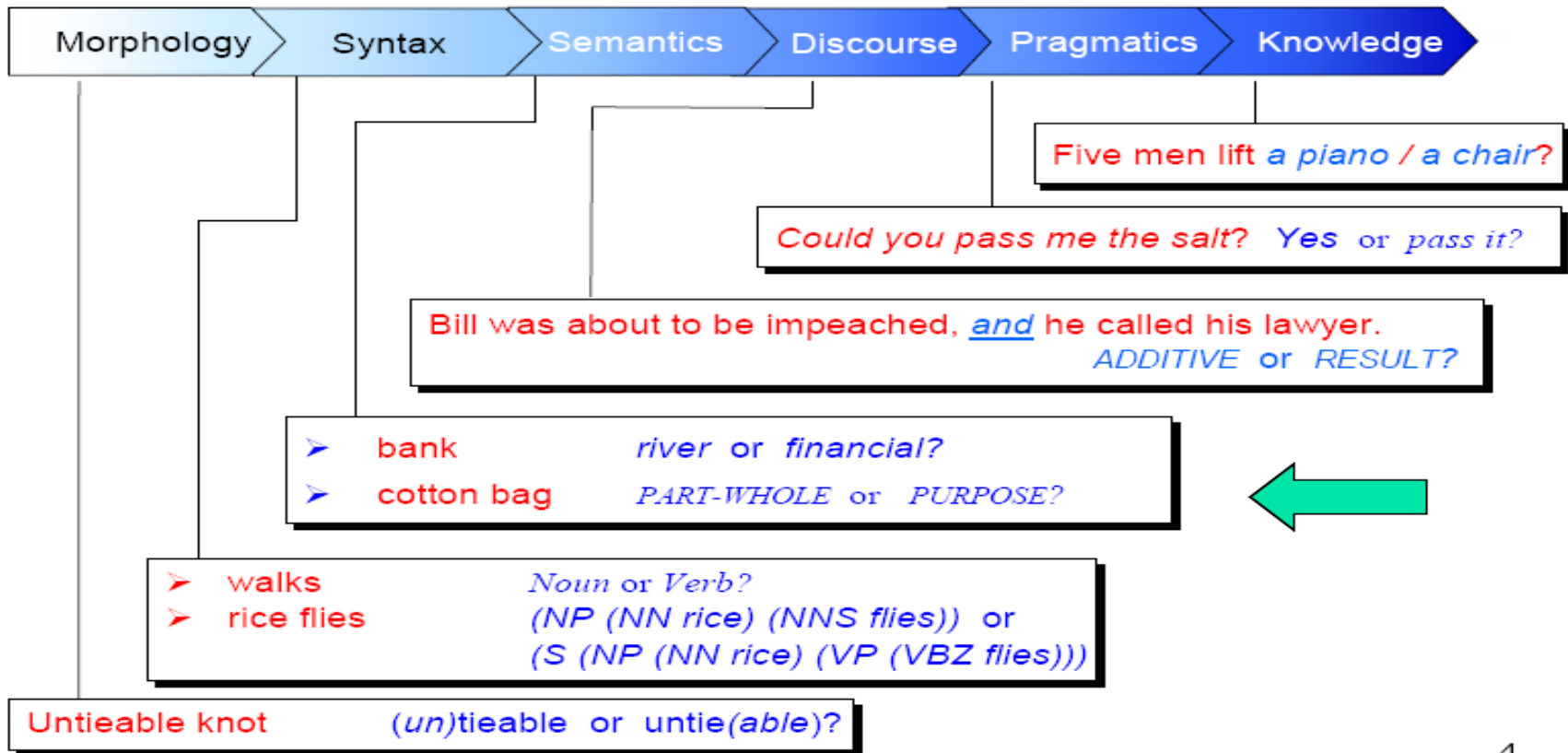


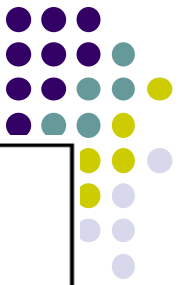
Eugene Agichtein, Luis Gravano (2000). Snowball: extracting relations from large plain-text collections, *ACM DL 2000*: 85-94

# Phát hiện quan hệ ngữ nghĩa



## Levels of Language Analysis - Computational challenges





## Lists of Semantic Relations: Approaches in Linguistics (10)

---

(Levi 1978):

- Two syntactic processes are used:
  - predicate nominalization:
    - those involving nominalizations, i.e., compounds whose heads are nouns derived from a verb, and whose modifiers are interpreted as arguments of the related verb
    - E.g.: "x such that x plans cities" => **city planner**;
  - predicate deletion:
    - List of relations: **cause, have, make, use, be, in, for, from, about**
    - E.g.: "**field mouse**" derived from "a mouse which is in the field" ("in" deletion);
    - Deleted predicates represent the only semantic relations which can underlie NCs not formed through predicate nominalization;

# Quan hệ ngữ nghĩa: XLNNTN



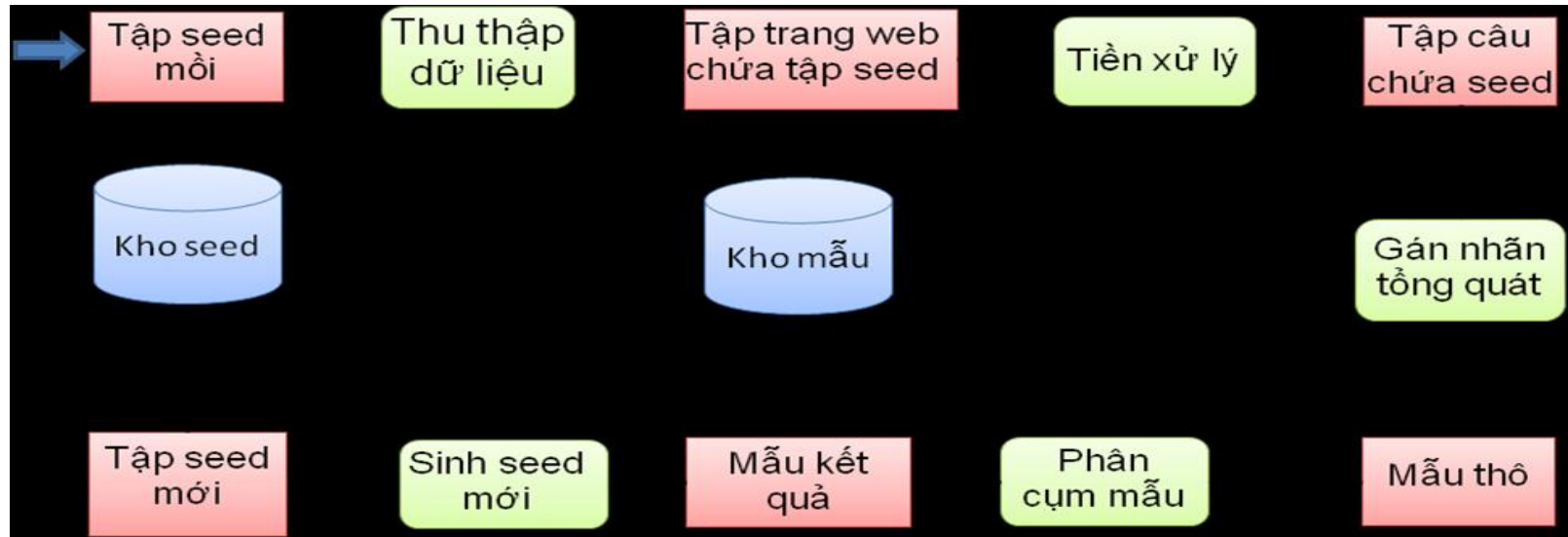
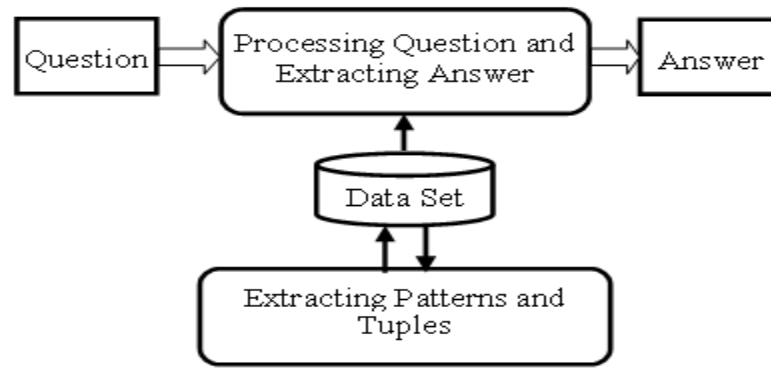
Semantic Relation	Definition/ Example
HYPERNYMY (IS-A)	an entity/event/state is a subclass of another; ( <i>daisy flower; large company, such as Microsoft</i> )
PART-WHOLE (MERONYMY)	an entity/event/state is a part of another entity/event/state; ( <i>door knob; the door of the car</i> );
CAUSE	an event/state makes another event/state to take place; ( <i>malaria mosquitos; "death by hunger"; "The earthquake generated a big Tsunami"</i> );
POSSESSION	an animate entity possesses (owns) another entity; ( <i>family estate; the girl has a new car.</i> )
KINSHIP	an animated entity related by blood, marriage, adoption or strong affinity to another animated entity; ( <i>boy's sister; Mary has a daughter</i> )
MAKE/PRODUCE	an animated entity creates or manufactures another entity; ( <i>honey bees; GM makes cars</i> )
INSTRUMENT	an entity used in an event as instrument; ( <i>pump drainage; He broke the box with a hammer.</i> )
TEMPORAL	time associated with an event; ( <i>5-0' clock tea; the store opens at 9 am</i> )
LOCATION/ SPACE	spacial relation between two entities or between an event and an entity; ( <i>field mouse; I left the keys in the car</i> )
PURPOSE	a state/activity intended to result from another state/event; ( <i>migraine drug; He was quiet in order not to disturb her.</i> )
SOURCE/FROM	place where an entity comes from; ( <i>olive oil</i> )

# Quan hệ ngữ nghĩa: XLNNTN



EXPERIENCER	an animated entity experiencing a state/feeling; ( <i>desire for chocolate; Mary's fear.</i> )
TOPIC	an object specializing another object; ( <i>they argued about politics</i> )
MANNER	a way in which an event is performed or takes place; ( <i>hard-working immigrants; performance with passion</i> )
MEANS	the means by which an event is performed or takes place; ( <i>bus service; I go to school by bus.</i> )
AGENT	the doer of an action; ( <i>the investigation of the police</i> )
THEME	the entity acted upon in an action/event ( <i>music lover</i> )
PROPERTY	characteristic or quality of an entity/event/state; ( <i>red rose; the juice has a funny color.</i> )
BENEFICIARY	an animated entity that benefits from the state resulting from an event; ( <i>customer service; I wrote Mary a letter.</i> )
MEASURE	an entity expressing quantity of another entity/event; ( <i>70-km distance; The jacket costs \$60; a cup of sugar</i> )
TYPE	a word/concept is a type of word/concept; ( <i>member state; framework law</i> )
DEPICTION-DEPICTED	an entity is represented in another; ( <i>the picture of the girl</i> )

# Phát hiện quan hệ ngữ nghĩa



Vu Tran, Vinh Nguyen, Uyen Pham, Oanh Tran and Quang Thuy Ha (2009). An Experimental Study of Vietnamese Question Answering System, *International Conference on Asian Language Processing (IALP 2009)*: 152-155, Dec 7-9, 2009, Singapore, <http://www.computer.org/portal/web/csdl/doi/10.1109/IALP.2009.39>

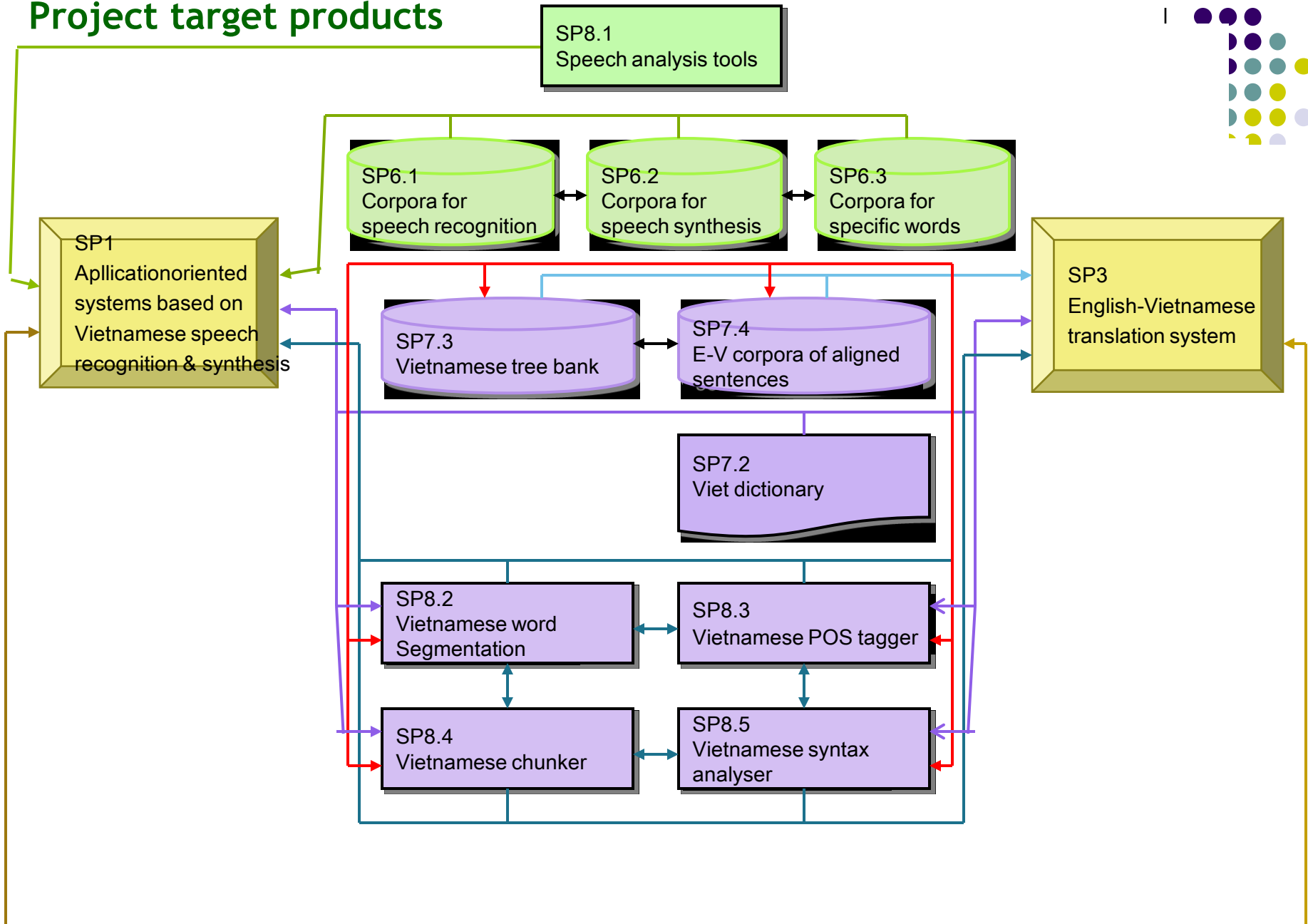


# Một số công cụ nguồn mở



- Chuyển từ trang Web sang văn bản
  - Bộ phân tích HTML (<http://jexpert.us>), Tác giả: Jose Solorzano
  - Một số công cụ tinh chế cho tiếng Việt (html2text.php, text2telex.php <http://203.113.130.205/~cuongnv/thesis/code/tools.tar.gz>). Tác giả: Nguyễn Việt Cường
- Một số bộ công cụ xử lý
  - Nhóm KPLD phát triển (PXHiếu, NCTú, NTTTrang)
    - ❖ Bộ công cụ xử lý Text trên Java: **JtextPro** (<http://jtextpro.sourceforge.net/>) và **JwebPro** (<http://jwebpro.sourceforge.net/>)
    - ❖ Phần mềm phân đoạn từ tiếng Việt: **JvnSegmenter** (<http://jvnsegmenter.sourceforge.net/>)
  - Sản phẩm tài nguyên và công cụ của Đề tài “*Nghiên cứu phát triển một số sản phẩm thiết yếu về xử lý tiếng nói và văn bản tiếng Việt*” mã số KC.01.01/06-10 do PGS, TS. Lương Chi Mai chủ trì.
    - ❖ <http://vlsp.vietlp.org:8080/demo/?page=home>
  - Một số tiện ích liên quan: <http://vnlp.net/blog/?p=229> và <http://vnlp.net/blog/wp-content/uploads/2010/08/Toolkits.pptx>

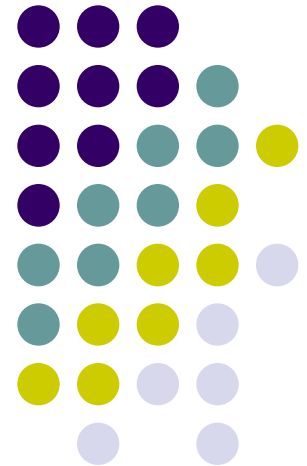
# Project target products



# BÀI GIẢNG KHAI PHÁ DỮ LIỆU WEB

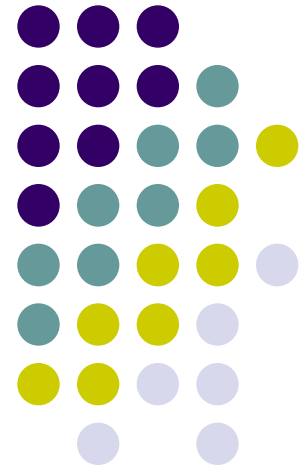
## CHƯƠNG 5. BIỂU DIỄN WEB

PGS. TS. HÀ QUANG THỤY  
HÀ NỘI 02-2011  
**TRƯỜNG ĐẠI HỌC CÔNG NGHỆ**  
ĐẠI HỌC QUỐC GIA HÀ NỘI



# Nội dung

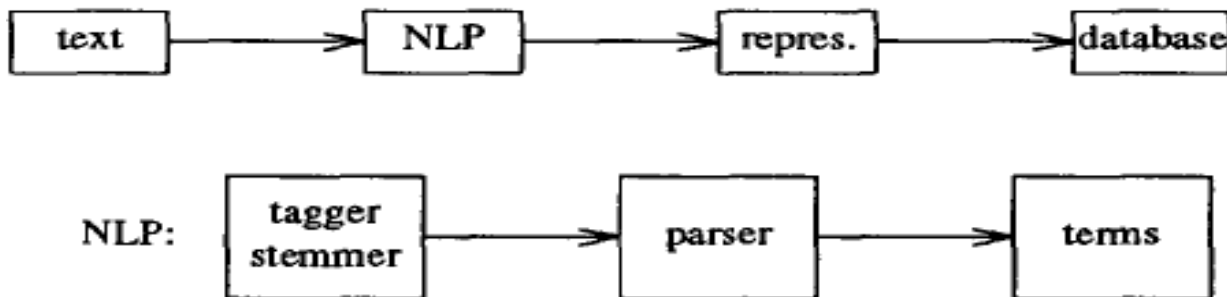
Giới thiệu  
Phân tích văn bản  
Biểu diễn Text  
Lựa chọn đặc trưng  
Thu gọn đặc trưng  
Biểu diễn Web



# Giới thiệu



- **Biểu diễn văn bản**
  - Là bước cần thiết đầu tiên trong xử lý văn bản
  - Phù hợp đầu vào của thuật toán khai phá dữ liệu
  - Tác động tới chất lượng kết quả của thuật toán KHDL
  - Thuật ngữ tiếng Anh: (document/text) (representation/indexing)
- **Phạm vi tác động của một phương pháp biểu diễn văn bản**
  - Không tồn tại phương pháp biểu diễn lý tưởng
  - Tồn tại một số phương pháp biểu diễn phổ biến
  - Chọn phương pháp biểu diễn phù hợp miền ứng dụng
- **Một sơ đồ sơ lược:** Tomek Strzalkowski: Document Representation in Natural Language Text Retrieval, *HLT 1994*: 364-369



# Nghiên cứu về biểu diễn văn bản



- **Nghiên cứu biểu diễn văn bản (Text + Web)**
  - Luôn là nội dung nghiên cứu thời sự
  - Biểu diễn Web bổ sung một số yếu tố cho biểu diễn Text
- **Số công trình liên quan**
  - "Document representation"
    - mọi nơi: 8000 bài; tiêu đề: 200 (60 bài từ 2006-nay)
  - "Document indexing"
    - mọi nơi: 5200 bài; tiêu đề: 220 (60 bài từ 2006-nay)
  - "Text representation"
    - mọi nơi: 9200 bài; tiêu đề: 240 (60 bài từ 2006-nay)
  - "Text indexing"
    - mọi nơi: 6800 bài; tiêu đề: 210 (60 bài từ 2006-nay)

Ghi chú: các bài “ở mọi nơi” phần đông thuộc vào các bài toán xử lý văn bản bao gồm bước trình bày văn bản

# Nghiên cứu về biểu diễn văn bản (2)



Research paper reference	Document Representation	Feature Selection	Learning algorithm
Apté et al. [6]	bag-of-words (freq)	stop list+ frequency	Decision Rules
Armstrong et al. [7]	bag-of-words	informativity	TFIDF Winnow, WordStat
Balabanović et al. [9]	bag-of-words (freq)	stop list+stemming+ keep 10 best words	TFIDF
Bartell et al. [11]	bag-of-words (freq)	latent semantic indexing using SVD	—
Berry et al. [12] Foltz and Dumais [28]	bag-of-words(freq)	latent semantic indexing using SVD	TFIDF
Cohen [21]	bag-of-words	infrequent words pruned	Decision Rules ILP
Joachims [40]	bag-of-words (freq)	infrequent words+ informativity	TFIDF, PrTFIDF, Naive Bayes
Lam et al. [60]	bag-of-words (freq)	mutual info.	Bayesian Network
Lewis et al. [66]	bag-of-words	log likelihood ratio	logistic regression with Naive Bayes
Maes [69]	bag-of-words+ header info.	mail/news header + selecting keywords	Memory-Based reasoning
Pazzani et al. [83, 84]	bag-of-words	stop list+ informativity	TFIDF, Naive Bayes, Nearest Neighbor, Neural Network, Decision Trees
Sorensen and Mc Elligott [97, 25]	n-gram graph (only bigrams)	weighting graph edges	connectionist combined with Genetic Algorithms
Yang [100]	bag-of-words	stop list	adapted k-Nearest Neighbor

Dunja Mladenic' (1998). Machine Learning on Non-homogeneous, Distributed Text Data. *PhD. Thesis*, University of Ljubljana, Slovenia.



# Phân tích văn bản

- **Mục đích biểu diễn văn bản (Keen, 1977 [Lew91])**
  - Từ được chọn liên quan tới chủ đề người dùng quan tâm
  - Gắn kết các từ, các chủ đề liên quan để phân biệt được từ ở các lĩnh vực khác nhau
  - Dự đoán được độ liên quan của từ với yêu cầu người dùng, với lĩnh vực và chuyên ngành cụ thể
- **Môi trường biểu diễn văn bản (đánh chỉ số)**
  - Thủ công / từ động hóa. Thủ công vẫn có hỗ trợ của công cụ máy tính và phần mềm
  - Điều khiển: chọn lọc từ làm đặc trưng (feature) biểu diễn) / không điều khiển: mọi từ đều được chọn.
  - Từ điển dùng để đánh chỉ số. Từ đơn và tổ hợp từ.



