

# Filtrado Basado en Contenido II

## IIC 3633 - Sistemas Recomendadores

Denis Parra

Profesor Asistente, DCC, PUC Chile

# Memo del Semestre

- **Tarea 1:** Deadline el Jueves 17 de Septiembre.
- **Lecturas en el semestre:** Chequear sitio Web curso. Daré plazo hasta el viernes 4 de Septiembre para que Uds. se inscriban en temas
  - Factorización Matricial: Nicolás Torres y Claudio Rojas (lecturas para el domingo)
  - Implicit Feedback: Alejandro Barrientos (lecturas para el domingo)
  - Active Learning RecSys: Javier Machin
  - Deep Learning RecSys: Gabriel de la Maggiore
  - Context-Aware Recsys Social: Julia Graller
  - Applications Music: Pierre Chaumier
  - Graph-based RecSys: Juan Pablo Salazar
- **Proyecto Final:**
  - Entrega de abstract con a lo más 3 ideas el martes 22 de Septiembre, el 29 de septiembre se debe entregar propuesta final.

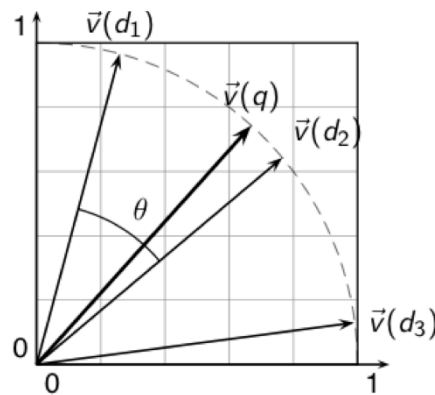
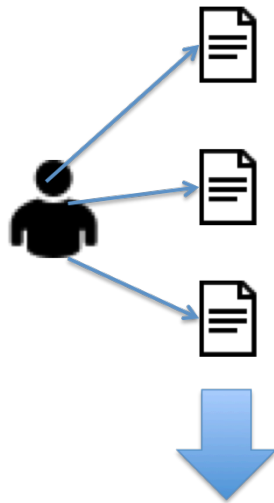
# TOC

## En esta clase

1. Representación y Aprendizaje del Modelo de Usuario
  1. Rocchio's Algorithm
  2. Revisión de Syskill & Webert
2. Tag-based recommendations
3. Actividad Práctica con R, RStudio y Packages tm() y ggplot2()

# Representación del Modelo de Usuario I

Bajo un escenario estático, podemos agregar el modelo del usuario en función de los contenidos por los cuales ha mostrado preferencias.



Doc\_1 = {w\_1, w\_2, ..., w\_3}



Doc\_2 = {w\_1, w\_2, ..., w\_3}



Doc\_3 = {w\_1, w\_2, ..., w\_3}



Doc\_n = {w\_1, w\_2, ..., w\_3}

user\_profile = {w\_1, w\_2, ..., w\_3} usando TF-IDF

- ¿Cómo incorporar información para actualizar el modelo?
- ¿Debiéramos "decaer" la importancia de ítems antiguos?
- ¿Distinta metadata debería integrarse en una sola bolsa de palabras o separarse?

# Representación del Modelo de Usuario II

## Modelo de Relevancia Rocchio

- Actualización de Query en base a feedback Positivo y Negativo

$$Q_{i+1} = \alpha Q_i + \beta \sum_{rel} \frac{D_i}{|D_i|} - \gamma \sum_{nonrel} \frac{D_i}{|D_i|}$$

- Clasificador que combina vectores de documentos para cada clase

$$\omega_{ki} = \beta \cdot \sum_{\{d_j \in POS_i\}} \frac{\omega_{kj}}{|POS_i|} - \gamma \cdot \sum_{\{d_j \in NEG_i\}} \frac{\omega_{kj}}{|NEG_i|}$$