# Luật Zipt



## Luật Zipt

- Cho dãy dữ liệu được xếp hạng  $x_1, x_2, \dots, x_n$

thì hạng tuân theo công thức

$C$  là hằng số, gần 1; kỳ vọng dạng loga

- Dạng hàm mật độ:

$$x_{(r)} = \frac{C}{r^\alpha}$$

$$E(\log x_{(r)}) = c - \alpha \log(r)$$

$$p(x) = \frac{C^{1/\alpha}}{\alpha n} \frac{1}{x^{(1/\alpha)+1}} = \frac{A}{x^\beta}$$

## Một số dạng khác

- Phân phối Yule

$$x_{(r)} = \frac{C}{r^\alpha B^r}$$

- Mô hình thống kê

$c = \log(C)$ ,  $b = \log(B)$

$$E(\log x_{(r)}) = c - \alpha \log(r) - b e^{\log(r)}$$

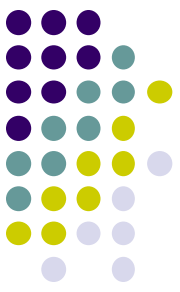
- Biến thể loga-chuẩn

$$E(\log x_{(r)}) = c - \alpha \log(r) - b(\log(r))^2$$

- Phân phối Weibull với  $0 < \alpha < 1$

$$E(\log x_{(r)}) = c - \alpha \log(r) - b e^{\beta \log(r)}$$

# Luật Zipt trong phân tích văn bản

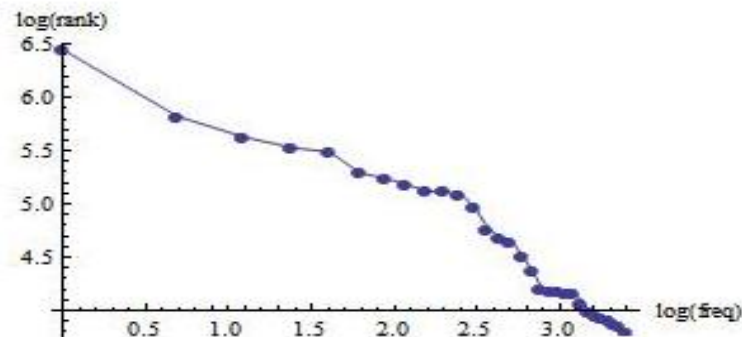


- Trọng số của từ trong biểu diễn văn bản (Luhn, 1958)
  - Dấu hiệu nhấn mạnh: một biểu hiện của độ quan trọng
    - thường viết lặp lại các từ nhất định khi phát triển ý tưởng
    - hoặc trình bày các lập luận,
    - phân tích các khía cạnh của chủ đề. ...
  - Các từ có tần suất xuất hiện cao nhất lại ít ngữ nghĩa. Từ xuất hiện trung bình lại có độ liên quan cao.

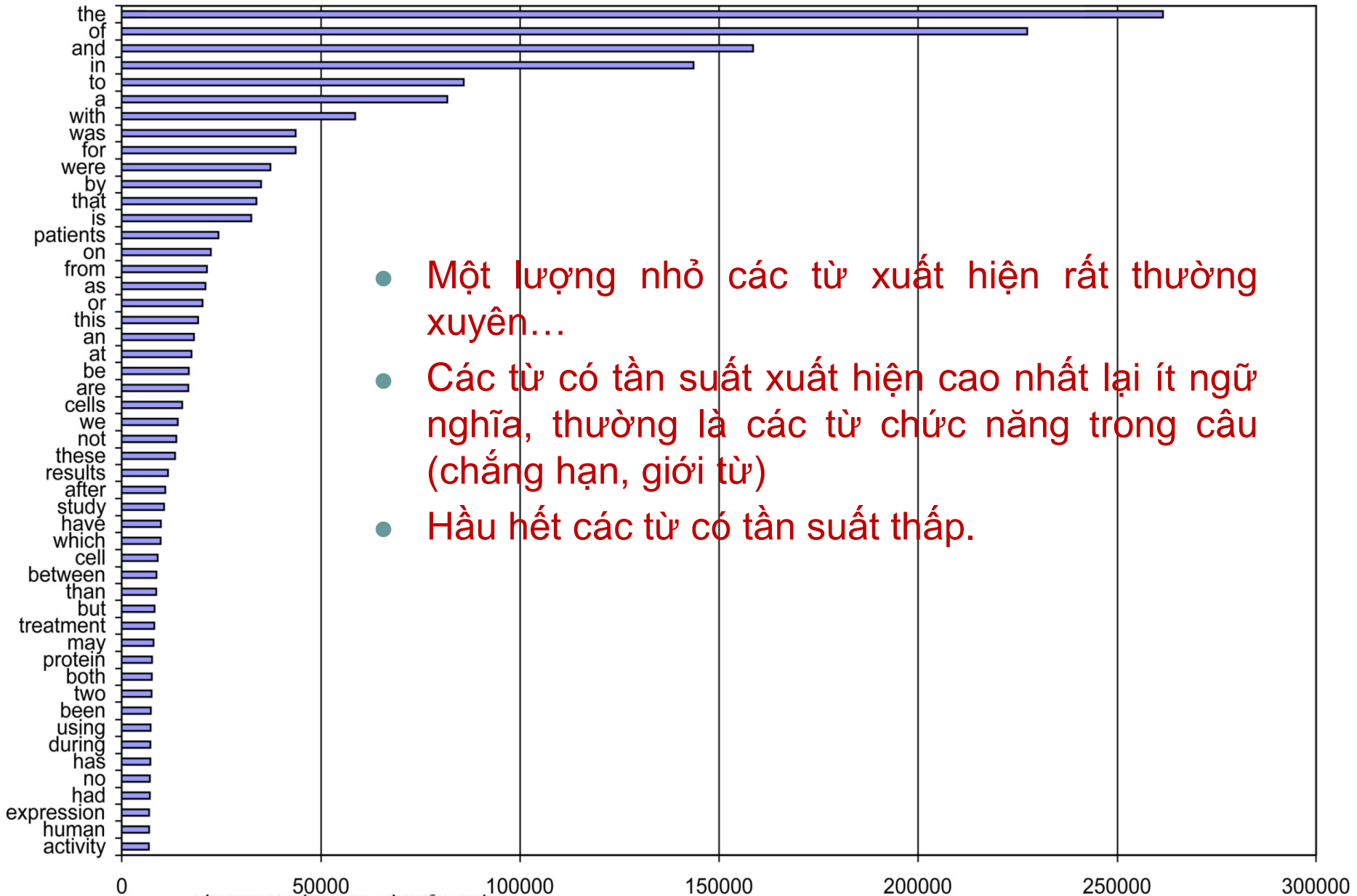
## Luật Zipt

- Là một quan sát hiện tượng mà không phải là luật thực sự: xem hình vẽ “Alice ở xứ sở mặt trời”
- $r_t * f_t = K$  (hằng số):  $r_t$  : độ quan trọng của từ  $t$ ;  $f_t$ : tần số xuất hiện từ  $t$ . Có thể logarith

the 632	and 338	a 278
to 252	she 242	of 199
it 189	i 178	was 167
alice 167	in 163	said 144
you 118	her 108	that 105
as 91	at 79	with 67
s 66	had 65	all 64
on 64	little 59	out 54
down 52	this 51	t 50
for 48	but 47	they 45



# Luật Zipt trong tiếng Anh

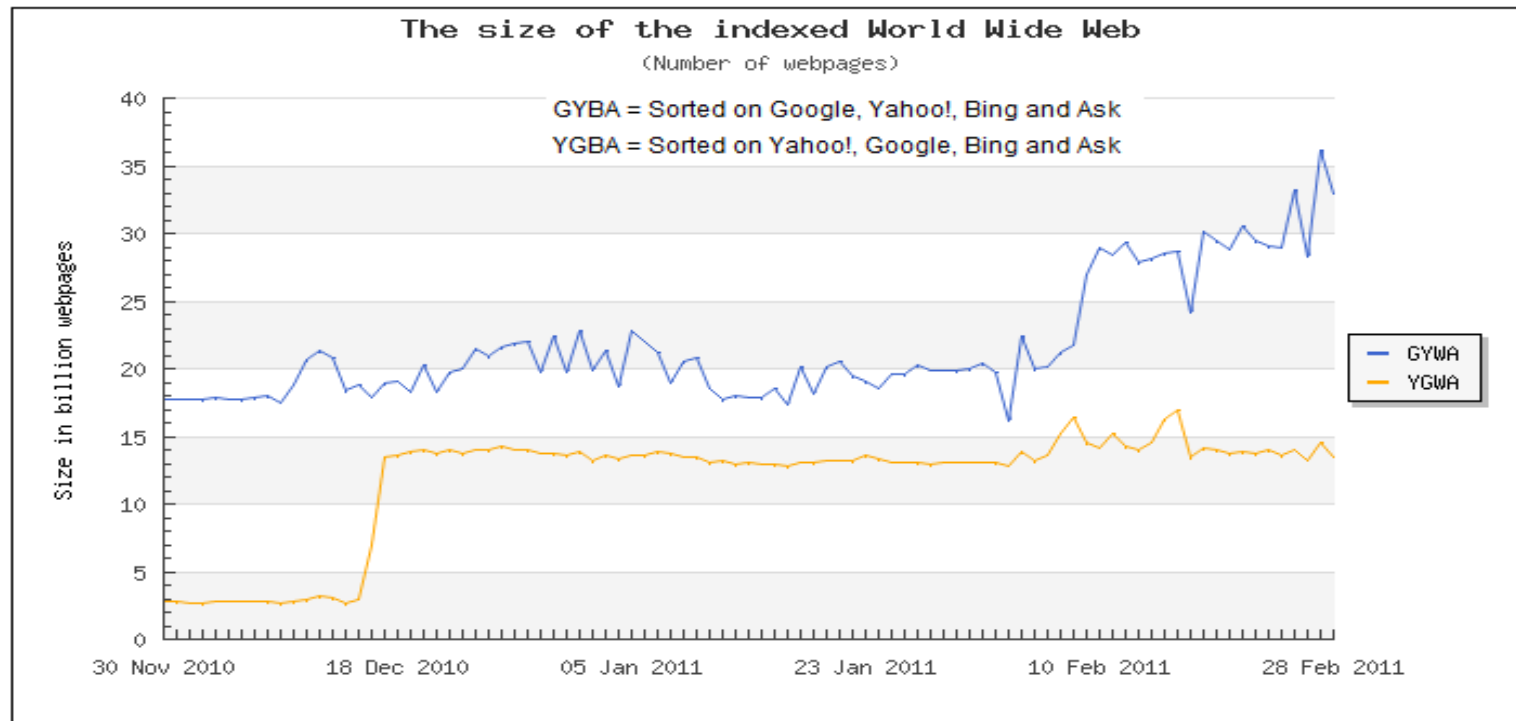


- Một lượng nhỏ các từ xuất hiện rất thường xuyên...
- Các từ có tần suất xuất hiện cao nhất lại ít ngữ nghĩa, thường là các từ chức năng trong câu (chẳng hạn, giới từ)
- Hầu hết các từ có tần suất thấp.

# Luật Zipt: ước lượng trang web được chỉ số



- Ước lượng tối thiểu lượng trang web chỉ số hóa
  - <http://www.worldwidewebsite.com/>
  - Luật Zipt: từ kho ngữ liệu DMOZ có hơn 1 triệu trang web
  - Dùng luật Zipt để ước tính lượng trang web chỉ số hóa.
  - Mỗi ngày: 50 từ (đều ở đoạn logarithm luật Zipt) gửi tới 4 máy tìm kiếm Google, Bing, Yahoo Search và Ask.
  - Trừ bớt phần giao ước tính giữa các công cụ tìm kiếm: làm già
  - Thứ tự trừ bớt phần giao → tổng (được làm non)



# Các mẫu luật Zipt khác



- Dân số thành phố
  - Dân số thành phố trong một quốc gia: có  $\sum = 1$ . Đã xác nhận ở 20 quốc gia.
  - Có thể mở rộng sang: dân cư khu đô thị, vùng lãnh thổ
- Lượt thăm trang web và mẫu giao vận Internet khác
  - Số lượt truy nhập trang web/tháng
  - Các hành vi giao vận Internet khác
- Quy mô công ty và một số số liệu kinh tế khác
  - Xếp hạng công ty theo: số nhân viên, lợi nhuận, thị trường
  - Các hành vi giao vận Internet khác
- ...

[Li02] Wentian Li (2002). Zipf's Law Everywhere, *Glottometrics* 5 (2002): 14-21

# Phương pháp lựa chọn từ Luhn58



- **Bài toán**

- Input: Cho một tập văn bản: có thể coi tất cả các văn bản trong miền ứng dụng; ngưỡng trên, ngưỡng dưới dương.
- Output: Tập từ được dùng để biểu diễn văn bản trong tập

- **Giải pháp**

- Tính tần số xuất hiện mỗi từ đơn nhất trong từng văn bản
- Tính tần số xuất hiện của các từ trong tập toàn bộ văn bản
- Sắp xếp các từ theo tần số giảm dần
- Loại bỏ các từ có tần số xuất hiện vượt quá ngưỡng trên hoặc nhỏ thua ngưỡng dưới.
- Các từ còn lại được dùng để biểu diễn văn bản
- “Từ” được mở rộng thành “đặc trưng”: n-gram, chủ đề..

- **Lưu ý**

- Chọn ngưỡng: ngưỡng cố định, ngưỡng được điều khiển
- Liên hệ vấn đề chọn lựa đặc trưng (mục sau).

# Phương pháp đánh trọng số của từ



- Bài toán

- Input: Cho một tập văn bản miền ứng dụng  $D$  và tập từ được chọn biểu diễn văn bản  $V$  (sau bước trước đây).
- Output: Đánh trọng số từ trong mỗi văn bản  $\Rightarrow$  Xây dựng ma trận  $\{w_{i,j}\}$  là trọng số của từ  $w_i$  trong văn bản  $d_j$ .

- Giải pháp

- Một số phương pháp điển hình
- Boolean
- dựa theo tần số xuất hiện từ khóa
- Dựa theo nghịch đảo tần số xuất hiện trong các văn bản

- Phương pháp Boolean

- Đơn giản: trọng số là xuất hiện/ không xuất hiện
- $w_{i,j} = 1$  nếu  $w_i$  xuất hiện trong văn bản  $d_j$ , ngược lại  $w_{i,j} = 0$ .

# Các phương pháp đánh trọng số của từ theo tần số



- **Dạng đơn giản: TF**
    - $w_{i,j} = f_{i,j}$ : trong đó  $f_{i,j}$  là số lần từ khóa  $w_i$  xuất hiện trong văn bản  $d_j$
  - **Một số phiên bản khác của dạng đơn giản**
    - Cân đối số lần xuất hiện các từ khóa: giảm chênh lệch số lần xuất hiện
    - Giảm theo hàm căn  $w_{i,j} = \sqrt{tf_{ij}}$
    - Tránh giá trị “0” và giảm theo hàm loga:  $w_{i,j} = 1 + \log(f_{i,j})$
  - **Nghịch đảo tần số xuất hiện trong tập văn bản: IDF**
    - Từ xuất hiện trong nhiều văn bản thì trọng số trong 1 văn bản sẽ thấp
    - $w_i = \log\left(\frac{m}{df_i}\right) = \log(m) - \log(df_i)$
- Trong đó  $m = |D|$ ,  $df_i$  là  $|d \in D : w_i \text{ xuất hiện trong } d|$





# Phương pháp TFIDF

- Tích hợp TF và IDF

- Dạng đơn giản:  $w_{i,j} = f_{i,j} * df_i / m$
- Dạng căn chỉnh theo hàm loga

$$w_{i,j} = \begin{cases} \log(tf_{ij}) \log\left(\frac{m}{df_i}\right) & : tf_{ij} > 0 \\ 0 & : tf_{ij} = 0 \end{cases}$$

- Ngoài ra, có một số dạng tích hợp trung gian khác

# Mô hình biểu diễn văn bản



- Bài toán

- Input: Cho tập văn bản miền ứng dụng  $D = \{d_j\}$ , tập đặc trưng được chọn biểu diễn văn bản  $V = \{w_i\}$ , ma trận trọng số  $W = (w_{i,j})$ .
- Output: Tìm biểu diễn của các văn bản  $d_j$  ■■■.

- Một số mô hình

- Mô hình Boolean
- Mô hình không gian vector
- Mô hình túi các từ (Mô hình xác suất)
- Các mô hình khác

- Mô hình Boolean

- Tập các từ thuộc  $V$  mà xuất hiện trong văn bản

# Mô hình không gian vector



- Nội dung chính

- Ánh xạ tập tài liệu vào không gian vector  $n = |V|$  chiều.
- Mỗi tài liệu được ánh xạ thành 1 vector

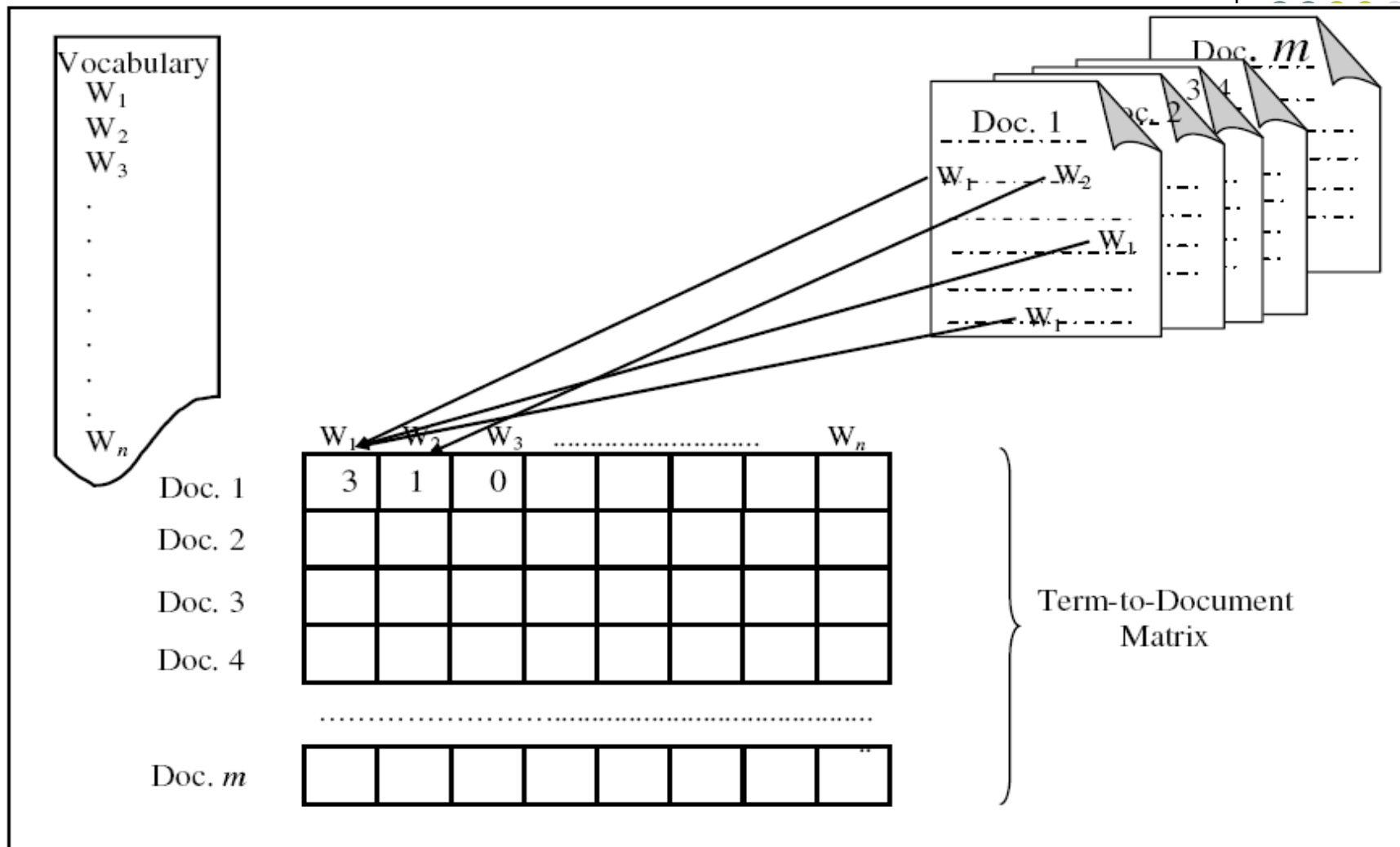
$$d_i \Leftrightarrow (w_{i1}, w_{i2}, \dots, w_{in})$$

- Độ đo tương tự nội dung văn bản

- Chuẩn hóa vector: đưa về độ dài 1
  - Độ “tương tự nội dung” giữa hai văn bản  $\Leftrightarrow$  độ tương tự giữa hai vector
  - Một số phương án sơ khai “các thành phần giống nhau”, “nghịch đảo khoảng cách”, ..
- Phổ biến là tính độ đo cosin của góc giữa hai vector: không yêu cầu chuẩn hóa

$$\text{sim}(d_1, d_2) = \frac{(v_1, v_2)}{\|v_1\| \|v_2\|} = \frac{\sum_i v_{1i} * v_{2i}}{\sqrt{\sum_i v_{1i}^2} * \sqrt{\sum_i v_{2i}^2}}$$

# Mô hình không gian vector



Khaled Shaban (2006). A semantic graph model for text representation and matching in document mining, *PhD Thesis*, University of Waterloo, Canada

# Mô hình xác suất



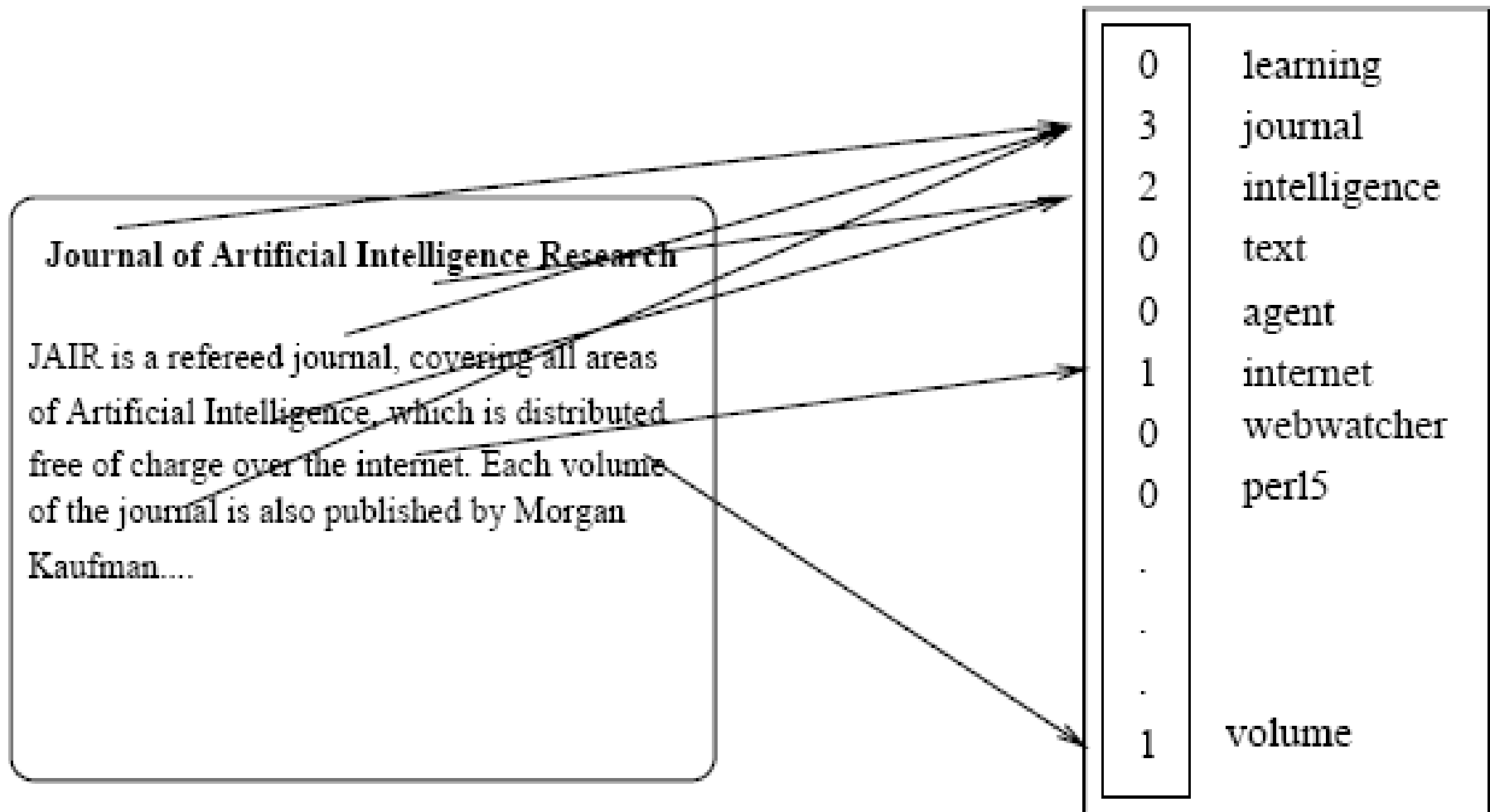
- **Giả thiết chính**

- Mô hình xác suất: cặp  $(Y, P)$  với  $Y$  là tập quan sát được và  $P$  là mô hình xác suất trên  $Y$  (có thể coi  $Y$  là quan sát được các từ/đặc trưng trên văn bản).
- Các từ xuất hiện trong văn bản thể hiện nội dung văn bản
- Sự xuất hiện của các từ là độc lập lẫn nhau và độc lập ngữ cảnh
- Dạng đơn giản: chỉ liệt kê từ, dạng chi tiết: liệt kê từ và số lần xuất hiện
- Lưu ý: Các giả thiết về tính độc lập không hoàn toàn đúng (độc lập lẫn nhau, độc lập ngữ cảnh) song mô hình thi hành hiệu quả trong nhiều trường hợp.

- **Độ đo tương tự nội dung văn bản**

- So sánh hai túi từ

# Mô hình túi từ (bag-of-words)



# Mô hình biểu diễn LSI và theo phân cụm



## ● Giới thiệu

- Tồn tại nhiều phương pháp biểu diễn khác
- Tồn tại nhiều phiên bản cho một phương pháp
- Gần đây có một số phương pháp mới
- Hai phương pháp phổ biến: LSI và theo phân cụm
- Lưu ý: Giá phải trả khi tiền xử lý dữ liệu

## ● Mô hình phân cụm

- Phân cụm các từ trong miền ứng dụng: ma trận trọng số
- Thay thế từ bằng cụm chứa nó

## ● Mô hình biểu diễn LSI

- LSI: Latent Semantic Indexing biểu diễn ngữ nghĩa ẩn
  - Nâng mức ngữ nghĩa (trừu tượng) của đặc trưng
  - Rút gọn tập đặc trưng, giảm số chiều không gian biểu diễn
  - Không gian từ khóa  $\Rightarrow$  không gian khái niệm (chủ đề).
- Phương pháp chuyển đổi
  - Ma trận trọng số  $\Rightarrow$  ma trận hạng nhỏ hơn
  - Phép biến đổi đó Từ khóa  $\Rightarrow$  khái niệm. Thay thế biểu diễn.

# Lựa chọn từ trong biểu diễn văn bản



- **Loại bỏ từ dừng**
  - Những từ được coi là không mang nghĩa
  - Có sẵn trong ngôn ngữ
- **Đưa về từ gốc**
  - Các ngôn ngữ có biến dạng từ: Anh, Nga...
  - Thay từ biến dạng về dạng gốc
- **Chọn đặc trưng n-gram**
  - Các âm tiết liền nhau n-gram
    - Uni-gram: chỉ chứa một âm tiết
    - Bigram: chứa không quá 2 âm tiết
    - Trigram: chứa không quá 3 âm tiết
    - N-gram: Thường không quá 4 gram
  - Một số đặc trưng
    - Chính xác hơn về ngữ nghĩa
    - Tăng số lượng đặc trưng
    - Tăng độ phức tạp tính toán



# Một số độ đo cho lựa chọn đặc trưng



- **Giới thiệu chung**
  - Lựa chọn đặc trưng: lợi thế chính xác, lợi thế tốc độ hoặc cả hai
  - Các độ đo giúp khẳng định lợi thế
- **Phân nhóm độ đo**
  - Hai nhóm: theo tần số và theo lý thuyết thông tin
- **Một số độ đo điển hình**
  - Xem hai trang sau

# Một số đo cho lựa chọn đặc trưng



$P(t_k, c_i)$  kí hiệu là xác suất của từ  $t_k$  có trong chủ đề  $c_i$  và  $P(t_k, \bar{c}_i)$  là xác suất của từ  $t_k$  không có trong chủ đề  $c_i$ .

1. DIA (Darmstadt Indexing Approach – Tiếp cận đánh chỉ số Darmstadt): Được đề xuất bởi Fuhn và đồng nghiệp [FHK91].

$$f(t_k, c_i) = z(t_k, c_i) = P(c_i|t_k)$$

2. Độ đo IG (Information Gain).

$$f(t_k, c_i) = IG(t_k, c_i) = \sum_{c \in \{c_i, \bar{c}_i\}} \sum_{t \in \{t_k, \bar{t}_k\}} P(t, c) \cdot \log \frac{P(t, c)}{P(t) \cdot P(c)}$$

3. Độ đo thông tin tương hỗ (mutual information).

$$f(t_k, c_i) = \log \frac{P(t_k, c_i)}{P(t_k) \cdot P(c_i)}$$

4. Độ đo Khi - bình phương (Chi-square).

$$f(t_k, c_i) = \chi^2(t_k, c_i) = \frac{|Tr| \cdot [P(t_k, c_i) \cdot P(\bar{t}_k, \bar{c}_i) - P(t_k, \bar{c}_i) \cdot P(\bar{t}_k, c_i)]^2}{P(t_k) \cdot P(\bar{t}_k) \cdot P(c_i) \cdot P(\bar{c}_i)}$$

5. Độ đo liên quan (Relevancy score).

$$f(t_k, c_i) = RS(t_k, c_i) = \log \frac{P(t_k|c_i) + d}{P(\bar{t}_k|\bar{c}_i) + d}$$

6. Tỷ lệ dư (Odd Ratio).

$$f(t_k, c_i) = OR(t_k, c_i) = \frac{P(t_k|c_i) \cdot (1 - P(t_k|\bar{c}_i))}{(1 - P(t_k|c_i)) \cdot P(t_k, |c_i)}$$

# Một số đo cho toàn bộ các lớp



Các độ đo trên là tính cho từng lớp. Độ đo cho toàn bộ các lớp trong tập hợp có thể được tính theo nhiều cách khác nhau,

$$f(t_k) = \sum_{i=1}^{|C|} f(t_k, c_i)$$

hoặc

$$f(t_k) = \sum_{i=1}^{|C|} P(c_i) f(t_k, c_i)$$

hoặc

$$f(t_k) = \max_{i=1}^{|C|} f(t_k, c_i).$$

# Thu gọn đặc trưng



- **Giới thiệu chung**

- “Tối ưu hóa” chọn tập đặc trưng
  - Số lượng đặc trưng nhỏ hơn
  - Hy vọng tăng tốc độ thi hành
  - Tăng cường chất lượng khai phá văn bản. ? Giảm đặc trưng đi là tăng chất lượng: có các đặc trưng “nhiều”
  - Hoặc cả hai mục tiêu trên

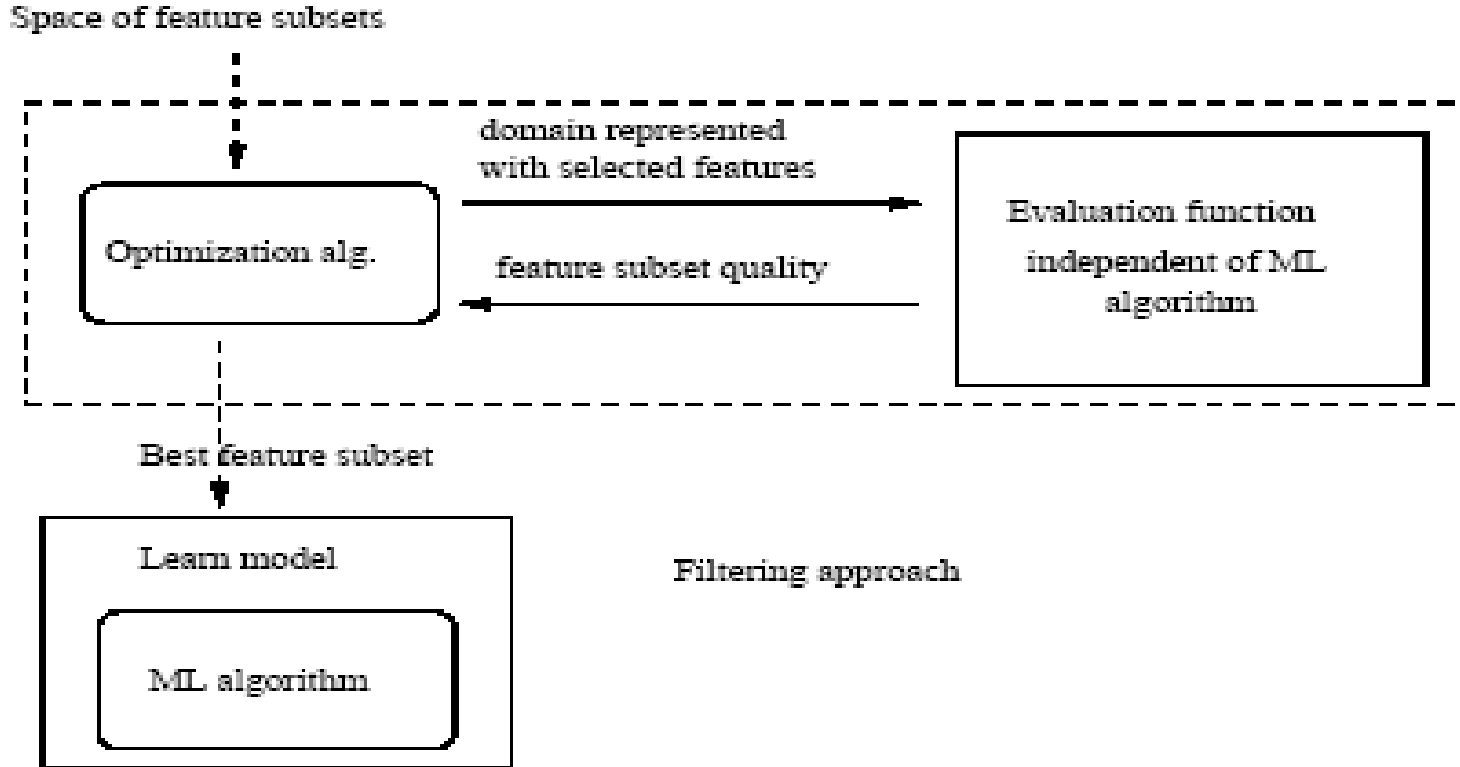
- **Hai tiếp cận điển hình**

- Tiếp cận lọc
- Tiếp cận bao gói

- **Với dữ liệu văn bản**

- Tập đặc trưng: thường theo mô hình vector
- Tính giá trị của từng đặc trưng giữ lại các đặc trưng được coi là “tốt”.

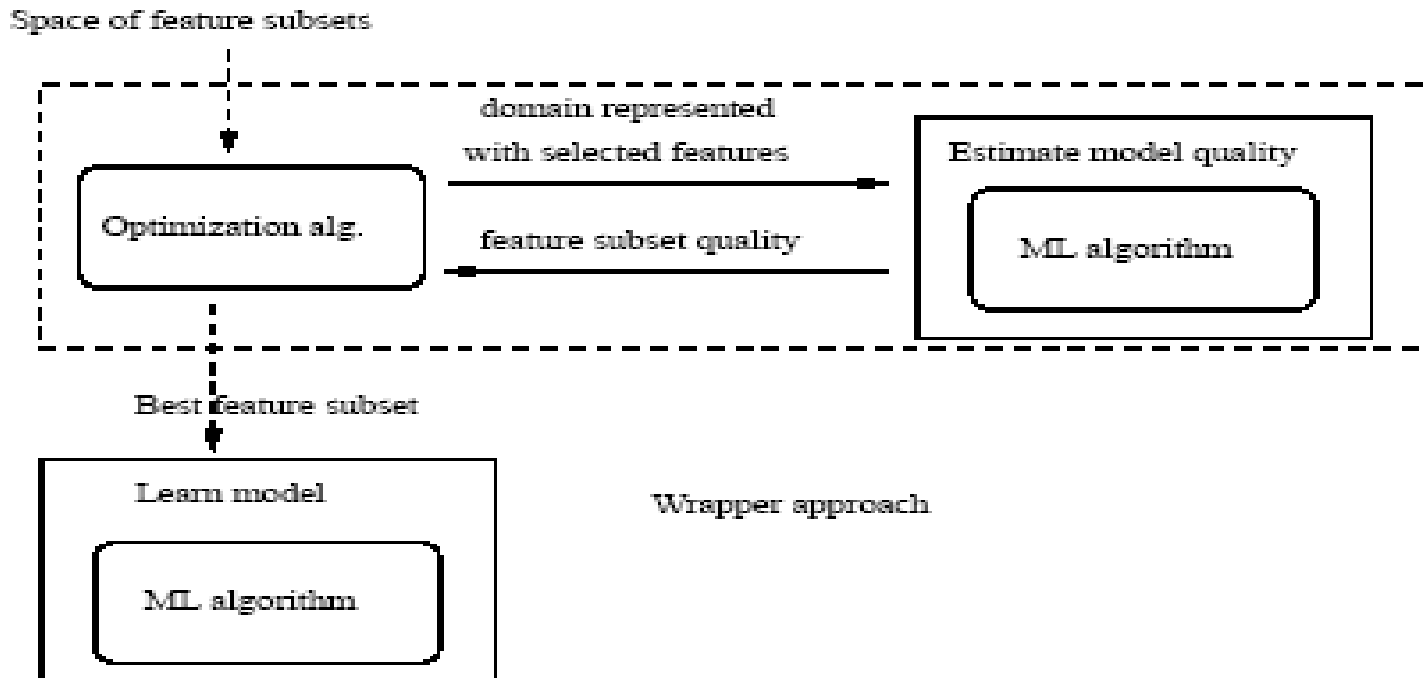
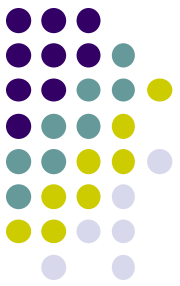
# Tiếp cận tổng quát: lọc



- **Tiếp cận lọc**

- Đầu vào: Không gian tập các tập đặc trưng
- Đầu ra: Tập con đặc trưng tốt nhất
- Phương pháp
  - Dò tìm “cải tiến” bộ đặc trưng: Thuật toán tối ưu hóa
  - Đánh giá chất lượng mô hình: độc lập với thuật toán học máy

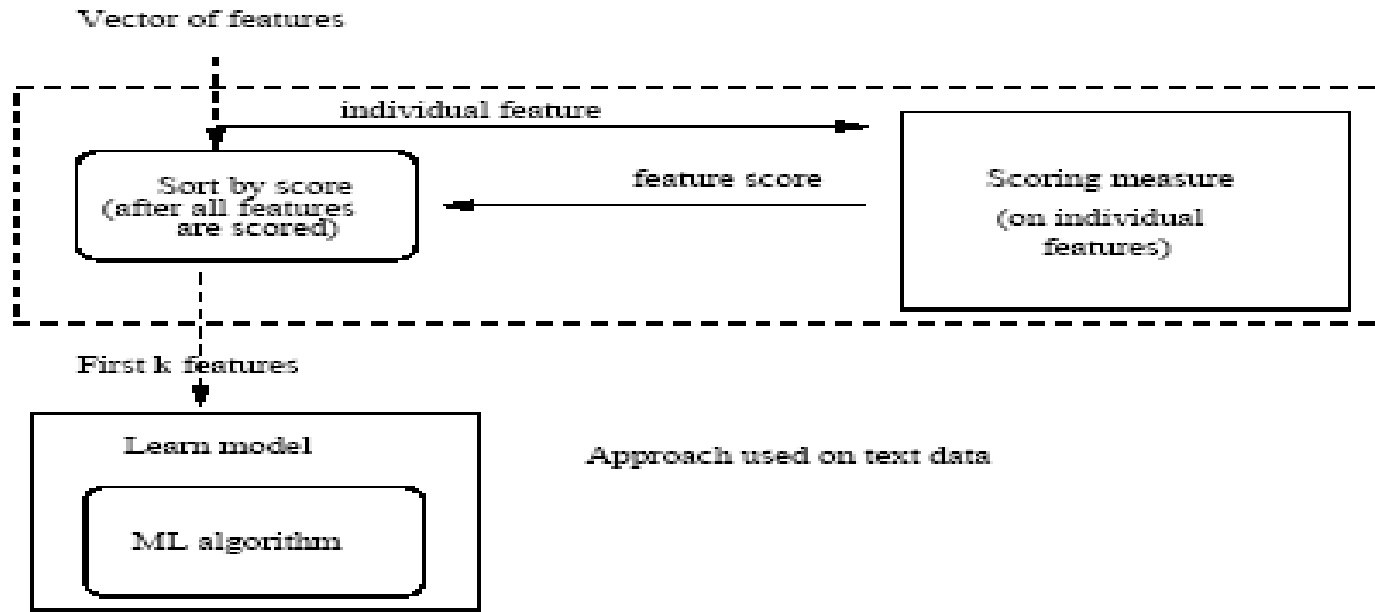
# Tiếp cận bao gói tổng quát



- **Tiếp cận bao gói**

- Đầu vào: Không gian tập các tập đặc trưng
- Đầu ra: Tập con đặc trưng tốt nhất
- Phương pháp
  - Dò tìm “cải tiến” bộ đặc trưng: Thuật toán tối ưu hóa
  - Đánh giá chất lượng mô hình: Dùng chính thuật toán học để đánh giá

# Thu gọn đặc trưng văn bản text



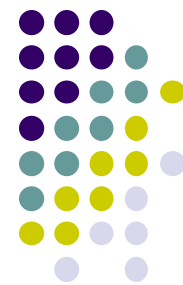
- **Thu gọn đặc trưng**

- Đầu vào: Vector đặc trưng
- Đầu ra: k đặc trưng tốt nhất
- Phương pháp (lùi)
  - Sắp xếp các đặc trưng theo độ “tốt” (để loại bỏ bớt)
  - Tính lại độ “tốt” của các đặc trưng
  - Chọn ra k-đặc trưng tốt nhất

- **Các kiểu phương pháp**

- Tiến / Tiến bậc thang (có xem xét thay thế khi tiến)
- Lùi / Lùi bậc thang (có xem xét thay thế khi lùi)

# Thu gọn đặc trưng phân lớp text nhị phân



```
SELECTFEATURES( $\mathbb{D}, c, k$ )
1   $V \leftarrow \text{EXTRACTVOCABULARY}(\mathbb{D})$ 
2   $L \leftarrow []$ 
3  for each  $t \in V$ 
4  do  $A(t, c) \leftarrow \text{COMPUTEFEATUREUTILITY}(\mathbb{D}, t, c)$ 
5     APPEND( $L, \langle A(t, c), t \rangle$ )
6  return FEATURESWITHLARGESTVALUES( $L, k$ )
```

- **Một thuật toán lựa chọn đặc trưng text**
  - $V$ : Bảng từ vựng có được từ tập văn bản  $D$
  - $c$ : lớp đang được quan tâm
  - giá trị  $A(t,c)$ : một trong ba thủ tục tính toán
- **Ba kiểu thủ tục tính toán  $A(t,c)$** 
  - Thông tin tương hỗ
  - Lựa chọn đặc trưng theo khi-bình phương (chi-square)
  - Lựa chọn đặc trưng theo tần suất



# Thu gọn đặc trưng: thông tin tương hỗ



$$I(U;C) = \sum_{e_t \in \{1,0\}} \sum_{e_c \in \{1,0\}} P(U = e_t, C = e_c) \log_2 \frac{P(U = e_t, C = e_c)}{P(U = e_t)P(C = e_c)}$$

$$I(U;C) = \frac{N_{11}}{N} \log_2 \frac{NN_{11}}{N_1.N_1} + \frac{N_{01}}{N} \log_2 \frac{NN_{01}}{N_0.N_1} + \frac{N_{10}}{N} \log_2 \frac{NN_{10}}{N_1.N_0} + \frac{N_{00}}{N} \log_2 \frac{NN_{00}}{N_0.N_0}$$

- Công thức MI (Mutual Information)

- Biến ngẫu nhiên U: từ khóa t xuất hiện/không xuất hiện
- Biến ngẫu nhiên c: tài liệu thuộc/không thuộc lớp c
- Ước lượng cho MI

- Ví dụ: Bộ dữ liệu Reuter-RCV1

- Lớp poultry, từ khóa export

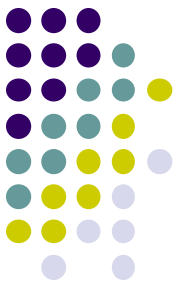
	$e_t = e_{\text{export}} = 1$	$e_t = e_{\text{export}} = 0$
$e_c = e_{\text{poultry}} = 1$	$N_{11} = 49$	$N_{01} = 141$
$e_c = e_{\text{poultry}} = 0$	$N_{10} = 27,652$	$N_{00} = 774,106$

$$I(U;C) = \frac{49}{801,948} \log_2 \frac{801,948 \cdot 49}{(49+27,652)(49+141)} + \frac{141}{801,948} \log_2 \frac{801,948 \cdot 141}{(141+774,106)(49+141)}$$

$$+ \frac{27,652}{801,948} \log_2 \frac{801,948 \cdot 27,652}{(49+27,652)(27,652+774,106)} + \frac{774,106}{801,948} \log_2 \frac{801,948 \cdot 774,106}{(141+774,106)(27,652+774,106)}$$

$$\approx 0.000105$$

# 10 đặc trưng tốt nhất cho 6 lớp



*UK*

london	0.1925
uk	0.0755
british	0.0596
stg	0.0555
britain	0.0469
plc	0.0357
england	0.0238
pence	0.0212
pounds	0.0149
english	0.0126

*China*

china	0.0997
chinese	0.0523
beijing	0.0444
yuan	0.0344
shanghai	0.0292
hong	0.0198
kong	0.0195
xinhua	0.0155
province	0.0117
taiwan	0.0108

*poultry*

poultry	0.0013
meat	0.0008
chicken	0.0006
agriculture	0.0005
avian	0.0004
broiler	0.0003
veterinary	0.0003
birds	0.0003
inspection	0.0003
pathogenic	0.0003

*coffee*

coffee	0.0111
bags	0.0042
growers	0.0025
kg	0.0019
colombia	0.0018
brazil	0.0016
export	0.0014
exporters	0.0013
exports	0.0013
crop	0.0012

*elections*

election	0.0519
elections	0.0342
polls	0.0339
voters	0.0315
party	0.0303
vote	0.0299
poll	0.0225
candidate	0.0202
campaign	0.0202
democratic	0.0198

*sports*

soccer	0.0681
cup	0.0515
match	0.0441
matches	0.0408
played	0.0388
league	0.0386
beat	0.0301
game	0.0299
games	0.0284
team	0.0264

Bộ dữ liệu Reuter-RCV1

# Thống kê khi-bình phương và tần số



$$X^2(\mathbb{D}, t, c) = \sum_{e_t \in \{0,1\}} \sum_{e_c \in \{0,1\}} \frac{(N_{e_t e_c} - E_{e_t e_c})^2}{E_{e_t e_c}}$$

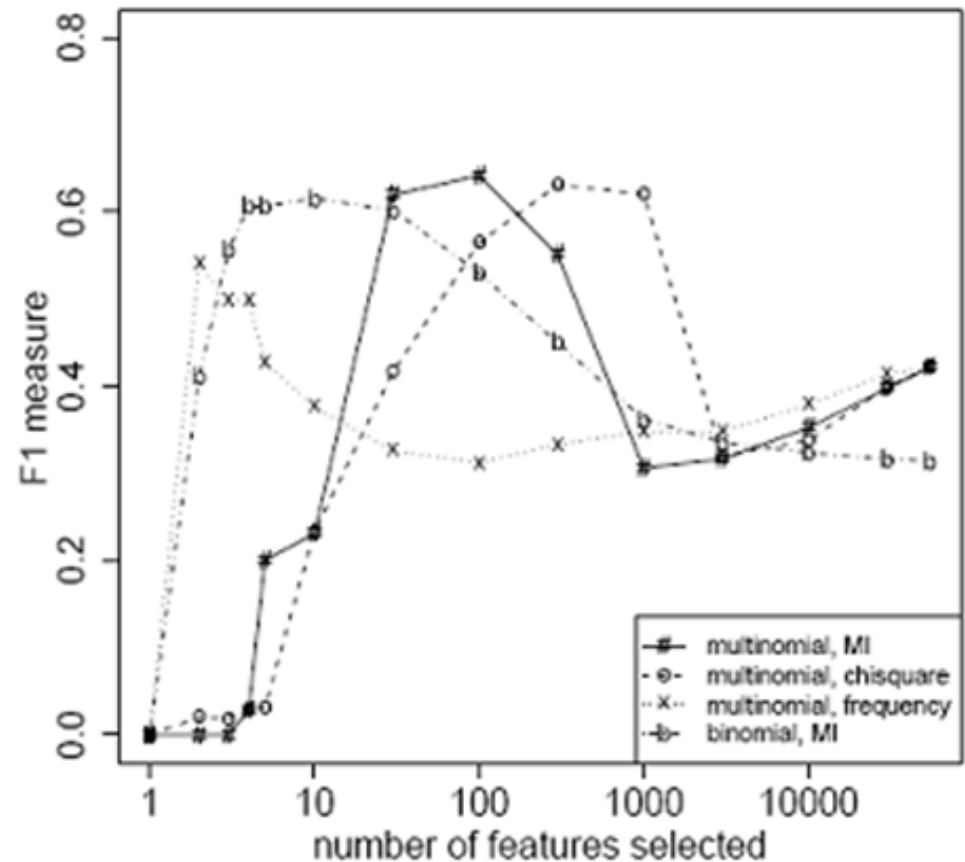
$$X^2(\mathbb{D}, t, c) = \frac{(N_{11} + N_{10} + N_{01} + N_{00}) \times (N_{11}N_{00} - N_{10}N_{01})^2}{(N_{11} + N_{01}) \times (N_{11} + N_{10}) \times (N_{10} + N_{00}) \times (N_{01} + N_{00})}$$

- **Thống kê khi-bình phương**

- Công thức xác suất:  $e_t$ ,  $e_c$  : như MI, các biến  $E$  là kỳ vọng,  $N$  là tần số quan sát được từ tập tài liệu  $D$
- Ước lượng cho MI: các giá trị  $N$  như MI

- **Tần số**

- Một ước lượng xác suất



# Thu gọn đặc trưng phân lớp text đa lớp



- Bài toán phân lớp đa lớp
  - Tập  $C = \{c_1, c_2, \dots, c_n\}$
  - Cần chọn đặc trưng tốt nhất cho bộ phân lớp đa lớp
- Phương pháp thống kê khi-bình phương
  - Mỗi từ khóa
    - Lập bảng xuất hiện/không xuất hiện các đặc trưng trong lớp văn bản
    - Tính giá trị thống kê khi-bình phương
  - Chọn k đặc trưng (từ khóa)
- Phương pháp lựa chọn từng lớp
  - Tính bộ đặc trưng tốt cho từng phân lớp thành phần
  - Kết hợp các bộ đặc trưng tốt
    - Tính toán giá trị kết hợp: trung bình (có trọng số xuất hiện) khi kết hợp
    - Chọn k-đặc trưng tốt nhất sau khi tính toán kết hợp

# Biểu diễn Web



- **Đồ thị Web**

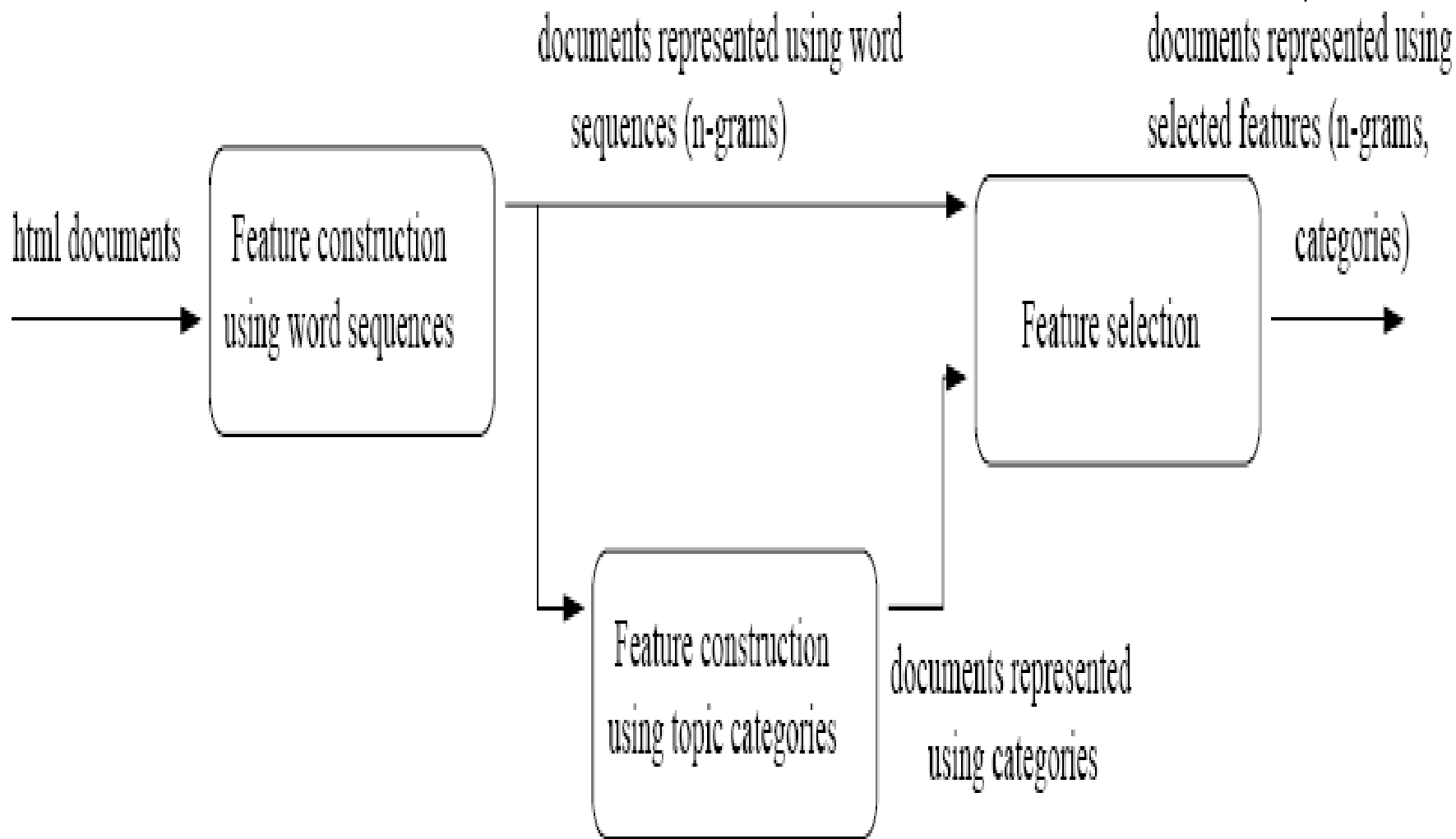
- Web có cấu trúc đồ thị
  - Đồ thị Web: nút  $\leftrightarrow$  trang Web, liên kết ngoài  $\leftrightarrow$  cung (có hướng, vô hướng).
  - Bản thân trang Web cũng có tính cấu trúc cây (đồ thị)
- Một vài bài toán đồ thị Web
  - Biểu diễn nội dung, cấu trúc
  - Tính hạng các đối tượng trong đồ thị Web: tính hạng trang, tính hạng cung..

Nghiên cứu về đồ thị Web (xem trang sau)

- **Đồ thị ngẫu nhiên**

- Tính ngẫu nhiên trong khai phá Web
  - WWW có tính ngẫu nhiên: mới, chỉnh sửa, loại bỏ
  - Hoạt động con người trên Web cũng có tính ngẫu nhiên
- Là nội dung nghiên cứu thời sự

# Một sơ đồ biểu diễn tài liệu Web



Dunja Mladenic' (1998). Machine Learning on Non-homogeneous, Distributed Text Data. *PhD. Thesis*, University of Ljubljana, Slovenia.

# Một sơ đồ biểu diễn tài liệu Web



Các biểu diễn vector trang Web

Phương pháp 1:

```
a b c d e f g
1 2 2 0 0 0 0
```

Phương pháp 2:

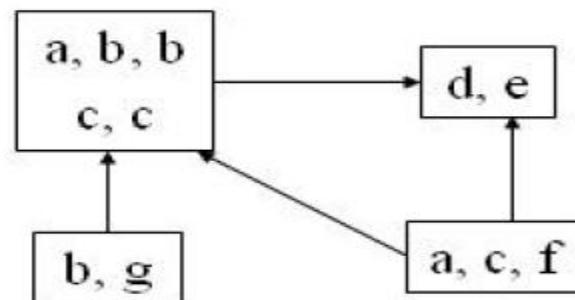
```
a b c d e f g
2 3 3 1 1 1 1
```

Phương pháp 3:

```
Đoạn 1      Đoạn 2
a b c d e f g a b c d e f g
1 2 2 0 0 0 0 1 1 1 1 1 1 1
```

Phương pháp 4:

```
Đoạn1      Đoạn 2      Đoạn 3      Đoạn 4
a b c d e f g a b c d e f g a b c d e f g a b c d e f g
1 2 2 0 0 0 0 0 0 0 1 1 0 0 1 0 1 0 0 1 0 0 1 0 0 0 0 0 1
1 2 2 0 0 0 0 1 0 1 0 0 1 0 0 1 0 0 0 0 0 1 0 0 0 1 1 0 0
```

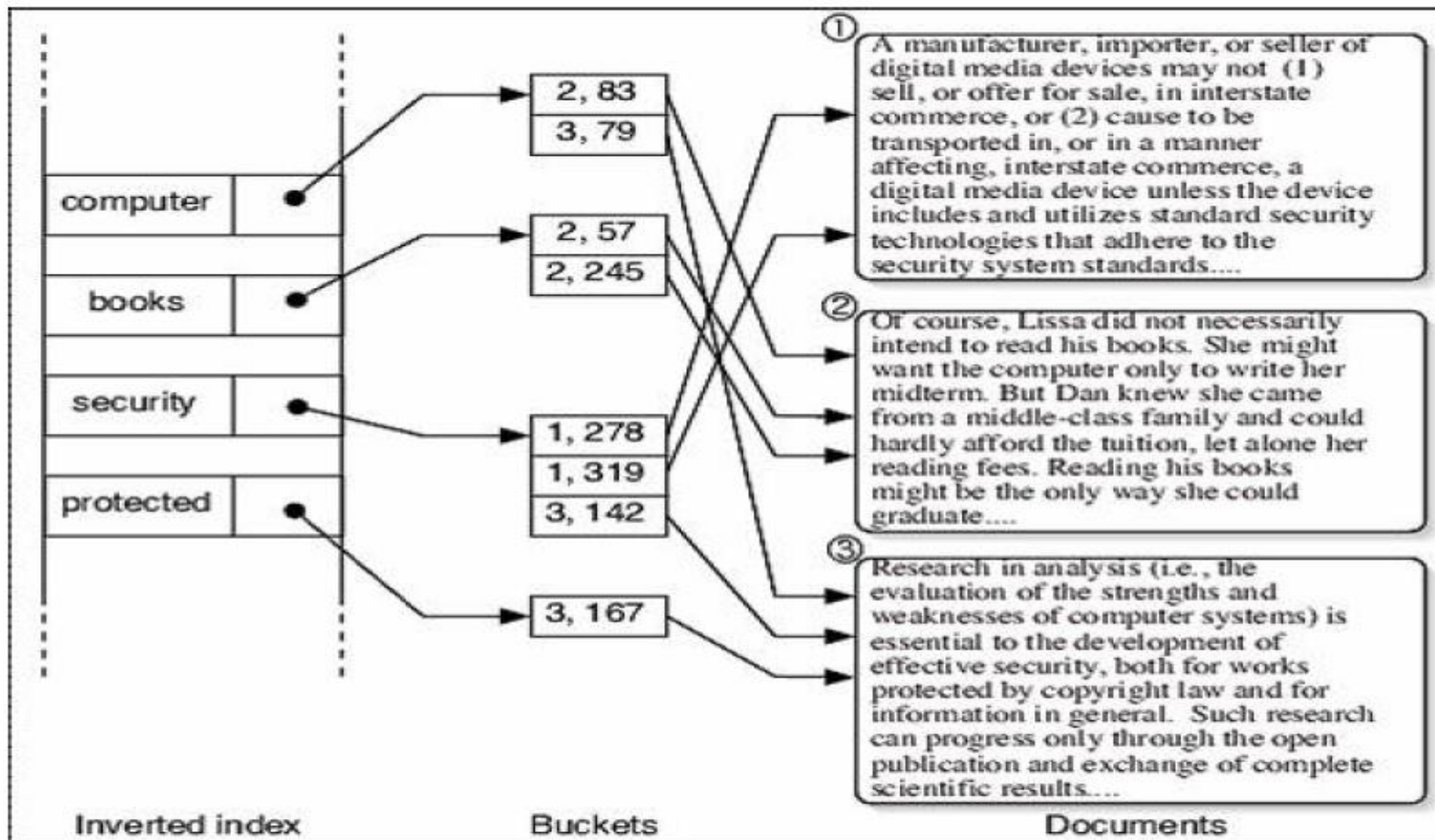




# Một sơ đồ biểu diễn tài liệu Web



Máy tìm kiếm từ khóa nhanh: Hệ thống chỉ số ngược (Inverted Index)

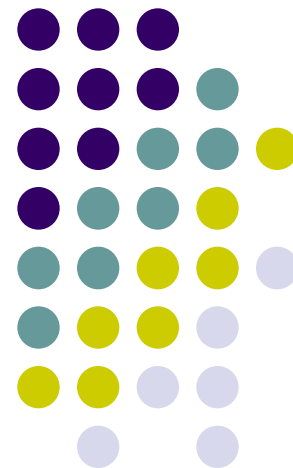




# BÀI GIẢNG NHẬP MÔN KHAI PHÁ DỮ LIỆU

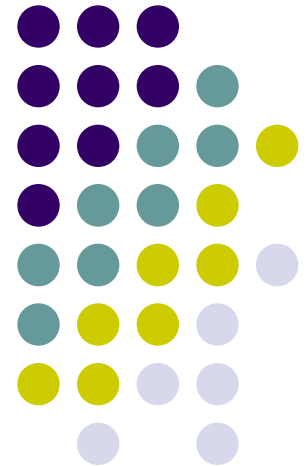
## CHƯƠNG 6. PHÂN CỤM DỮ LIỆU

PGS. TS. HÀ QUANG THỤY  
HÀ NỘI 9-2011  
**TRƯỜNG ĐẠI HỌC CÔNG NGHỆ**  
ĐẠI HỌC QUỐC GIA HÀ NỘI



# Nội dung

Giới thiệu phân cụm  
Thuật toán phân cụm k-min  
Thuật toán phân cụm phân cấp  
Gán nhãn cụm  
Đánh giá phân cụm



# 1. Bài toán phân cụm Web



## ● Bài toán

- Tập dữ liệu  $D = \{d_i\}$
- Phân các dữ liệu thuộc  $D$  thành các cụm
  - Các dữ liệu trong một cụm: “tương tự” nhau (gần nhau)
  - Dữ liệu hai cụm: “không tương tự” nhau (xa nhau)
- Đo “tương tự” (gần) nhau ?
  - *Tiên đề phân cụm*: Nếu người dùng lựa chọn một đối tượng  $d$  thì họ cũng lựa chọn các đối tượng cùng cụm với  $d$
  - Khai thác “cách chọn lựa” của người dùng
  - Đưa ra một số độ đo “tương tự” theo biểu diễn dữ liệu

## ● Một số nội dung liên quan

- Xây dựng độ đo tương tự
- Khai thác thông tin bổ sung
- Số lượng cụm cho trước, số lượng cụm không cho trước

# Sơ bộ tiếp cận phân cụm



- **Phân cụm mô hình và phân cụm phân vùng**
  - Mô hình: Kết quả là mô hình biểu diễn các cụm tài liệu
  - Vùng: Danh sách cụm và vùng tài liệu thuộc cụm
- **Phân cụm đơn định và phân cụm xác suất**
  - Đơn định: Mỗi tài liệu thuộc duy nhất một cụm
  - Xác suất: Danh sách cụm và xác suất một tài liệu thuộc vào các cụm
- **Phân cụm phẳng và phân cụm phân cấp**
  - Phẳng: Các cụm tài liệu không giao nhau
  - Phân cấp: Các cụm tài liệu có quan hệ phân cấp cha- con
- **Phân cụm theo lô và phân cụm tăng**
  - Lô: Tại thời điểm phân cụm, toàn bộ tài liệu đã có
  - Tăng: Tài liệu tiếp tục được bổ sung trong quá trình phân cụm

# Các phương pháp phân cụm



- Các phương pháp phổ biến

- Phân vùng, phân cấp, dựa theo mật độ, dựa theo lưới, dựa theo mô hình, và mờ

- Phân cụm phân vùng

- Xây dựng từng bước phân hoạch các cụm và đánh giá chúng theo các tiêu chí tương ứng
- Độ đo tương tự / khoảng cách
- K-mean, k-mediod
- CLARANS, ...

- Phân cụm phân cấp

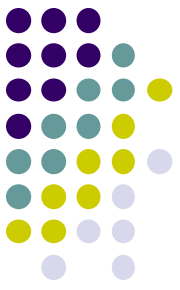
- Xây dựng hợp (tách) dần các cụm tạo cấu trúc phân cấp và đánh giá theo các tiêu chí tương ứng
- Độ đo tương tự / khoảng cách
- HAC: Hierarchical agglomerative clustering
- CHAMELEON, BIRRCH và CURE, ...

# Các phương pháp phân cụm



- **Phân cụm dựa theo mật độ**
  - ❑ Hàm mật độ: Tìm các phần tử chính tại nơi có mật độ cao
  - ❑ Hàm liên kết: Xác định cụm là lân cận phần tử chính
  - ❑ DBSCAN, OPTICS...
- **Phân cụm dựa theo lưới**
  - ❑ Sử dụng lưới các ô cùng cỡ
  - ❑ Tạo phân cấp ô lưới theo một số tiêu chí: số lượng đối tượng trong ô
  - ❑ STING, CLIQUE, WaveCluster...
- **Phân cụm dựa theo mô hình**
  - ❑ Sử dụng một số mô hình giả thiết được phân cụm
  - ❑ Xác định mô hình tốt nhất phù hợp với dữ liệu
  - ❑ MCLUST...
- **Phân cụm mờ**
  - ❑ Giả thiết: không có phân cụm “cứng” cho dữ liệu và đối tượng có thể thuộc một số cụm
  - ❑ Sử dụng hàm mờ từ các đối tượng tới các cụm
  - ❑ FCM (Fuzzy CMEANS),...

# Chế độ và đặc điểm phân cụm web



## ● Hai chế độ

- Trực tuyến: phân cụm kết quả tìm kiếm người dùng
- Ngoại tuyến: phân cụm tập văn bản cho trước

## ● Đặc điểm

- Chế độ trực tuyến: tốc độ phân cụm
  - Web số lượng lớn, tăng nhanh và biến động lớn
  - Quan tâm tới phương pháp gia tăng
- Một lớp quan trọng: phân cụm liên quan tới câu hỏi tìm kiếm
  - Trực tuyến
  - Ngoại tuyến

Carpineto C., Osinski S., Romano G., Weiss D. (2009). A survey of web clustering engines, *ACM Comput. Surv.* , **41**(3), Article 17, 38 pages.

# Thuật toán K-mean gán cứng



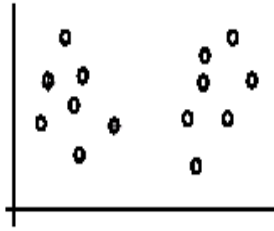
1. Khởi động: Chọn ngẫu nhiên  $k$  dữ liệu trong  $S$  làm trọng tâm (đại diện) cho các cụm  $S_i = \{c_i: c_i \in S\}, \forall i=1, \dots, k$
2. Bước lặp:
  - 2.1.  $S_i = \emptyset$  // Các cụm mới là rỗng
  - 2.2.  $\forall d \in S$ :
    - 2.2.1. Tính  $\text{sim}(d, c_i), \forall i=1, \dots, k$
    - 2.2.2.  $S_i = S_i \cup \{d\}$  nếu  $\text{sim}(d, c_i) = \max \{\text{sim}(d, c_i) | i=1, \dots, k\}$
  - 2.3.  $\forall i=1, \dots, k$ , tính lại trọng tâm các cụm  $S_i: c_i = \frac{1}{\|S_i\|} \sum_{d \in S_i} d$
3. Nếu chưa gặp Điều kiện dừng thì quay lại bước 2, ngược lại Dừng

## ● Một số lưu ý

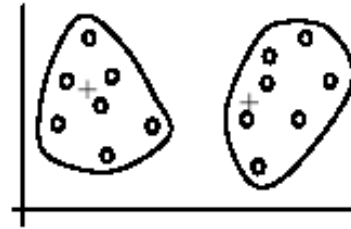
- Điều kiện dừng
  - Sau bước 2 không có sự thay đổi cụm
  - Điều kiện dừng cưỡng bức
    - ❖ Khống chế số lần lặp
    - ❖ Giá trị mục tiêu đủ nhỏ
- Vấn đề chọn tập đại diện ban đầu ở bước Khởi động
- Có thể dùng độ đo khoảng cách thay cho độ đo tương tự



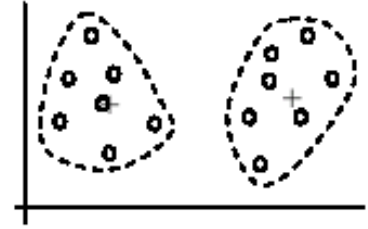
# Thuật toán K-mean gán cứng



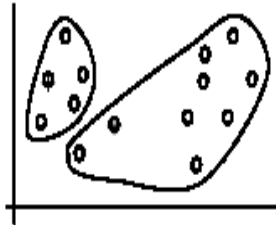
(A). Random selection of  $k$  centers



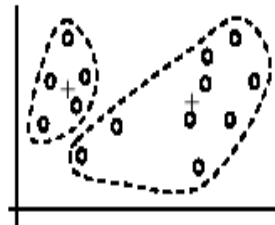
Iteration 2: (D). Cluster assignment



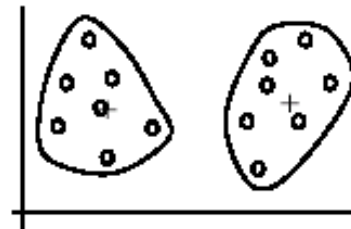
(E). Re-compute centroids



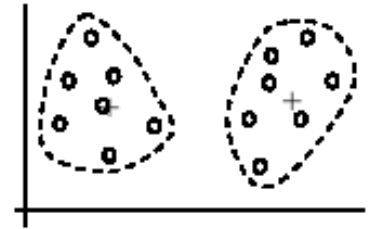
Iteration 1: (B). Cluster assignment



(C). Re-compute centroids



Iteration 3: (F). Cluster assignment



(G). Re-compute centroids

## ● Một số lưu ý (tiếp) và ví dụ

- ❑ Trong bước 2: các trọng tâm có thể không thuộc  $S$
- ❑ Thực tế: số lần lặp  $\approx 50$
- ❑ Thi hành k-mean với dữ liệu trên đĩa
  - Toàn bộ dữ liệu quá lớn: không thể ở bộ nhớ trong
  - Với mỗi vòng lặp: duyệt CSDL trên đĩa 1 lần
    - ❖ Tính được độ tương tự của  $d$  với các  $c_i$ .
    - ❖ Tính lại  $c_i$  mới: bước 2.1 khởi động (tổng, bộ đếm); bước 2.2 cộng và tăng bộ đếm; bước 2.3 chỉ thực hiện  $k$  phép chia.

# Thuật toán K-mean dạng mềm



- **Input**

- Số nguyên  $k > 0$ : số cụm biết trước
- Tập tài liệu  $D$  (cho trước)

- **Output**

- Tập  $k$  “đại diện cụm” ■ làm tối ưu lỗi “lượng tử”  $\sum_d \min_c |d - \mu_c|^2$

- **Định hướng**

- Tinh chỉnh ■ dần với tỷ lệ học ■ (learning rate)  $\mu_c \leftarrow \mu_c + \Delta\mu_c$

$$\Delta\mu_c = \sum_d \begin{cases} \eta (d - \mu_c) & \text{nếu } \mu_c \text{ gần } d \text{ nhất} \\ 0 & \text{các trường hợp khác} \end{cases}$$

$$\Delta\mu_c = \eta \frac{1/|d - \mu_c|^2}{\sum_\gamma 1/|d - \mu_\gamma|^2} (d - \mu_c) \quad \Delta\mu_c = \eta \frac{\exp(-|d - \mu_c|^2)}{\sum_\gamma \exp(-|d - \mu_\gamma|^2)} (d - \mu_c)$$

# Thuật toán K-mean



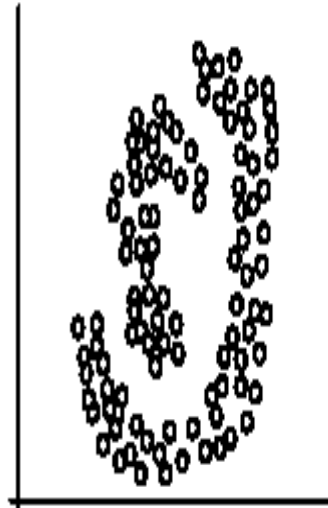
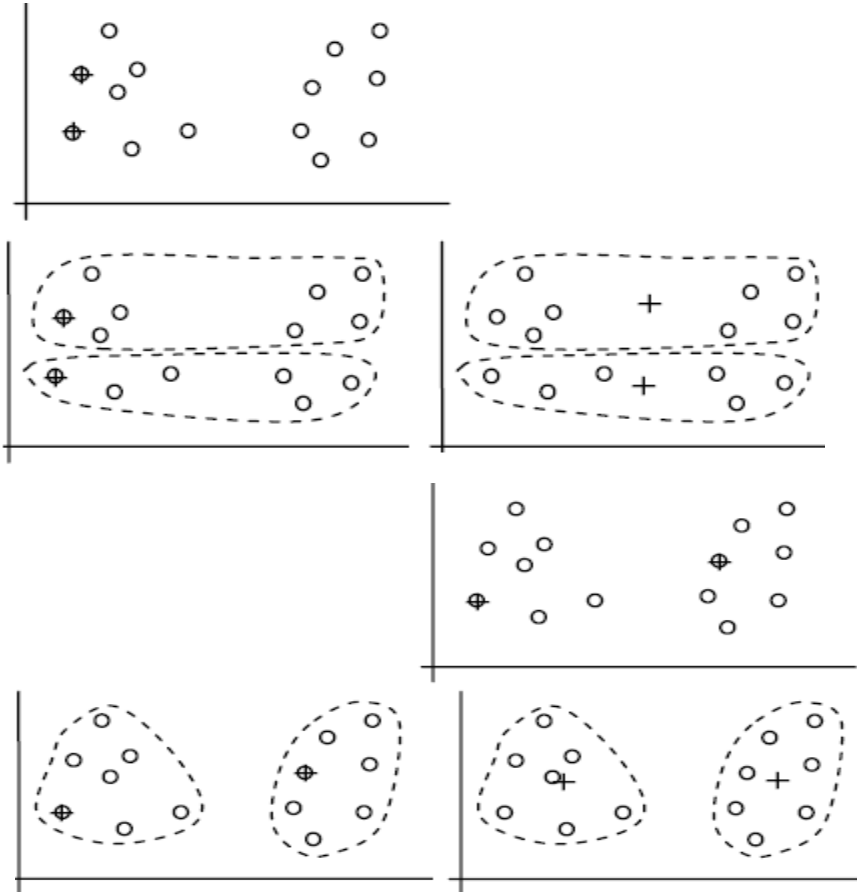
## ● Ưu điểm

- ❑ Đơn giản, dễ sử dụng
- ❑ Hiệu quả về thời gian: tuyến tính  $O(tkn)$ ,  $t$  số lần lặp,  $k$  số cụm,  $n$  là số phần tử
- ❑ Một thuật toán phân cụm phổ biến nhất
- ❑ Thường cho tối ưu cục bộ. Tối ưu toàn cục rất khó tìm

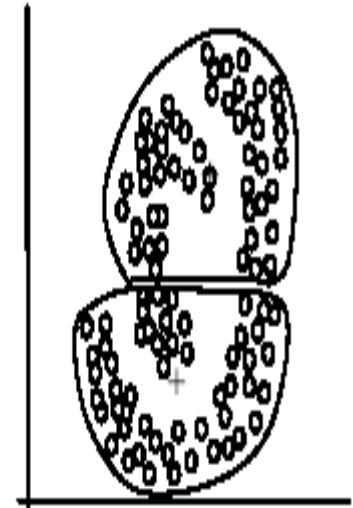
## ● Nhược điểm

- ❑ Phải “tính trung bình được”: dữ liệu phân lớp thì dựa theo tần số
- ❑ Cần cho trước  $k$  : số cụm
- ❑ Nhạy cảm với ngoại lệ (cách xa so với đại đa số dữ liệu còn lại): ngoại lệ thực tế, ngoại lệ do quan sát sai (làm sạch dữ liệu)
- ❑ Nhạy cảm với mẫu ban đầu: cần phương pháp chọn mẫu thô tốt
- ❑ Không thích hợp với các tập dữ liệu không siêu-ellip hoặc siêu cầu (các thành phần con không ellip/cầu hóa)

# Thuật toán K-mean



(A): Two natural clusters



(B):  $k$ -means clusters

Trái: Nhạy cảm với chọn mẫu ban đầu

Phải: Không thích hợp với bộ dữ liệu không siêu ellip/cầu hóa

# 3. Phân cụm phân cấp từ dưới lên



- **HAC:** Hierarchical agglomerative clustering
- **Một số độ đo phân biệt cụm**
  - Độ tương tự hai tài liệu
  - Độ tương tự giữa hai cụm
    - Độ tương tự giữa hai đại diện
    - Độ tương tự cực đại giữa hai tài liệu thuộc hai cụm: **single-link**
    - Độ tương tự cực tiểu giữa hai tài liệu thuộc hai cụm: **complete-link**
    - Độ tương tự trung bình giữa hai tài liệu thuộc hai cụm
- **Sơ bộ về thuật toán**
  - Đặc điểm: Không cho trước số lượng cụm  $k$ , cho phép đưa ra các phương án phân cụm theo các giá trị  $k$  khác nhau
  - Lưu ý:  $k$  là một tham số  $\Rightarrow$  “tìm  $k$  tốt nhất”
  - Tinh chỉnh: Từ cụ thể tới khái quát

# Phân cụm phân cấp từ dưới lên

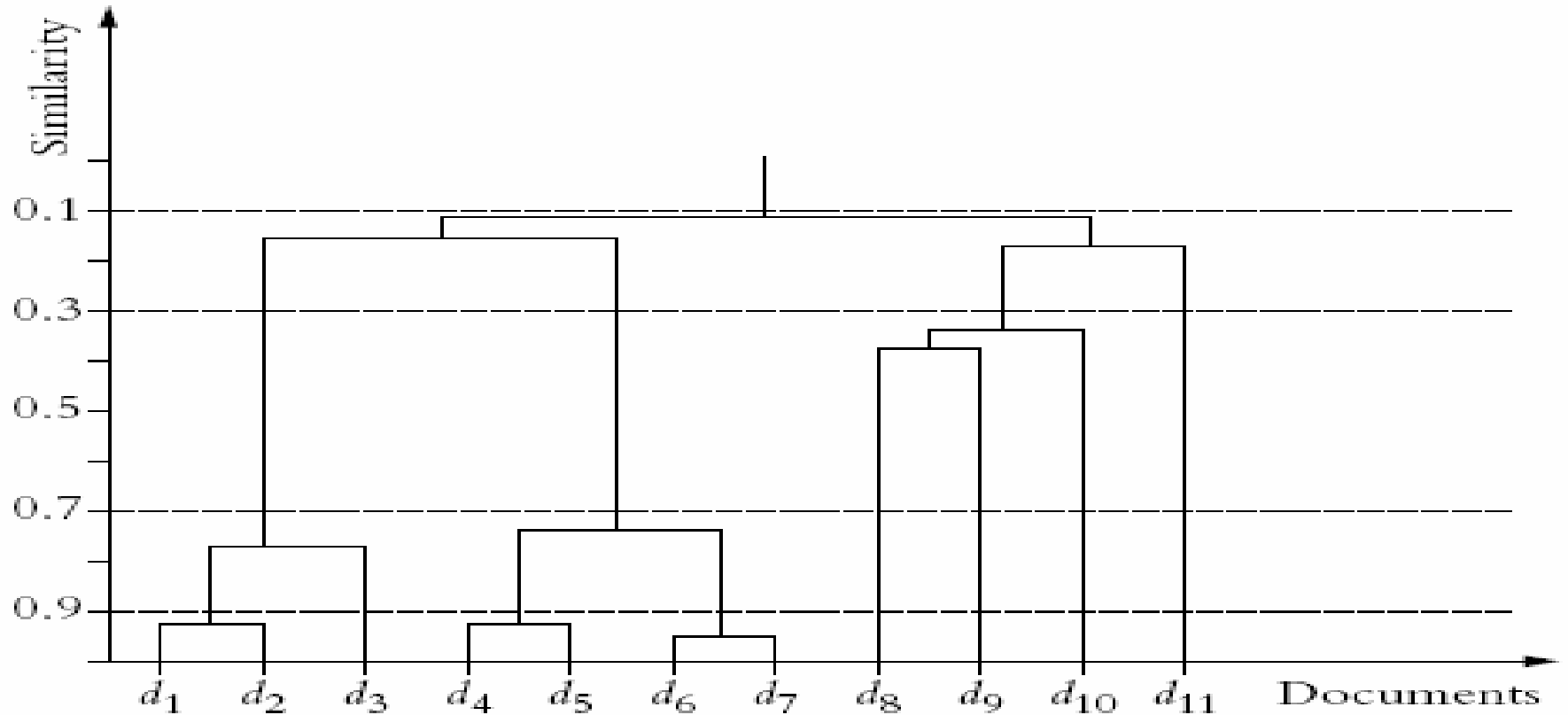


1.  $G \leftarrow \{ \{d\} \mid d \in S \}$  (khởi tạo  $G$  là tập các cụm chỉ gồm một trang web trong tập  $S$ ).
2. Nếu  $|G| < k$  thì dừng thuật toán (đã đạt được số lượng cụm mong muốn).
3. Tìm hai cụm  $S_i, S_j \in G$  sao cho  $(i, j) = \arg \max_{(i, j)} \text{sim}(S_i, S_j)$  (tìm hai cụm có độ tương tự lớn nhất).
4. Nếu  $\text{sim}(S_i, S_j) < q$  thì dừng thuật toán (độ tương tự của 2 cụm nhỏ hơn ngưỡng cho phép).
5. Loại bỏ  $S_i, S_j$  khỏi  $G$ .
6.  $G = G \cup \{S_i, S_j\}$  (ghép hai cụm  $S_i, S_j$  và đưa vào trong tập  $G$ ).
7. Nhảy đến bước 2.

## ● Giải thích

- $G$  là tập các cụm trong phân cụm
- Điều kiện  $|G| < k$  có thể thay thế bằng  $|G|=1$

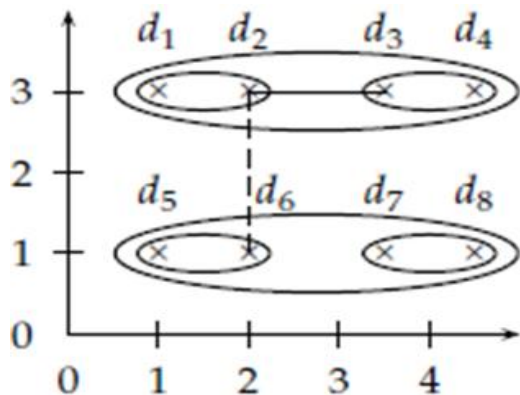
# Phân cụm phân cấp từ dưới lên



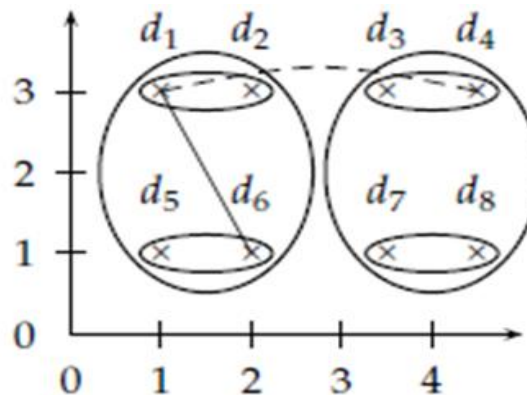
## ● Hoạt động HAC

- Cho phép với mọi  $k$
- Chọn phân cụm theo “ngưỡng” về độ tương tự

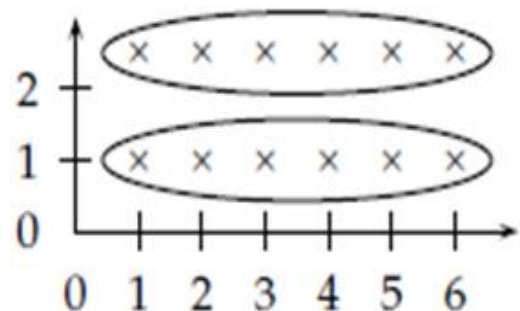
# HAC với các độ đo khác nhau



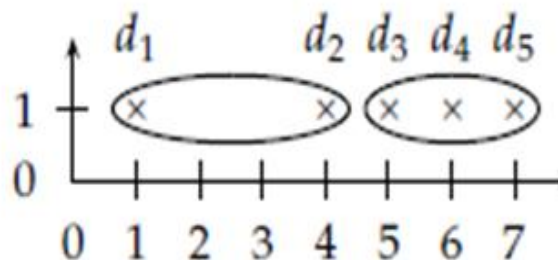
(a) single-link



(b) complete-link



(a) single-link



(b) complete-link

- Ảnh hưởng của các độ đo

- Trên: Hoạt động thuật toán khác nhau theo các độ đo khác nhau: độ tương tự cực tiểu (complete-link) có tính cầu hơn so với cực đại
- Dưới: Độ tương tự cực đại (Single-link) tạo cụm chuỗi dòng



# 4. Biểu diễn cụm và gán nhãn



- Các phương pháp biểu diễn điển hình
  - Theo đại diện cụm
    - Đại diện cụm làm tâm
    - Tính bán kính và độ lệch chuẩn để xác định phạm vi của cụm
    - Cụm không ellip/cầu hóa: không tốt
  - Theo mô hình phân lớp
    - Chỉ số cụm như nhãn lớp
    - Chạy thuật toán phân lớp để tìm ra biểu diễn cụm
  - Theo mô hình tần số
    - Dùng cho dữ liệu phân loại
    - Tần số xuất hiện các giá trị đặc trưng cho từng cụm
- Lưu ý
  - Dữ liệu phân cụm ellip/cầu hóa: đại diện cụm cho biểu diễn tốt
  - Cụm hình dạng bất thường rất khó biểu diễn

# Gán nhãn cụm tài liệu



- Phân biệt các cụm (MU)

$$I(U; C) = \frac{N_{11}}{N} \log_2 \frac{NN_{11}}{N_1 N_1} + \frac{N_{01}}{N} \log_2 \frac{NN_{01}}{N_0 N_1} \\ + \frac{N_{10}}{N} \log_2 \frac{NN_{10}}{N_1 N_0} + \frac{N_{00}}{N} \log_2 \frac{NN_{00}}{N_0 N_0}$$

- Chọn từ khóa đặc trưng tương quan cụm
- $N_{xy}$  (x có từ khóa t, y tài liệu thuộc C)
  - $N_{11}$  : số tài liệu chứa t thuộc cụm C
  - $N_{10}$  : số tài liệu chứa t không thuộc cụm C
  - $N_{01}$  : số tài liệu không chứa t thuộc cụm C
  - $N_{00}$  : số tài liệu không chứa t không thuộc cụm C
  - N: Tổng số tài liệu
- Hướng “trọng tâm” cụm
  - Dùng các từ khóa tần số cao tại trọng tâm cụm
- Tiêu đề
  - Chọn tiêu đề của tài liệu trong cụm gần trọng tâm nhất

# Gán nhãn cụm tài liệu



	# docs	labeling method		
		centroid	mutual information	title
4	622	oil plant mexico production crude power 000 refinery gas bpd	plant oil production barrels crude bpd mexico dolly capacity petroleum	MEXICO: Hurricane Dolly heads for Mexico coast
9	1017	police security russian people military peace killed told grozny court	police killed military security peace told troops forces rebels people	RUSSIA: Russia's Lebed meets rebel chief in Chechnya
10	1259	00 000 tonnes traders futures wheat prices cents september tonne	delivery traders futures tonne tonnes desk wheat prices 000 00	USA: Export Business - Grain/oilseeds complex

## ● Ví dụ

- Ba phương pháp chọn nhãn cụm đối với 3 cụm là cụm 4 (622 tài liệu), cụm 9 (1017 tài liệu), cụm 10 (1259 tài liệu) khi phân cụm 10000 tài liệu đầu tiên của bộ Reuters-RCV1
- centroid: các từ khóa có tần số cao nhất trong trọng tâm; mutual information (MU): thông tin liên quan phân biệt các cụm; title: tiêu đề tài liệu gần trọng tâm nhất.

# 5. Đánh giá phân cụm



- **Đánh giá chất lượng phân cụm là khó khăn**

- Chưa biết các cụm thực sự

- **Một số phương pháp điển hình**

- Người dùng kiểm tra

- Nghiên cứu trọng tâm và miền phủ
- Luật từ cây quyết định
- Đọc các dữ liệu trong cụm

- Đánh giá theo các độ đo tương tự/khoảng cách

- Độ phân biệt giữa các cụm
- Phân ly theo trọng tâm

- Dùng thuật toán phân lớp

- Coi mỗi cụm là một lớp
- Học bộ phân lớp đa lớp (cụm)
- Xây dựng ma trận nhầm lẫn khi phân lớp

- Tính các độ đo: entropy, tinh khiết, chính xác, hồi tưởng, độ đo F và đánh giá theo các độ đo này

# Đánh giá theo độ đo tương tự



## ● Độ phân biệt các cụm

- ❑ Cực đại hóa tổng độ tương tự nội tại của các cụm
- ❑ Cực tiểu hóa tổng độ tương tự các cặp cụm khác nhau
- ❑ Lấy độ tương tự cực tiểu (complete link), cực đại (single link)

$$J_e = \frac{1}{2} \sum_{i=1}^k \frac{1}{|S_i|} \sum_{d_j, d_l \in S_i} \|d_j - d_l\|^2$$

$$J_s = \frac{1}{2} \sum_{i=1}^k \frac{1}{|S_i|} \sum_{d_j, d_l \in S_i} \text{sim}(d_j, d_l) = \frac{1}{2} \sum_{i=1}^k |S_i| \text{sim}(S_i)$$

## ● Một số phương pháp điển hình

- ❑ Phân lý theo trọng tâm

$$J_e = \sum_{i=1}^k \sum_{d \in S_i} \|d - c_i\|^2$$

# Ví dụ



Bảng 7.2 Dữ liệu mẫu dành cho phân cụm phẳng

Tên trang web	A1	A2	A3	A4	A5	A6
Anthropology	0	0.537	0.477	0	0.673	0.177
Art	0	0	0	0.961	0.195	0.196
Biology	0	0.347	0.924	0	0.111	0.112
Chemistry	0	0.975	0	0	0.155	0.158
Communication	0	0	0	0.78	0.626	0
Computer Science	0	0.989	0	0	0.13	0.067
Criminal Justice	0	0	0	0	1	0
Economics	0	0	1	0	0	0
English	0	0	0	0.98	0	0.199
Geography	0	0.849	0	0	0.528	0
History	0.991	0	0	0.135	0	0
Mathematics	0	0.616	0.549	0.49	0.198	0.201
Modern Languages	0	0	0	0.928	0	0.373
Music	0.97	0	0	0	0.17	0.172
Philosophy	0.741	0	0	0.658	0	0.136
Physics	0	0	0.894	0	0.315	0.318
Political Science	0	0.933	0.348	0	0.062	0.063
Psychology	0	0	0.852	0.387	0.313	0.162
Sociology	0	0	0.639	0.57	0.459	0.237
Theatre	0	0	0	0	0.967	0.254

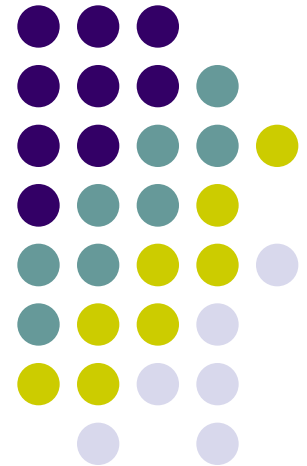
Bảng 7.6 Giá trị của hàm đánh giá dựa trên độ đo tương tự với giải thuật k-means

$k=2$	$k=3$	$k=4$
<p><b>1 [8.53381]</b></p> <p>Anthropology, Biology, Chemistry, Computer Science, Economics, Geography, Mathematics, Physics, Political Science, Psychology, Sociology</p> <p><b>2 [6.12743]</b></p> <p>Art, Communication, Criminal Justice, English, History, Modern Languages, Music, Philosophy, Theatre</p> <p><math>\Sigma = [12.0253]</math></p>	<p><b>1 [2.83806]</b></p> <p>History, Music, Philosophy</p> <p><b>2 [6.09107]</b></p> <p>Anthropology, Biology, Chemistry, Computer Science, Geography, Mathematics, Political Science</p> <p><b>3 [7.12119]</b></p> <p>Art, Communication, Criminal Justice, Economics, English, Modern Languages, Physics, Psychology, Sociology, Theatre</p> <p><math>\Sigma = [12.0253]</math></p>	<p><b>1 [3.81771]</b></p> <p>Art, Communication, English, Modern Languages</p> <p><b>2 [5.44416]</b></p> <p>Biology, Economics, Mathematics, Physics, Psychology, Sociology</p> <p><b>3 [2.83806]</b></p> <p>History, Music, Philosophy</p> <p><b>4 [5.64819]</b></p> <p>Anthropology, Chemistry, Computer Science, Criminal Justice, Geography, Political Science, Theatre</p> <p><math>\Sigma = [12.0253]</math></p>

# BÀI GIẢNG KHAI PHÁ DỮ LIỆU WEB

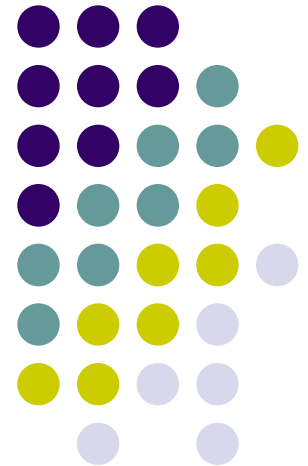
## CHƯƠNG 7. PHÂN LỚP WEB

PGS. TS. HÀ QUANG THỤY  
HÀ NỘI 10-2010  
**TRƯỜNG ĐẠI HỌC CÔNG NGHỆ**  
ĐẠI HỌC QUỐC GIA HÀ NỘI



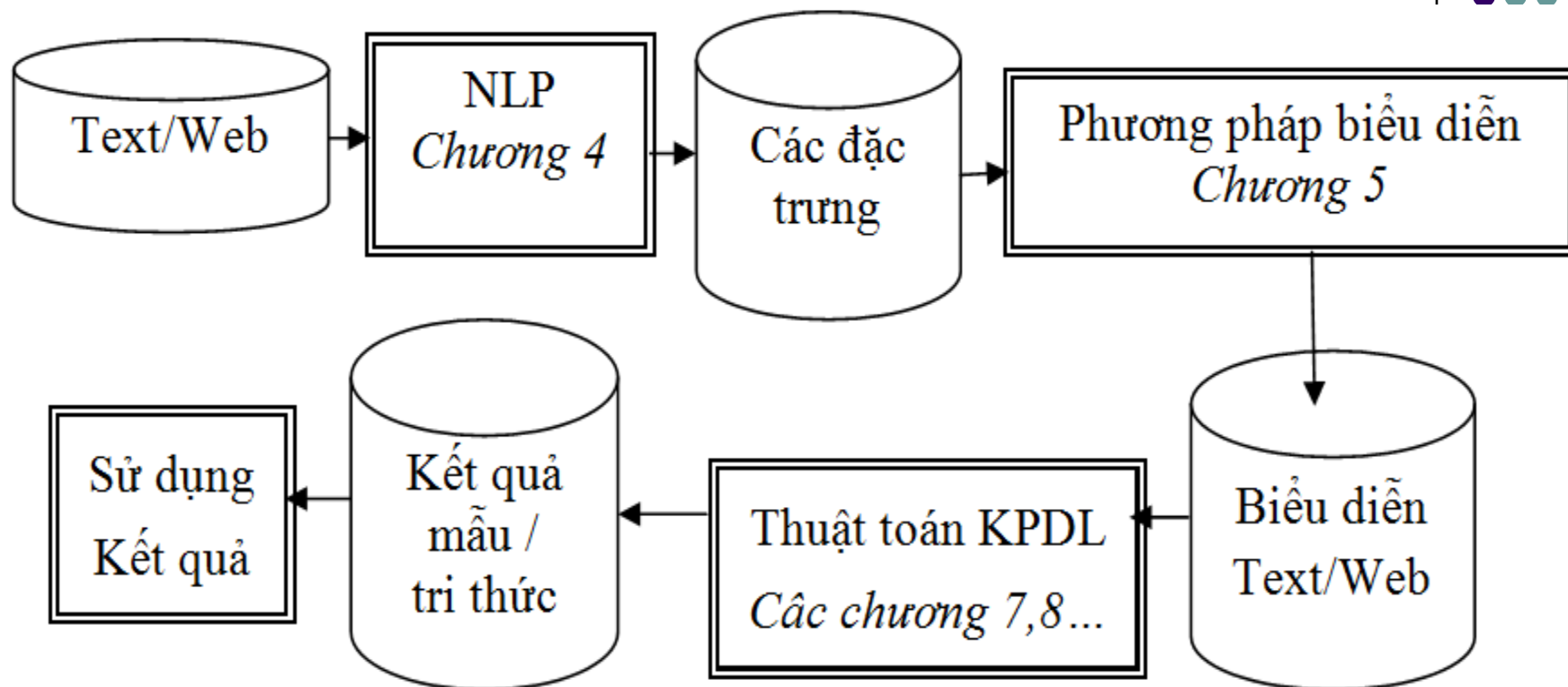
# Nội dung

Giới thiệu phân lớp Web  
Phân lớp học giám sát  
Phân lớp học bán giám sát





# Giới thiệu: Sơ đồ khai phá Web



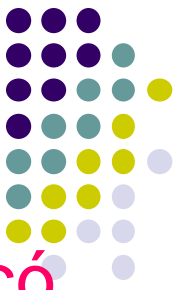
- Thuật toán KPD: phân lớp, phân cụm, tóm tắt... Sử dụng các thuật toán KPD chung (phân lớp, phân cụm...)
- Chọn các đặc trưng, chọn cách biểu diễn Web đóng vai trò quan trọng trong KPD Web: Chương 4 và Chương 5.
- Các chương: phát biểu bài toán và một số thuật toán KPD điển hình

# Bài toán phân lớp Web



- Đầu vào
  - Tập tài liệu web  $D = \{d_i\}$
  - Tập các lớp  $C_1, C_2, \dots, C_k$  mỗi tài liệu  $d$  thuộc một lớp  $C_i$
  - Tập ví dụ  $D_{\text{exam}} = D_1 + D_2 + \dots + D_k$  với  $D_i = \{d \in D_{\text{exam}} : d \text{ thuộc } C_i\}$
  - Tập ví dụ  $D_{\text{exam}}$  đại diện cho tập  $D$
- Đầu ra
  - Mô hình phân lớp: ánh xạ từ  $D$  sang  $C$
- Sử dụng mô hình
  - $d \in D \setminus D_{\text{exam}}$  : xác định lớp của tài liệu  $d$
- Ví dụ
  - Crawler hướng chủ đề: Chủ đề  $C_i$
  - Phân lớp/phân cụm tập trang Web trả về “chủ đề/lớp”

# Phân lớp: Quá trình hai pha



## ● Xây dựng mô hình: Tìm mô tả cho tập lớp đã có

- Cho trước tập lớp  $C = \{C_1, C_2, \dots, C_k\}$
- Cho ánh xạ (chưa biết) từ miền  $D$  sang tập lớp  $C$
- Có tập ví dụ  $D_{\text{exam}} = D_1 + D_2 + \dots + D_k$  với  $D_i = \{d_{\text{exam}}^i : d_i\}$   
 $D_{\text{exam}}$  được gọi là tập ví dụ mẫu.
- Xây dựng ánh xạ (mô hình) phân lớp trên: Dạy bộ phân lớp.
- Mô hình: Luật phân lớp, cây quyết định, công thức toán học...

## ● Pha 1: Dạy bộ phân lớp

- Tách  $D_{\text{exam}}$  thành  $D_{\text{train}}$  (2/3) +  $D_{\text{test}}$  (1/3).  $D_{\text{train}}$  và  $D_{\text{test}}$  “tính đại diện” cho miền ứng dụng
- $D_{\text{train}}$  : xây dựng mô hình phân lớp (xác định tham số mô hình)
- $D_{\text{test}}$  : đánh giá mô hình phân lớp (các độ đo hiệu quả)
- Chọn mô hình có chất lượng nhất

## ● Pha 2: Sử dụng bộ phân lớp

- $d \in D \setminus D_{\text{exam}}$  : xác định lớp của  $d$ .

# Ví dụ phân lớp: Bài toán cho vay



<i>Tid</i>	<i>Refund</i>	<i>Marital Status</i>	<i>Taxable Income</i>	<i>Cheat</i>
1	No	Single	75K	No
2	Yes	Married	50K	No
3	No	Single	75K	No
4	No	Married	150K	Yes
5	No	Single	40K	No
6	No	Married	80K	Yes
7	No	Single	75K	No
8	Yes	Married	50K	No
9	Yes	Married	50K	No
10	No	Married	150K	Yes
11	No	Single	40K	No
12	No	Married	150K	Yes
13	No	Married	80K	Yes
14	No	Single	40K	No
15	No	Married	80K	Yes

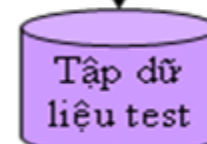
B

# Phân lớp: Quá trình hai pha

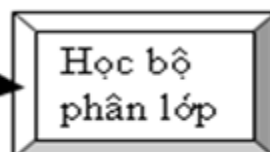


Tid	Refund	Marital Status	Taxable Income	Cheat
1	No	Single	75K	No
2	Yes	Married	50K	No
4	No	Married	150K	Yes
5	No	Single	40K	No
6	No	Married	80K	Yes
7	No	Single	75K	No
9	Yes	Married	50K	No
10	No	Married	150K	Yes
11	No	Single	40K	No
13	No	Married	80K	Yes

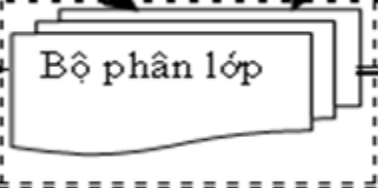
Tid	Refund	Marital Status	Taxable Income	Cheat
3	No	Single	75K	No
8	Yes	Married	50K	No
12	No	Married	150K	Yes
14	No	Single	40K	No
15	No	Married	80K	Yes



*Pha 1. Học bộ phân lớp*



Refund	Marital Status	Taxable Income	Cheat
No	Married	75K	?



Refund	Marital Status	Taxable Income	Cheat
No	Married	75K	Y/N

*Pha 2. Sử dụng bộ phân lớp*

# Phân lớp: Quá trình hai pha

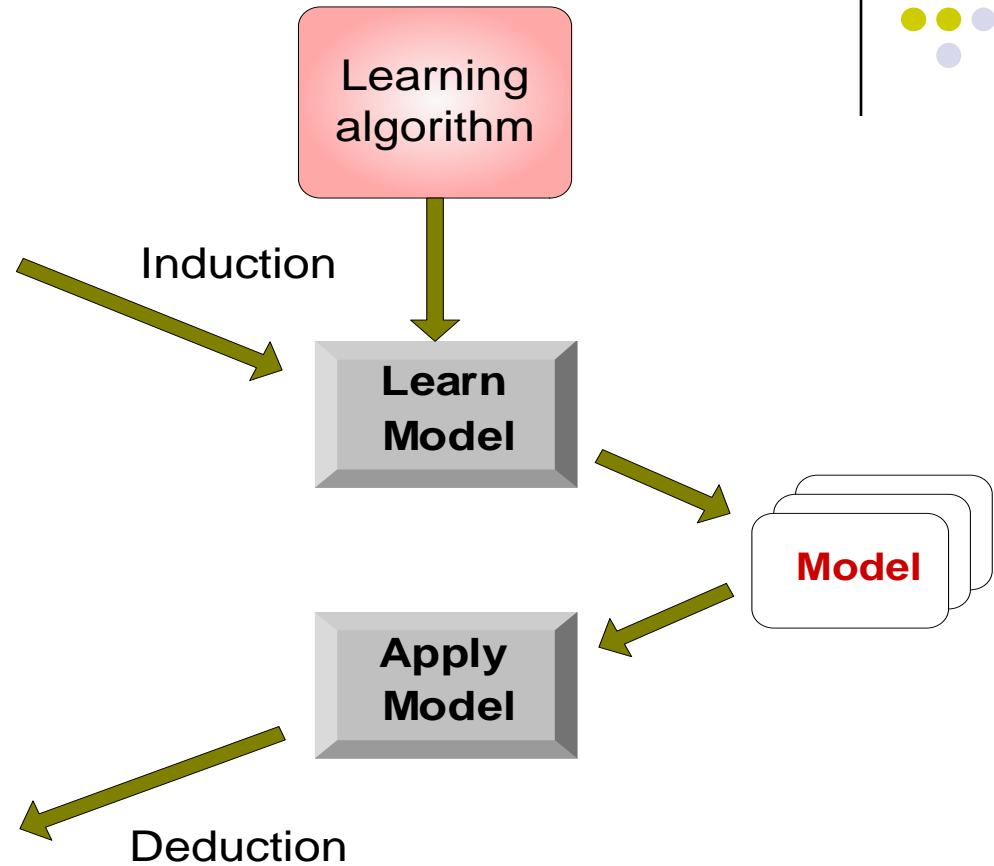


Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



# Các loại phân lớp



- Phân lớp nhị phân/ đa lớp:
  - $|C|=2$ : phân lớp nhị phân.
  - $|C|>2$ : phân lớp đa lớp.
- Phân lớp đơn nhãn/ đa nhãn:
  - Đơn nhãn: mỗi tài liệu được gán vào chính xác một lớp.
  - Đa nhãn: một tài liệu có thể được gán nhiều hơn một lớp.
  - Phân cấp: lớp này là cha/con của lớp kia

# Các vấn đề đánh giá mô hình



- Các phương pháp đánh giá hiệu quả

Câu hỏi: Làm thế nào để đánh giá được hiệu quả của một mô hình?

- Độ đo để đánh giá hiệu quả

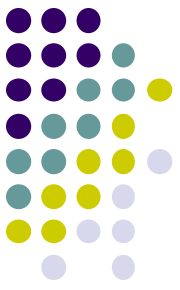
Câu hỏi: Làm thế nào để có được ước tính đáng tin cậy?

- Phương pháp so sánh mô hình

Câu hỏi: Làm thế nào để so sánh hiệu quả tương đối giữa các mô hình có tính cạnh tranh?



# Đánh giá phân lớp nhị phân



- Theo dữ liệu test
- Giá trị thực: P dương / N âm; Giá trị qua phân lớp: T đúng/F sai. : còn gọi là *ma trận nhầm lẫn*
- Sử dụng các ký hiệu TP (true positives), TN (true negatives), FP (false positives), FN (false negatives)
  - TP: số ví dụ dương P mà thuật toán phân lớp cho giá trị đúng T
  - TN: số ví dụ âm N mà thuật toán phân lớp cho giá trị đúng T
  - FP: số ví dụ dương P mà thuật toán phân lớp cho giá trị sai F
  - FN: số ví dụ âm N mà thuật toán phân lớp cho giá trị sai F
- Độ hồi tưởng  $\blacksquare$  độ chính xác  $\blacksquare$  các độ đo  $F_1$  và  $F_\beta$

$$\blacksquare = \frac{TP}{TP + FP}$$

$$\blacksquare = \frac{TP}{TP + FN}$$

$$f_\beta = \frac{(\beta^2 + 1)\pi\rho}{\beta^2\pi + \rho}$$

$$f_1 = \frac{2\pi\rho}{\pi + \rho}$$

# Đánh giá phân lớp nhị phân



- Phương án khác đánh giá mô hình nhị phân theo độ chính xác (accuracy) và hệ số lỗi (Error rate)
- *Ma trận nhầm lẫn*

		Lớp dự báo	
		Lớp = 1	Lớp = 0
Lớp thực sự	Lớp = 1	$f_{11}$	$f_{10}$
	Lớp = 0	$f_{01}$	$f_{00}$

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}} = \frac{f_{11} + f_{00}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

$$\text{Error rate} = \frac{\text{Number of wrong predictions}}{\text{Total number of predictions}} = \frac{f_{10} + f_{01}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

# So sánh hai phương án



- Tập test có 9990 ví dụ lớp 0 và 10 ví dụ lớp 1. Kiểm thử: mô hình dự đoán cả 9999 ví dụ là lớp 0 và 1 ví dụ lớp 1 (chính xác: TP)
  - Theo phương án (precision, recall) có
    - $1/10=0.1$ ; ■  $1/1=1$ ;  $f_1 = 2*0.1/(0.1+1.0)= 0.18$
  - Theo phương án (accuracy, error rate) có
    - accuracy=0.9991; error rate =  $9/10000 = 0.0009$
    - Được coi là rất chính xác !
  - $f_1$  thể hiện việc đánh giá nhạy cảm với giá dữ liệu

# Đánh giá phân lớp đa lớp



- Bài toán ban đầu:  $C$  gồm có  $k$  lớp
- Đối với mỗi lớp  $C_i$ , cho thực hiện thuật toán với các dữ liệu thuộc  $D_{\text{test}}$  nhận được các đại lượng  $TP_i$ ,  $TF_i$ ,  $FP_i$ ,  $FN_i$  (như bảng dưới đây)

Lớp $C_i$		Giá trị thực	
		Thuộc lớp $C_i$	Không thuộc lớp $C_i$
Giá trị qua bộ phân lớp đa lớp	Thuộc lớp $C_i$	$TP_i$	$TN_i$
	Không thuộc lớp $C_i$	$FP_i$	$FN_i$

# Đánh giá phân lớp đa lớp



- Tương tự bộ phân lớp hai lớp (nhị phân)
  - Độ chính xác  $Pr_i$  của lớp  $C_i$  là tỷ lệ số ví dụ dương được thuật toán phân lớp cho giá trị đúng trên tổng số ví dụ được thuật toán phân lớp vào lớp  $C_i$ :

$$Pr_i = \frac{TP_i}{TP_i + FN_i}$$

- Độ hồi tưởng  $Re_i$  của lớp  $C_i$  là tỷ lệ số ví dụ dương được thuật toán phân lớp cho giá trị đúng trên tổng số ví dụ dương thực sự thuộc lớp  $C_i$ :

$$Re_i = \frac{TP_i}{TP_i + FP_i}$$



# Đánh giá phân lớp đa lớp

- Các giá trị  $\bar{K}_c$  và  $\bar{K}_c^*$ : độ hồi phục và độ chính xác đối với lớp  $C_i$ .
- Đánh giá theo các độ đo
  - vi trung bình-microaveraging (được ưa chuộng)  $\bar{K}_c$  và  $\bar{K}_c^*$
  - trung bình lớn-macroaveraging  $\bar{K}_c$  và  $\bar{K}_c^*$

$$\bar{K}_c = \frac{1}{K_c} \sum_{i=1}^K \frac{TP_c}{(TP_c + FP_c)}$$

$$\bar{K}_c^* = \frac{1}{K_c} \sum_{i=1}^K \frac{TP_c}{(TP_c + FN_c)}$$

# Các kỹ thuật phân lớp



- Các phương pháp cây quyết định  
Decision Tree based Methods
- Các phương pháp dựa trên luật  
Rule-based Methods
- Các phương pháp Bayes «ngây thơ» và mạng tin cậy Bayes  
Naïve Bayes and Bayesian Belief Networks
- Các phương pháp máy vector hỗ trợ  
Support Vector Machines
- Lập luận dựa trên ghi nhớ  
Memory based reasoning
- Các phương pháp mạng nơon  
Neural Networks
- Một số phương pháp khác

# Phân lớp cây quyết định



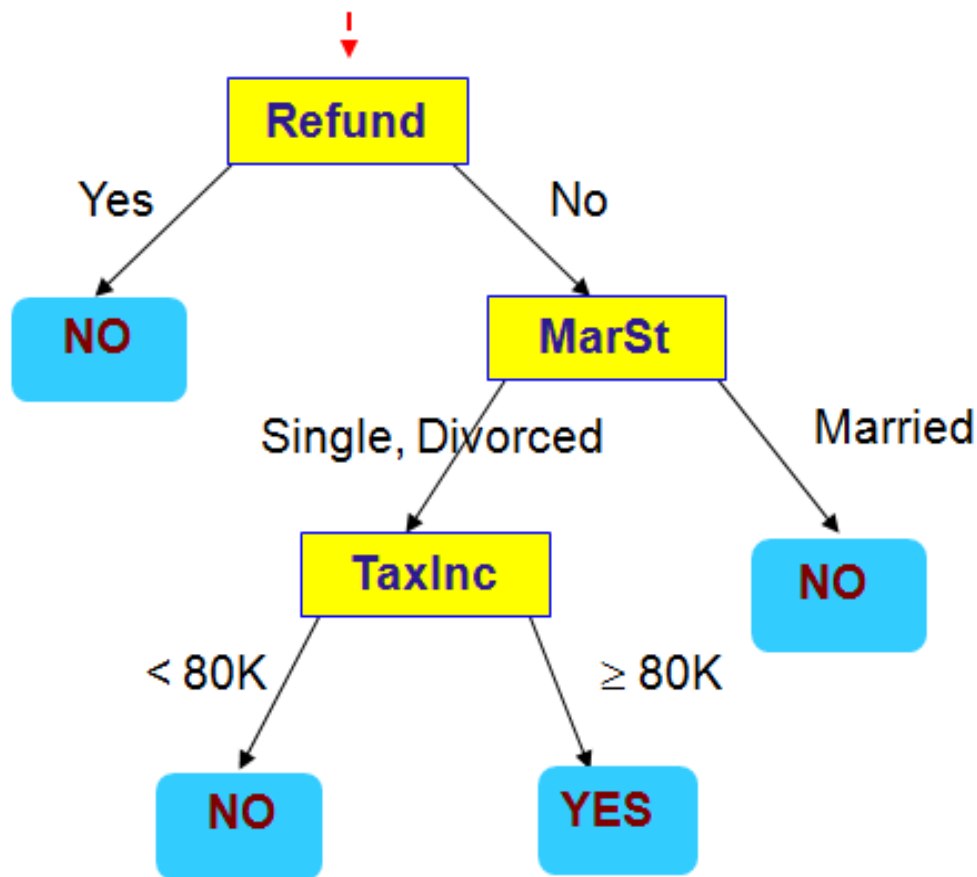
- Mô hình phân lớp là cây quyết định
- Cây quyết định
  - Gốc: **tên thuộc tính**; không có cung vào + không/một số cung ra
  - Nút trong: **tên thuộc tính**; có chính xác một cung vào và một số cung ra (gắn với điều kiện kiểm tra giá trị thuộc tính của nút)
  - Lá hoặc nút kết thúc: **giá trị lớp**; có chính xác một cung vào + không có cung ra.
  - Ví dụ: xem trang tiếp theo
- Xây dựng cây quyết định
  - Phương châm: “chia để trị”, “chia nhỏ và chế ngự”. Mỗi nút tương ứng với một tập các ví dụ học. **Gốc: toàn bộ dữ liệu học**
  - Một số thuật toán phổ biến: Hunt, họ ID3+C4.5+C5.x
- Sử dụng cây quyết định
  - Kiểm tra từ gốc theo các điều kiện



# Ví dụ cây quyết định và sử dụng



Bắt đầu từ gốc của cây



## Test Data

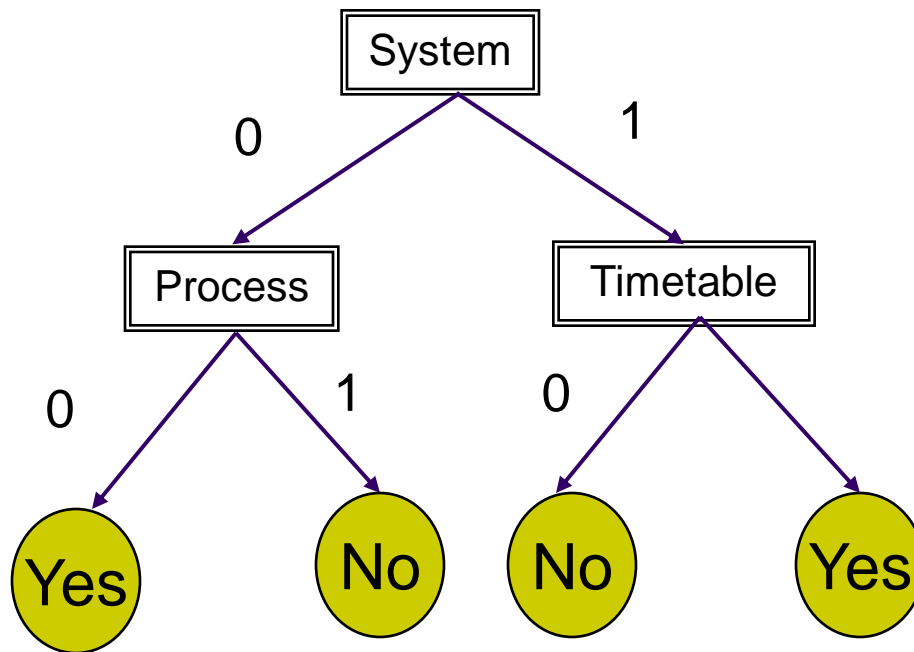
Refund	Marital Status	Taxable Income	Cheat
No	Married	80K	?

Kết luận: Gán giá trị **YES** vào trường **Cheat** cho bản ghi



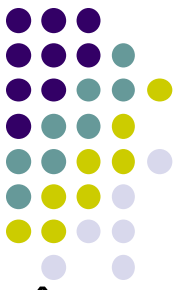
# Ví dụ cây quyết định phân lớp văn bản

- Phân lớp văn bản vào lớp AI : trí tuệ nhân tạo
- Dựa vào các từ khóa có trong văn bản: System, Process, Timetable (Phân tích miền ứng dụng)



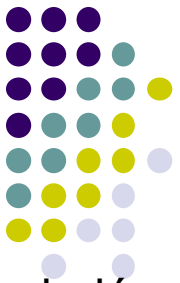
1. **If System=0 and Process=0 then Class AI = Yes.**
2. **If System=0 and Process=1 then Class AI = No.**
3. **If System=1 and Timetable=1 then Class AI = Yes.**
4. **If System=1 and Timetable=0 then Class AI = No.**

# Dựng cây quyết định: thuật toán Hunt



- Thuật toán dựng cây quyết định sớm nhất, đệ quy theo nút của cây, bắt đầu từ gốc
- **Input**
  - Cho nút  $t$  trên cây quyết định đang được xem xét
  - Cho tập các ví dụ học  $D_t$ .
  - Cho tập nhãn lớp (giá trị lớp)  $y_1, y_1, \dots, y_k$ . ( $k$  lớp)
- **Output**
  - Xác định nhãn nút  $t$  và các cung ra (nếu có) của  $t$
- **Nội dung**
  - 1: Nếu mọi ví dụ trong  $D_t$  đều thuộc vào một lớp  $y$  thì nút  $t$  là một lá và được gán nhãn  $y$ .
  - 2: Nếu  $D_t$  chứa các ví dụ thuộc nhiều lớp thì
    - 2.1. **Chọn 1 thuộc tính  $A$**  để phân hoạch  $D_t$  và gán nhãn nút  $t$  là  $A$
    - 2.2. Tạo phân hoạch  $D_t$  theo tập giá trị của  $A$  thành các tập con
    - 2.3. Mỗi tập con theo phân hoạch của  $D_t$  tương ứng với một nút con  $u$  của  $t$ : cung nối  $t$  tới  $u$  là miền giá trị  $A$  theo phân hoạch, tập con nói trên được xem xét với  $u$  tiếp theo. Thực hiện thuật toán với từng nút con  $u$  của  $t$ .

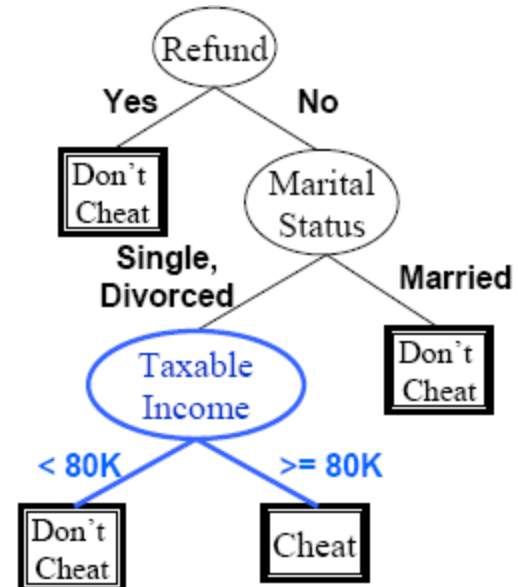
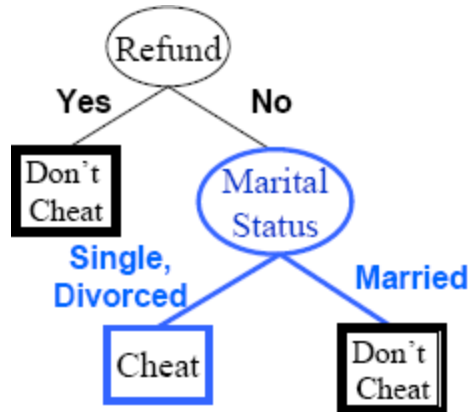
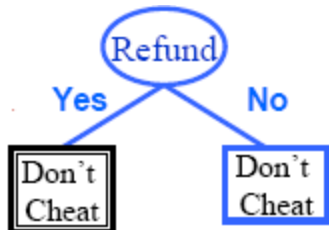
# Ví dụ: thuật toán Hunt



Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

## Giải thích

- Xuất phát từ gốc với 10 bản ghi
- Thực hiện bước 2: **chọn thuộc tính Refund** có hai giá trị Yes, No. Chia thành hai tập gồm 3 bản ghi có Refund = Yes và 7 bản ghi có Refund = No
- Xét hai nút con của gốc từ trái sang phải. **Nút trái** có 3 bản ghi cùng thuộc lớp Cheat=No (Bước 1) nên là lá gán **No (Don't cheat)**. **Nút phải** có 7 bản ghi có cả No và Yes nên áp dụng bước 2. **Chọn thuộc tính Marital Status** với phân hoạch Married và hai giá trị kia...



# Thuật toán cây quyết định ID3



ID3 (*Examples*, *Target\_attribute*, *Attributes*)

Ở đây: *Examples* là tập ví dụ học; *Target\_attribute* là các thuộc tính đầu ra (lớp) cho cây quyết định dự đoán; *Attributes* là danh sách các thuộc tính khác tham gia trong quá trình học của cây quyết định. Kết quả thủ tục trả về cây quyết định phân lớp đúng các mẫu ví dụ đưa ra.

1. Tạo một nút gốc *Root* cho cây quyết định.
2. Nếu toàn bộ *Examples* đều là các ví dụ thuộc cùng một lớp thì trả lại cây *Root* một nút đơn với nhãn + (nếu các ví dụ thuộc lớp +) hoặc với nhãn - (nếu các ví dụ thuộc lớp -).
3. Nếu *Attributes* là rỗng thì trả lại cây *Root* một nút đơn với nhãn gán bằng giá trị phổ biến nhất của *Target\_attribute* trong *Examples*.
4. Còn lại

Begin

4.1. Gán  $A \leftarrow$  thuộc tính từ tập *Attributes* mà phân lớp tốt nhất tập *Examples*.

4.2. Thuộc tính quyết định cho  $Root \leftarrow A$

4.3. Lặp với các giá trị có thể  $v_i$  của  $A$ ,

- Cộng thêm một nhánh cây con ở dưới *Root*, phù hợp với biểu thức kiểm tra  $A = v_i$ .

- Đặt  $Examples_{v_i}$  là một tập con của tập các ví dụ có giá trị  $v_i$  cho  $A$

- Nếu  $Examples_{v_i}$  rỗng

+ Thì dưới mỗi nhánh mới thêm một nút lá với nhãn = giá trị phổ biến nhất của *Target\_attribute* trong tập *Examples*.

+ Ngược lại thì dưới nhánh mới này thêm một cây con

ID3( $Examples_{v_i}$ , *Target\_attribute*, *Attribute* - { $A$ }).

End

5. Return *Root*.



# Thuộc tính tốt nhất: Độ đo Gini

- Bước 4.1. chọn thuộc tính A tốt nhất gán cho nút t.
- Tồn tại một số độ đo: Gini, Information gain...
- **Độ đo Gini**

- Đo tính hỗn tạp của một tập ví dụ mẫu
- Công thức tính độ đo Gini cho nút t:

$$Gini(t) = 1 - \sum_j p(j|t)^2$$

Trong đó  $p(j|t)$  là tần suất liên quan của lớp  $j$  tại nút  $t$

- Gini (t) lớn nhất =  $1 - 1/n_c$  (với  $n_c$  là số các lớp tại nút t): khi các bản ghi tại t phân bố đều cho  $n_c$  lớp; tính hỗn tạp cao nhất, không có phân biệt giữa các lớp
  - Gini (t) nhỏ nhất = 0 khi tất cả các bản ghi thuộc một lớp duy nhất.
- **Ví dụ:** Bốn trường hợp

C1	<b>0</b>
C2	<b>6</b>
<b>Gini=0.000</b>	

C1	<b>1</b>
C2	<b>5</b>
<b>Gini=0.278</b>	

C1	<b>2</b>
C2	<b>4</b>
<b>Gini=0.444</b>	

C1	<b>3</b>
C2	<b>3</b>
<b>Gini=0.500</b>	

# Chia tập theo độ đo Gini



- Dùng trong các thuật toán CART, SLIQ, SPRINT
- Khi một nút  $t$  được phân hoạch thành  $k$  phần ( $k$  nút con của  $t$ ) thì chất lượng của việc chia tính bằng

$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i)$$

trong đó

- $n$  là số bản ghi của tập bản ghi tại nút  $t$ ,
- $n_i$  là số lượng bản ghi tại nút con  $i$  (của nút  $t$ ).



# Chia tập theo độ đo Gini: Ví dụ

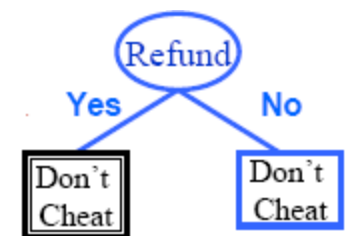
- Tính toán GINI cho Refund (Yes, No), Marital Status (Single&Divorced, Married) và Taxable Income (<80K, 80K-120K, >120K).
- Refund:  $3/10 * (0) + 7/10 * (1-(3/7)^2 - (4/7)^2) = 7/10*(24/49) = 24/70$
- Marital Status:  $4/10 * 0 + 6/10 * (1- (3/6)^2 - (3/6)^2) = 6/10 * 1/2 = 3/10$
- Taxable Income: thuộc tính liên tục cần chia khoảng (tồn tại một số phương pháp theo Gini, kết quả 2 thùng và 80K là mốc)  
 $3/10 * (0) + 7/10 * (1-(3/7)^2 - (4/7)^2) = 7/10*(24/49) = 24/70$

Như vậy, Gini của Refund và Taxable Income bằng nhau (24/70) và lớn hơn Gini của Marital Status (3/10) nên chọn Refund cho gốc cây quyết định.

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

$$GINI_{split} = \sum_i p_i GINI(i)$$

$$Gini(t) = \sum_j p_j Gini(j|t)$$





# Chọn thuộc tính: Information Gain



- Độ đo Information Gain

- Thông tin thu được sau khi phân hoạch tập ví dụ
- Dùng cho các thuật toán ID3, họ C4.5

- Entropy

- Công thức tính entropy nút t:

$$Entropy(t) = - \sum_j p(j|t) \log p(j|t)$$

Trong đó  $p(j|t)$  là tần suất liên quan của lớp  $j$  tại nút  $t$  độ không đồng nhất tại nút  $t$ .

- Entropy (t) lớn nhất =  $\log(n_c)$  (với  $n_c$  là số các lớp tại nút t): khi các bản ghi tại t phân bố đều cho  $n_c$  lớp; tính hỗn tạp cao nhất, không có phân biệt giữa các lớp
- Entropy (t) nhỏ nhất = 0 khi tất cả các bản ghi thuộc một lớp duy nhất.
- Lấy loga cơ số 2 thay cho loga tự nhiên

- Tính toán entropy (t) cho một nút tương tự như Gini (t)

# Chọn thuộc tính: Information Gain



- Độ đo Information Gain

$$Gain_{chia} = entropy(t) - \sum_{i=1}^k \frac{n_i}{n} entropy(i)$$

Trong đó,  $n$  là số lượng bản ghi tại nút  $t$ ,  $k$  là số tập con trong phân hoạch,  $n_i$  là số lượng bản ghi trong tập con thứ  $i$ .

Độ đo giảm entropy sau khi phân hoạch: chọn thuộc tính làm cho Gain đạt lớn nhất.

C4.5 là một trong 10 thuật toán KPD L phổ biến nhất.

- Hạn chế: Xu hướng chọn phân hoạch chia thành nhiều tập con

- Cải tiến

$$GainRATIO = \frac{Gain_{chia}}{SplitINFO} \quad SplitINFO = \sum_{i=1}^k \frac{n_i}{n} \log \frac{n_i}{n}$$

- Dùng GainRatio để khắc phục xu hướng chọn phân hoạch nhiều tập con

- Áp dụng: Tự tiến hành

# Phân lớp dựa trên luật



- Giới thiệu

- Phân lớp các bản ghi dựa vào tập các luật “kiểu” if ... then

- Luật

- Luật: <điều kiện> [redacted]

Trong đó:

<điều kiện> là sự kết nối các thuộc tính (còn gọi là tiên đề/điều kiện của luật: LHS bên trái)

y là nhãn lớp (còn gọi là kết quả của luật: RHS bên phải).

- Ví dụ

Refund = ‘Yes’ [redacted] heat = ‘No’

(Refund = ‘No’) [redacted] Marital Status = ‘Married’) [redacted] heat = ‘No’

- Sử dụng luật

- Một luật được gọi là “bảo đảm” thể hiện r (bản ghi) nếu các thuộc tính của r đáp ứng điều kiện của luật.
- Khi đó, vế phải của luật cũng được áp dụng cho thể hiện.

# Xây dựng luật phân lớp



- **Giới thiệu**

- Trực tiếp và gián tiếp

- **Trực tiếp**

- Trích xuất luật trực tiếp từ dữ liệu
- Ví dụ: RIPPER, CN2, Holte's 1R
- Trích xuất luật trực tiếp từ dữ liệu
  1. Bắt đầu từ một tập rỗng
  2. Mở rộng luật bằng hàm Học\_một\_luật
  3. Xóa mọi bản ghi “bảo đảm” bởi luật vừa được học
  4. Lặp các bước 2-3 cho đến khi gặp điều kiện dừng

- **Gián tiếp**

- Trích xuất luật từ mô hình phân lớp dữ liệu khác, chẳng hạn, mô hình cây quyết định, mô hình mạng nơ ron, ...
- Ví dụ: C4.5Rule

# Mở rộng luật: một số phương án



- Sử dụng thống kê

- Thống kê các đặc trưng cho ví dụ
- Tìm đặc trưng điển hình cho từng lớp

- Thuật toán CN2

- Khởi đầu bằng liên kết rỗng: {}
- Bổ sung các liên kết làm cực tiểu entropy: {A}, {A, B}...
- Xác định kết quả luật theo đa số của các bản ghi đảm bảo luật

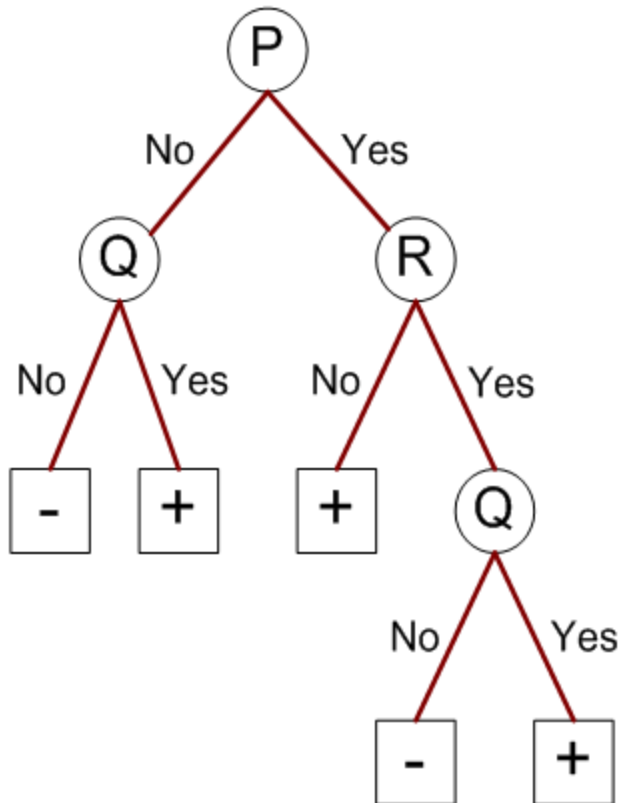
- Thuật toán RIPPER

- Bắt đầu từ một luật rỗng: {}
- Bổ sung các liên kết làm cực đại lợi ích thông tin FAIL
- R0: {} => lớp (luật khởi động)
- R1: {A} => lớp (quy tắc sau khi thêm liên kết)
- Gain (R0, R1) = t [log (p1 / (p1 + n1)) - log (p0 / (p0 + n0))]

với t: số thể hiện đúng đảm bảo cả hai R0 và R1

- p0: số thể hiện đúng được bảo đảm bởi R0
- n0: số thể hiện sai được đảm bảo bởi R0
- P1: số thể hiện đúng được bảo đảm bởi R1
- n 1: số trường hợp sai được đảm bảo bởi R1

# Luật phân lớp: từ cây quyết định



## Tập luật

Liệt kê các đường đi từ gốc

r1: (P=No,Q=No) ==> -

r2: (P=No,Q=Yes) ==> +

r3: (P=Yes,R=No) ==> +

r4: (P=Yes,R=Yes,Q=No) ==> -

r5: (P=Yes,R=Yes,Q=Yes) ==> +

# Sinh luật gián tiếp: C4.5rules



- Trích xuất luật từ cây quyết định chưa cắt tỉa
- Với mỗi luật,  $r: A \rightarrow y$ 
  - Xem xét luật thay thế  $r': A' \rightarrow y$ , trong đó  $A'$  nhận được từ  $A$  bằng cách bỏ đi một liên kết
  - So sánh tỷ lệ lỗi  $r$  so với các  $r'$
  - Loại bỏ các  $r'$  có lỗi thấp hơn  $r$
  - Lặp lại cho đến khi không cải thiện được lỗi tổng thể
- Thay thế sắp xếp theo luật bằng sắp xếp theo tập con của luật (thứ tự lớp)
  - Mỗi tập con là một tập các luật với cùng một kết quả (lớp)
  - Tính toán độ dài mô tả của mỗi tập con
  - Độ dài mô tả =  $L(\text{lỗi}) + g * L(\text{mô hình})$
  - $g$  : tham số đếm sự hiện diện của các thuộc tính dư thừa trong một tập luật (giá trị chuẩn,  $g=0.5$ )

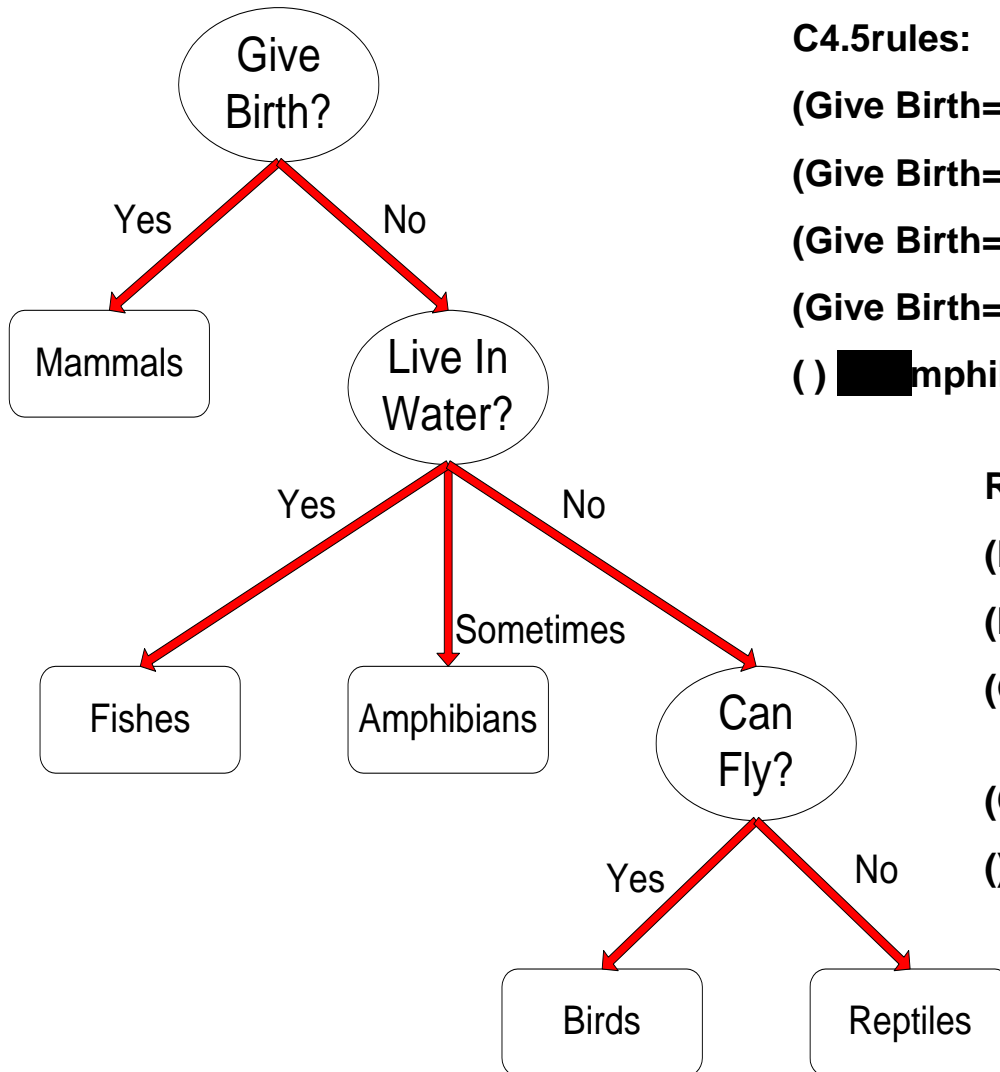
# C4.5rules: Ví dụ



Name	Give Birth	Lay Eggs	Can Fly	Live in Water	Have Legs	Class
human	yes	no	no	no	yes	mammals
python	no	yes	no	no	no	reptiles
salmon	no	yes	no	yes	no	fishes
whale	yes	no	no	yes	no	mammals
frog	no	yes	no	sometimes	yes	amphibians
komodo	no	yes	no	no	yes	reptiles
bat	yes	no	yes	no	yes	mammals
pigeon	no	yes	yes	no	yes	birds
cat	yes	no	no	no	yes	mammals
leopard shark	yes	no	no	yes	no	fishes
turtle	no	yes	no	sometimes	yes	reptiles
penguin	no	yes	no	sometimes	yes	birds
porcupine	yes	no	no	no	yes	mammals
eel	no	yes	no	yes	no	fishes
salamander	no	yes	no	sometimes	yes	amphibians
gila monster	no	yes	no	no	yes	reptiles
platypus	no	yes	no	no	yes	mammals
owl	no	yes	yes	no	yes	birds
dolphin	yes	no	no	yes	no	mammals
eagle	no	yes	yes	no	yes	birds



# C4.5rules: Ví dụ



C4.5rules:

(Give Birth=No, Can Fly=Yes) ■■■irds

(Give Birth=No, Live in Water=Yes) ■■■ishes

(Give Birth=Yes) ■■■ammals

(Give Birth=No, Can Fly=No, Live in Water=No) ■■■eptiles

() ■■■mphibians

RIPPER:

(Live in Water=Yes) ■■■ishes

(Have Legs=No) ■■■eptiles

(Give Birth=No, Can Fly=No, Live In Water=No) ■■■eptiles

(Can Fly=Yes, Give Birth=No) ■■■irds

() ■■■ammals



# Phân lớp Bayes

- Giới thiệu

- Khung xác suất để xây dựng bộ phân lớp
- Xác suất có điều kiện

Hai biến cố A và C

$$P(C | A) = \frac{P(A, C)}{P(A)}$$

$$P(A | C) = \frac{P(A, C)}{P(C)}$$

- Định lý Bayes:

$$P(c|x) = P(x|c).P(c)/P(x)$$

- $P(x)$  bằng nhau cho tất cả các lớp
- Tìm c sao cho  $P(c|x)$  lớn nhất  $\Leftrightarrow$  Tìm c sao cho  $P(x|c).P(c)$  lớn nhất
- $P(c)$ : tần suất xuất hiện của các tài liệu thuộc lớp c
- Vấn đề: làm thế nào để tính  $P(x|c)$ ?



# Định lý Bayes: Ví dụ

- Một bác sỹ biết
  - Bệnh nhân viêm màng não có triệu chứng cứng cổ S|M: 50%
  - Xác suất một bệnh nhân bị viêm màng não M là 1/50.000
  - Xác suất một bệnh nhân bị cứng cổ S là 1/20
- Một bệnh nhân bị cứng cổ hỏi xác suất anh/cô ta bị viêm màng não ?

$$P(M | S) = \frac{P(S | M)P(M)}{P(S)} = \frac{0.5 / 50000}{1/20} = 0.0002$$

# Phân lớp Bayes



- Các thuộc tính (bao gồm nhãn lớp) là các biến ngẫu nhiên.
- Cho một bản ghi với các giá trị thuộc tính  $(A_1, A_2, \dots, A_n)$ 
  - Cần dự báo nhãn  $c$
  - Tìm lớp  $c$  để cực đại xác suất  $P(C|A_1, A_2, \dots, A_n)$
- Có thể tính xác suất  $P(C|A_1, A_2, \dots, A_n)$  từ dữ liệu học?



# Phân lớp văn bản Naïve Bayes

- Giả thiết Naïve Bayes:
  - giả thiết độc lập: xác suất xuất hiện của một từ khóa trong văn bản độc lập với ngữ cảnh và vị trí của nó trong văn bản:

$$p(c | x, T) = p(c | x) \prod_{T \in \text{words}(x)} p(T | c)$$

$$P(\mathbf{x}_1, \dots, \mathbf{x}_k | C) = P(\mathbf{x}_1 | C) \cdot \dots \cdot P(\mathbf{x}_k | C)$$

# Phân lớp văn bản Naïve Bayes



## • Cho

- Tập ví dụ  $D_{\text{exam}} = D_{\text{learn}} + D_{\text{test}}$
- Tập từ vựng  $V = \{f_1, f_2, \dots, f_{||V||}\}$
- Tập lớp  $C = \{C_1, C_2, \dots, C_n\}$  với mỗi  $C_i$  một ngưỡng  $\theta_i > 0$

## • Tính xác suất tiên nghiệm

- Trên tập ví dụ học  $D_{\text{learn}}$
- $p(C_i) = M_i/M$ ,  $M = ||D_{\text{learn}}||$ ,  $M_i = ||\text{Doc} \in C_i||$
- Xác suất một đặc trưng (từ)  $f_j$  thuộc lớp  $C$ :

$$P(f_j | C) = \frac{1 + TF(f_j, C)}{|V| + \sum_i TF(f_j, C_i)}$$

## • Cho tài liệu Doc mới

- Tính xác suất hậu nghiệm
- Nếu  $P(C|Doc) > \theta_i$  thì Doc  $\in C_i$ !

$$P(C | Doc) = \frac{p(C) * \prod_{F_j} (F_j | C)^{TF(F_j, Doc)}}{\sum_i p(C_i) * \prod_{F_j} (F_j | C_i)^{TF(F_j, Doc)}}$$

# Công thức phân lớp Bayes thứ hai



$$P(F_j|C) = \frac{1 + TF(F_j, C)}{|V| + \sum_{i=1}^n TF(F_i, C)}$$

$$C|Doc) = \frac{P(C) \times \prod_{F_j \in Doc} P(F_j|C)^{TF(F_j, Doc)}}{\sum_{i=1}^n P(C_i) \times \prod_{F_i \in Doc} P(F_i|C_i)^{TF(F_i, Doc)}}$$



# Phân lớp k-NN

$$Sm(Doc, D_i) = \frac{\text{Cos}(Doc, D_i) * Y_l}{\sqrt{\frac{\sum_{l=1}^L \text{Doc}_l^2}{l} + \frac{\sum_{l=1}^L D_i_l^2}{l}}}$$

- **Cho trước**

- Một tập D các tài liệu biểu diễn bản ghi các đặc trưng
- Một đo đo khoảng cách (Ơclit) hoặc tương tự (như trên)
- Một số  $k > 0$  (láng giềng gần nhất)

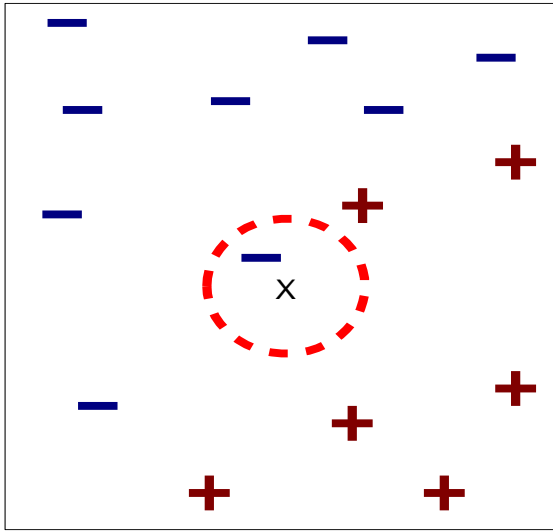
- **Phân lớp tài liệu mới Doc được biểu diễn**

- Tính khoảng cách (độ tương tự) từ Doc tới tất cả tài liệu thuộc D
- Tìm k tài liệu thuộc D gần Doc nhất
- Dùng nhãn lớp của k-láng giềng gần nhất để xác định nhãn lớp của Doc: nhãn nhiều nhất trong k-láng giềng gần nhất

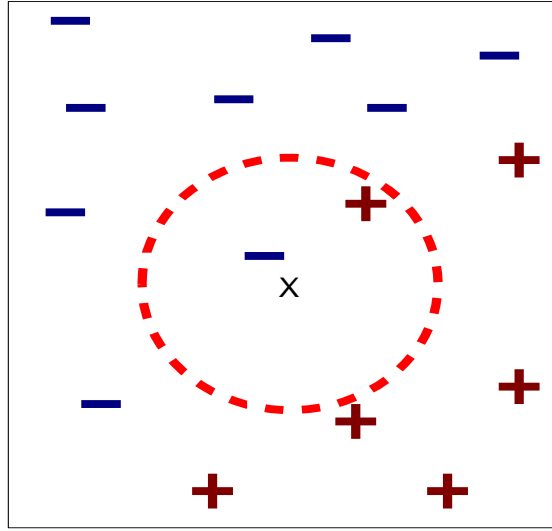




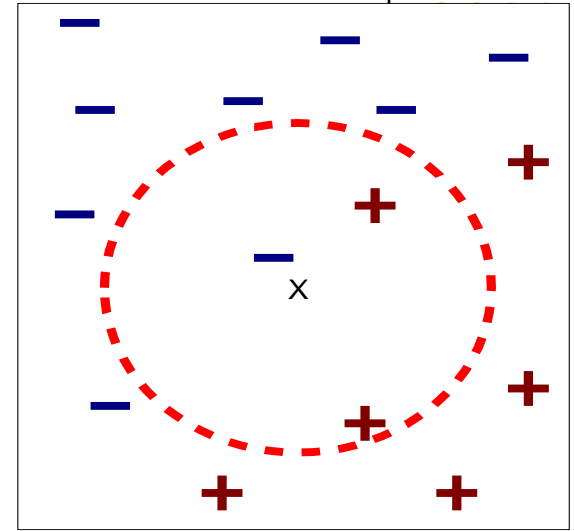
# Phân lớp k-NN: Ví dụ



(a) 1-nearest neighbor



(b) 2-nearest neighbor



(c) 3-nearest neighbor

- **Ba trường hợp như hình vẽ**

- 1-NN: Chọn lớp "-": lát giềng có nhãn "-" là nhiều nhất
- 2-NN: Chọn lớp "-": hai nhãn có số lượng như nhau, chọn nhãn có tổng khoảng cách gần nhất
- 3-NN: Chọn lớp "+": lát giềng có nhãn "+" là nhiều nhất

# Thuật toán SVM



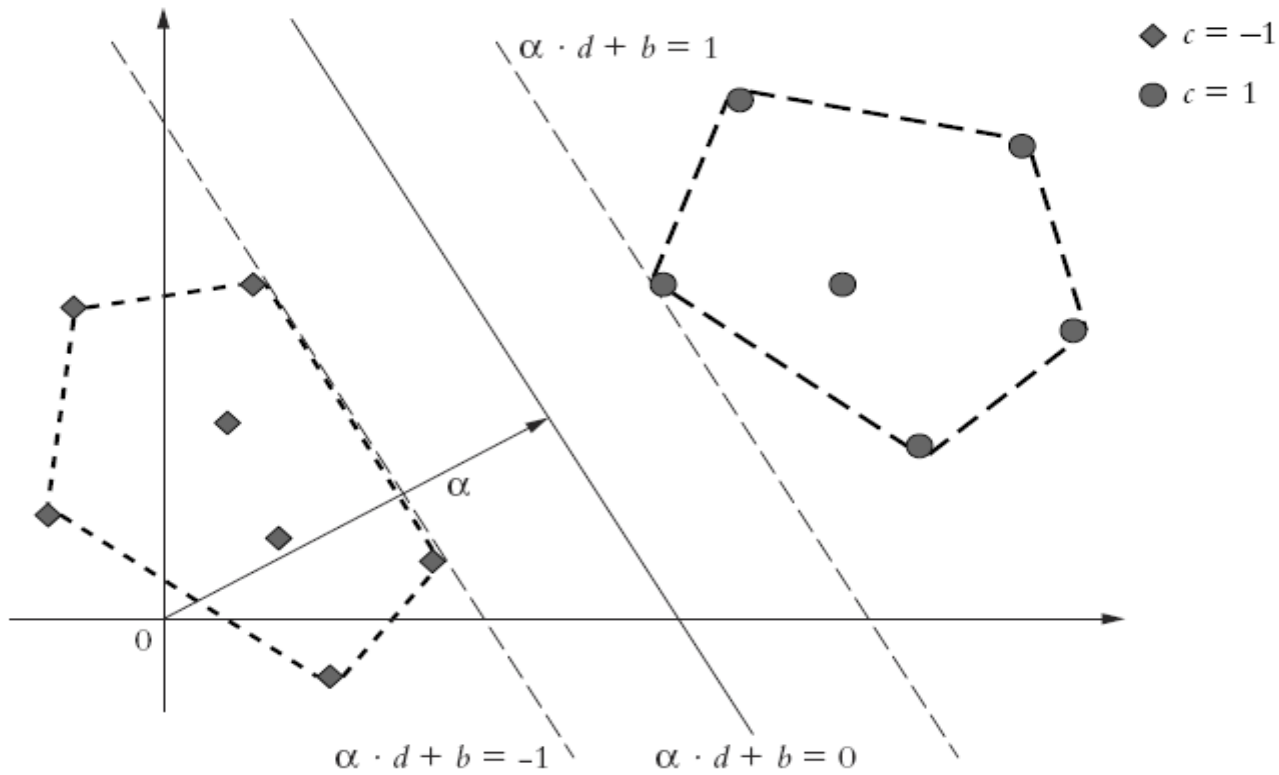
- Thuật toán máy vector hỗ trợ (Support Vector Machine – SVM): được Corters và Vapnik giới thiệu vào năm 1995.
- SVM rất hiệu quả để giải quyết các bài toán với dữ liệu có số chiều lớn (như các vector biểu diễn văn bản).

# Thuật toán SVM



- Tập dữ liệu học:  $D = \{(X_i, C_i), i=1, \dots, n\}$ 
  - $C_i \in \{-1, 1\}$  xác định dữ liệu dương hay âm
- Tìm một siêu phẳng:  $\alpha_{SVM} \cdot \mathbf{d} + \mathbf{b}$  phân chia dữ liệu thành hai miền.
- Phân lớp một tài liệu mới: xác định dấu của
  - $f(d) = \alpha_{SVM} \cdot \mathbf{d} + \mathbf{b}$
  - Thuộc lớp dương nếu  $f(d) > 0$
  - Thuộc lớp âm nếu  $f(d) < 0$

# Thuật toán SVM



# Thuật toán SVM



- Nếu dữ liệu học là tách rời tuyến tính:

- Cực tiểu:

$$\frac{1}{2}$$

- Thỏa mãn:

$$c_i$$

(2)

- Nếu dữ liệu học không tách rời tuyến tính: thêm biến  $\{\xi_1 \dots \xi_n\}$ :

- Cực tiểu:

$$\frac{1}{2}$$

- Thỏa mãn:

$$c_i$$

(4)

# Phân lớp Web bán giám sát



- Giới thiệu phân lớp bán giám sát web
  - Khái niệm sơ bộ
  - Tại sao học bán giám sát
- Nội dung phân lớp bán giám sát web
  - Một số cách tiếp cận cơ bản
  - Các phương án học bán giám sát phân lớp web
- Phân lớp bán giám sát trong NLP

# Học bán giám sát: Tài liệu tham khảo



1. Xiaojin Zhu ([2006](#) \*\*\*). Semi-Supervised Learning Literature Survey, 1-2006. (Xiao Zhu [1])  
[http://www.cs.wisc.edu/~jerryzhu/pub/ssl\\_survey.pdf](http://www.cs.wisc.edu/~jerryzhu/pub/ssl_survey.pdf)
- Zhou, D., Huang, J., & Scholkopf, B. ([2005](#)). Learning from labeled and unlabeled data on a directed graph. *ICML05, 22nd International Conference on Machine Learning*. Bonn, Germany.
- Zhou, Z.-H., & Li, M. ([2005](#)). Semi-supervised regression with co-training. *International Joint Conference on Artificial Intelligence (IJCAI)*.
- Zhu, X. ([2005](#)). *Semi-supervised learning with graphs*. Doctoral dissertation, Carnegie Mellon University (mã số CMU-LTI-05-192).
1. Olivier Chapelle, Mingmin Chi, Alexander Zien ([2006](#)) A Continuation Method for Semi-Supervised SVMs. Proceedings of the 23rd International Conference on Machine Learning, Pittsburgh, PA, 2006.  
và [các tài liệu khác](#)

# Sơ bộ về học bán giám sát

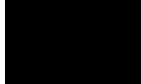


- Học bán giám sát là gì ? Xiao Zhu [1] FQA
  - Học giám sát: tập ví dụ học đã được gán nhãn (ví dụ gán nhãn) là tập các cặp (tập thuộc tính, nhãn)
  - ví dụ gán nhãn
    - Thủ công: khó khăn █████ chuyên gia █████ tốn thời gian, tiền
    - Tự động: như tự động sinh corpus song hiệu quả chưa cao
  - ví dụ chưa gán nhãn
    - Dễ thu thập █████ nhiều
      - xử lý tiếng nói: bài nói nhiều, xây dựng tài nguyên đòi hỏi công phu
      - xử lý văn bản: trang web vô cùng lớn, ngày càng được mở rộng
    - Có sẵn █████ điều kiện tiến hành tự động gán nhãn
  - Học bán giám sát: dùng cả ví dụ có nhãn và ví dụ chưa gán nhãn
    - Tạo ra bộ phân lớp tốt hơn so với chỉ dùng học giám sát: học bán giám sát đòi hỏi điều kiện về dung lượng khối lượng



# Cơ sở của học bán giám sát



- Biểu diễn dữ liệu chưa mô tả hết ánh xạ gán nhãn trên dữ liệu
  - chẳng hạn, nghịch lý “hiệu quả như nhau” trong biểu diễn văn bản
- Ánh xạ gán nhãn có liên quan mô hình dữ liệu (mô hình / đặc trưng/ nhân / hàm tương tự)   
mô hình đã có theo tự nhiên hoặc giả thiết dữ liệu tuân theo.



# Hiệu lực của học bán giám sát

- **Dữ liệu chưa nhãn không luôn luôn hiệu quả**
  - Nếu giả thiết mô hình không phù hợp ██████ ảnh hưởng hiệu quả
  - Một số phương pháp cần điều kiện về miền quyết định: tránh miền có mật độ cao:
    - Transductive SVM (máy hỗ trợ vector lan truyền)
    - Information Regularization (quy tắc hóa thông tin)
    - mô hình quá trình Gauss với nhiễu phân lớp bằng không
    - phương pháp dựa theo đồ thị với trọng số cạnh là khoảng cách
  - “Tồi” khi dùng phương pháp này song lại “tốt” khi dùng phương pháp khác

# Phương pháp học bán giám sát



- Các phương pháp học bán giám sát điển hình
  - EM với mô hình trộn sinh
  - Self-training
  - Co-training
  - TSVM
  - Dựa trên đồ thị
  - ...
- So sánh các phương pháp
  - Đòi hỏi các giả thiết mô hình mạnh. Giả thiết mô hình phù hợp cấu trúc dữ liệu: khó kiểm nghiệm
  - Một số định hướng lựa chọn
    - Lớp ██████ phân cụm tốt: dùng EM với mô hình sinh trộn.
    - Đặc trưng phân thành hai phần riêng rẽ: co-training
    - Nếu hai điểm tương tự hướng tới một lớp: dựa trên đồ thị
    - Đã sử dụng SVM thì mở rộng TSVM
    - Khó nâng cấp học giám sát đã có: dùng self-training

# Phương pháp học bán giám sát



- **Dùng dữ liệu chưa gán nhãn**
  - Hoặc biến dạng hoặc thay đổi thứ tự giả thiết thu nhờ chỉ dữ liệu có nhãn
  - Mô tả chung
    - Giả thiết dưới dạng  $p(y|x)$  còn dữ liệu chưa có nhãn  $p(x)$
    - Mô hình sinh có tham số chung phân bố kết nối  $p(x, y)$
    - Mô hình trộn với EM mở rộng thêm self-training
    - Nhiều phương pháp là phân biệt: TSVM, quy tắc hóa thông tin, quá trình Gauss, dựa theo đồ thị
  - Có dữ liệu không nhãn: nhận được xác suất  $p(x)$
- **Phân biệt “học lan truyền” với “học bán giám sát”**
  - Đa dạng về cách gọi. Hạn chế bài toán phân lớp.
  - “Bán giám sát”
    - dùng ví dụ có / không có nhãn,
    - “học dữ liệu nhãn/không nhãn,
    - “học dữ liệu phân lớp/có nhãn bộ phận”.
    - Có cả lan truyền hoặc quy nạp.
  - Lan truyền để thu hẹp lại cho quy nạp: học chỉ dữ liệu sẵn. Quy nạp: có thể liên quan tới dữ liệu chưa có.

# Mô hình sinh: Thuật toán EM



## ● Sơ bộ

- Mô hình sớm nhất, phát triển lâu nhất
- Mô hình có dạng  $p(x,y) = p(y)*p(x|y)$
- Với số lượng nhiều dữ liệu chưa nhãn cho  $P(x|y)$  mô hình trộn đồng nhất. Miền tài liệu được phân thành các thành phần,
- Lý tưởng hóa tính "Đồng nhất": chỉ cần một đối tượng có nhãn cho mỗi thành phần

## ● Tính đồng nhất

- Là tính chất cần có của mô hình
- Cho họ phân bố  $\{p_{\theta}\}$  là đồng nhất nếu  $\theta_1, \theta_2$  thì  $p_{\theta_1}$  cho tới một hoán đổi vị trí các thành phần  $\theta_2$  khả tách của phân bố tới các thành phần



# Mô hình sinh: Thuật toán EM

- Tính xác thực của mô hình
  - Giả thiết mô hình trộn là chính xác [REDACTED] dữ liệu không nhãn sẽ làm tăng độ chính xác phân lớp
  - Chú ý cấu trúc tốt mô hình trộn: nếu tiêu đề được chia thành các tiêu đề con thì nên mô hình hóa thành đa chiều thay cho đơn chiều
- Cực đại EM địa phương
  - Miền áp dụng
    - Khi mô hình trộn chính xác
  - Ký hiệu
    - $D$ : tập ví dụ đã có (có nhãn /chưa có nhãn)
    - $D^K$ : tập ví dụ có nhãn trong  $D$  ( $|D^K| \ll |D|$ )

# Mô hình sinh: Thuật toán EM



- Nội dung thuật toán

1: Cố định tập tài liệu không nhãn  $D^U \subseteq D \setminus D^K$  dùng trong E-bước và M-bước

2: dùng  $D^K$  xây dựng mô hình ban đầu

3: **for**  $i = 0, 1, 2, \dots$  cho đến khi kết quả đảm bảo **do**

4: **for** mỗi tài liệu  $d \in D^U$  **do**

5: E-bước: dùng phân lớp Bayes thứ nhất xác định  $P(c|d)$

6: **end for**

7: **for** mỗi lớp  $c$  và từ khóa  $t$  **do**

8: M-bước: xác định  $\theta_{c,t}$  dùng công thức (\*) để xây dựng mô hình  $i+1$

9: **end for**

10: **end for**

$$P(d|c) = P(L = \ell_d | c) \binom{\ell_d}{\{n(d, t)\}} \prod_{t \in d} \theta_t^{n(d, t)}$$

$$\theta_{c,t} = \frac{1 + \sum_{d \in D} P(c|d) n(d, t)}{|W| + \sum_{d \in D} \sum_{\tau} P(c|d) n(d, \tau)} \quad P(c) = \frac{1}{|D|} \sum_{d \in D} P(c|d)$$

# Mô hình sinh: Thuật toán EM



- Một số vấn đề với EM
  - Phạm vi áp dụng: mô hình trộn chính xác
  - Nếu cực trị địa phương khác xa cực trị toàn cục thì khai thác dữ liệu không nhãn không hiệu quả
  - "Kết quả đảm bảo yêu cầu": đánh giá theo các độ đo hồi tưởng, chính xác,  $F_1$ ...
  - Một số vấn đề khác cần lưu ý:
    - Thuật toán nhân là Bayes naive: có thể chọn thuật toán cơ bản khác
    - Chọn điểm bắt đầu bằng học tích cực



# Mô hình sinh: Thuật toán khác



- **Phân cụm - và - Nhãn**

- Sử dụng phân cụm cho toàn bộ ví dụ
  - cả dữ liệu có nhãn và không có nhãn
  - dành tập  $D_{\text{test}}$  để đánh giá
- Độ chính xác phân cụm cao
  - Mô hình phân cụm phù hợp dữ liệu
  - Nhãn cụm (nhãn dữ liệu có nhãn) làm nhãn dữ liệu khác

- **Phương pháp nhân Fisher cho học phân biệt**

- Phương pháp nhân là một phương pháp điển hình
- Nhân là gốc của mô hình sinh
- Các ví dụ có nhãn được chuyển đổi thành vector Fisher để phân lớp

# Self-Training



- **Giới thiệu**

- Là kỹ thuật phổ biến trong SSL
  - EM địa phương là dạng đặc biệt của self-training

- **Nội dung**

## ***Gọi***

L : Tập các dữ liệu gán nhãn.

U : Tập các dữ liệu chưa gán nhãn

## ***Lặp*** (cho đến khi U = [REDACTED])

Huấn luyện bộ phân lớp giám sát h trên tập L

Sử dụng h để phân lớp dữ liệu trong tập U

Tìm tập con U' [REDACTED] có độ tin cậy cao nhất:

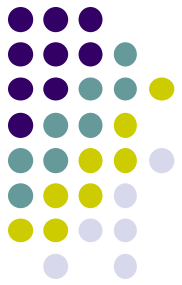
$$L + U' \quad \text{[REDACTED]}$$

$$U - U' \quad \text{[REDACTED]}$$

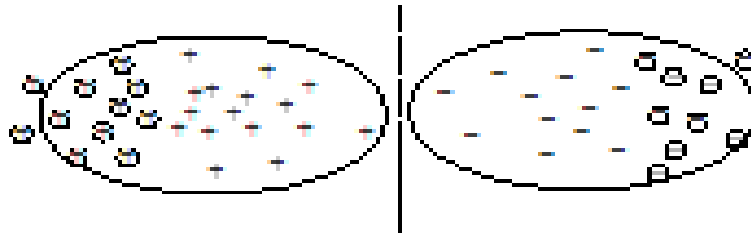
Vấn đề tập U' có "độ tin cậy cao nhất"

- Thủ tục "bootstrapping"
- Thường được áp dụng cho các bài toán NLP

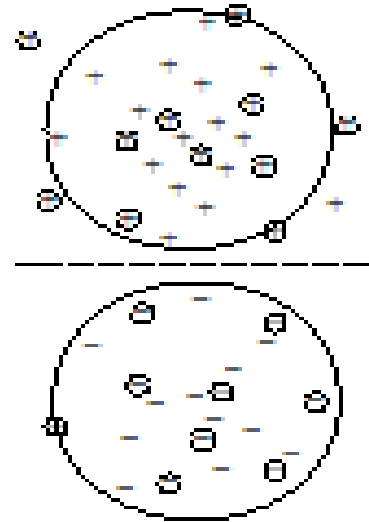
# Co-Training



- Tư tưởng
  - Một dữ liệu có hai khung nhìn
  - Ví dụ, các trang web
    - Nội dung văn bản
    - Tiêu đề văn bản



(a)  $x^1$  view

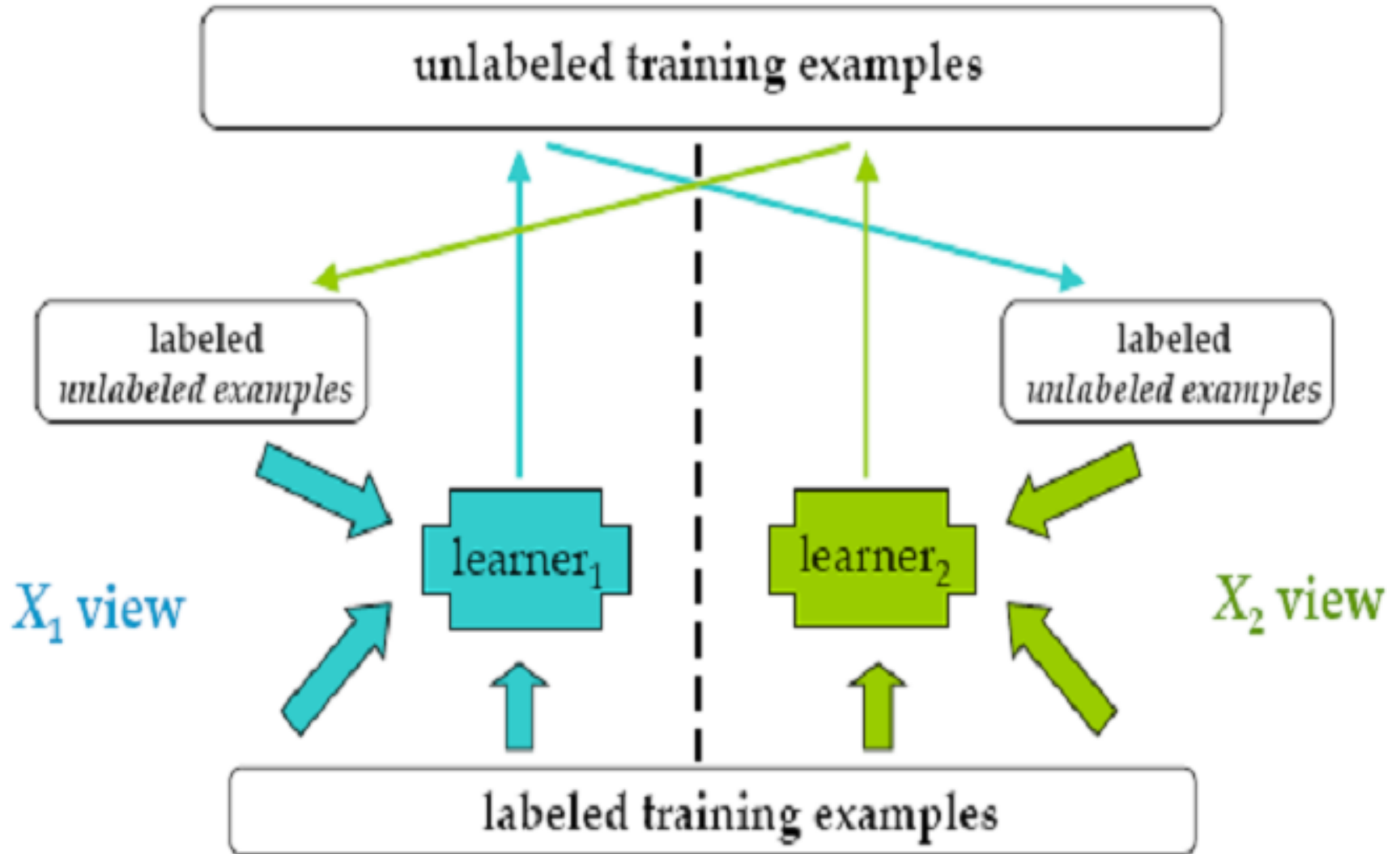


(b)  $x^2$  view

# Co-Training



- Mô hình thuật toán



# Co-Training



- Điều kiện dừng
  - hoặc tập dữ liệu chưa gán nhãn là rỗng
  - hoặc số vòng lặp đạt tới ngưỡng được xác định trước
- Một số lưu ý
  - Tập dữ liệu gán nhãn có ảnh hưởng lớn đến co-training
    - Quá ít: không hỗ trợ co-training
    - Quá nhiều: không thu lợi từ co-training
  - Cơ sở tăng hiệu quả co-training: thiết lập tham số
    - Kích cỡ tập dữ liệu gán nhãn
    - Kích cỡ tập dữ liệu chưa gán nhãn
    - Số các mẫu thêm vào sau mỗi vòng lặp
  - Bộ phân lớp thành phần rất quan trọng

# Chặn thay đổi miền dày đặc



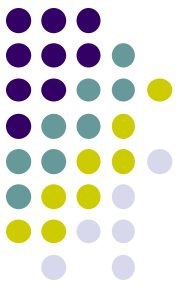
- Transductive SVMs (S3VMs)
  - Phương pháp phân biệt làm việc trên  $p(y|x)$  trực tiếp
  - Khi  $p(x)$  và  $p(y|x)$  không tương thích ████████ đưa  $p(x)$  ra khỏi miền dày đặc
- Quá trình Gauxơ)

# Mô hình đồ thị



- Biểu diễn dữ liệu chưa mô tả hết ánh xạ gán nhãn trên dữ liệu (chẳng hạn, nghịch lý “hiệu quả như nhau” trong biểu diễn văn bản)
- Ánh xạ gán nhãn có liên quan mô hình dữ liệu (mô hình / đặc trưng/ nhân / hàm tương tự)            mô hình đã có theo tự nhiên hoặc giả thiết dữ liệu tuân theo.

# Học bán giám sát với dữ liệu Web



- Tài liệu tham khảo
  - Soumen Chakrabarti (2003). *Mining the Web: Discovering Knowledge from Hypertext Data*. Morgan Kaufmann Publishers. Chương 6. SEMISUPERVISED LEARNING)
  - Các tài liệu về học máy [tài liệu chưa gán nhãn](#).
  - Pierre Baldi, Paolo Frasconi, Padhraic Smyth (2003). *Modeling the Internet and the Web: Probabilistic Methods and Algorithms*. Wiley, 2003, ISBN: 0-470-84906-1 ([Tài liệu giảng dạy 2](#)).



# Học bán giám sát với dữ liệu Web



- Một số thuật toán điển hình (xem [chương 6])
  - Expectation Maximization
    - *Experimental Results*
    - *Reducing the Belief in Unlabeled Documents*
    - *Modeling Labels Using Many Mixture Components*
  - Labeling Hypertext Graphs
    - *Absorbing Features from Neighboring Pages*
    - *A Relaxation Labeling Algorithm*
    - *A Metric Graph- Labeling Problem*
  - Co- training