# Representación del Modelo de Usuario II

## Modelo de Relevancia Rocchio II

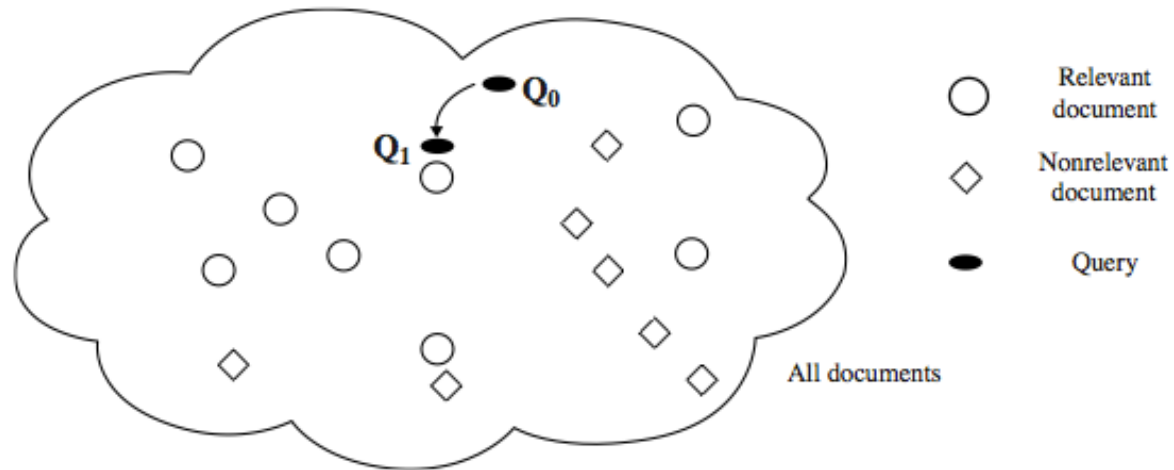


Figure 3.2. Relevance feedback. After feedback, the original query is moved toward the cluster of the relevant documents; see also Manning et al. (2008).

- En ausencia de feedback del usuario, también se suele usar pseudo relevance feedback.

# Representación del Modelo de Usuario III

## Modelo Bayesiano de Syskill & Webert

- El problema es recomendar páginas web de forma personalizada, con un algoritmo que aprenda a medida que el usuario muestra preferencias por ciertos ítems.

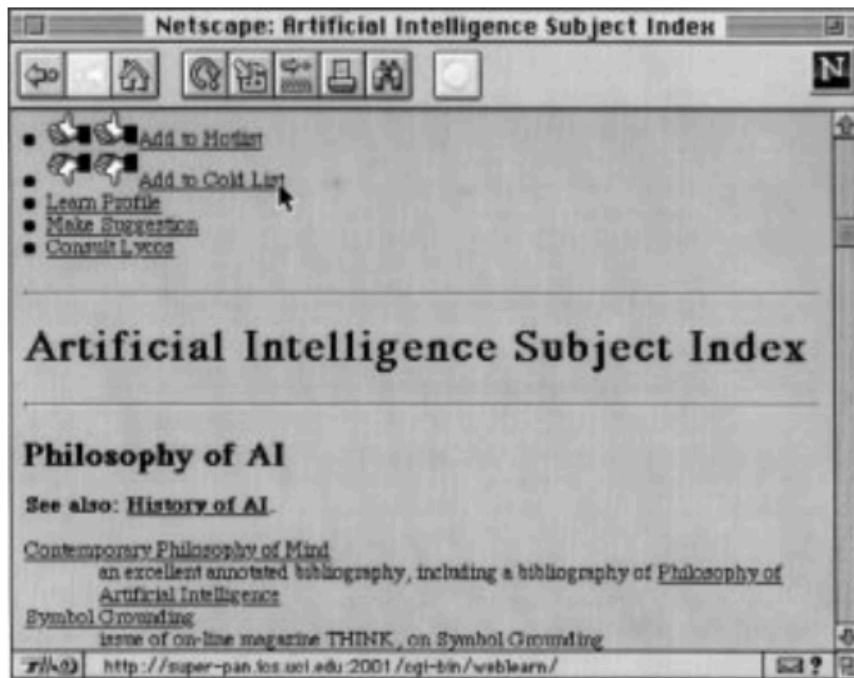


Figure 1. Syskill and Webert interface for rating pages.

# Representación del Modelo de Usuario III

## Syskill & Webert: Selección de Features

- No se usaron todas las palabras de la páginas, si no que las que proveían mayor information gain con respecto a clases predefinidas.

$$E(W, S) = I(S) - [P(W = \textit{present})I(S_{w=\textit{present}}) + P(W = \textit{absent})I(S_{w=\textit{absent}})]$$

donde

$$I(S) = \sum_{c \in \{\textit{hot}, \textit{cold}\}} -p(S_c) \log_2(p(S_c))$$

- Algunas de las features con mayor info gain.

<b>nirvana</b>	<b>suite</b>	<b>lo</b>
<b>pop</b>	<b>records</b>	<b>rockin</b>
<b>july</b>	<b>jams</b>	<b>songwriting</b>
<b>following</b>	<b>today</b>	<b>vocals</b>
<b>island</b>	<b>tribute</b>	<b>previous</b>
<b>favorite</b>	<b>airplay</b>	<b>noise</b>

**Table 1: Some of the words used as features**



# Representación del Modelo de Usuario III

## Modelo Bayesiano de Syskill & Webert III

- Se usa un clasificador Bayesiano

$$P(C_i | A_1 = V_{1j} \& \dots \& A_n = V_{nj}) \quad \text{donde} \quad P(C_i) \prod_k P(A_k = V_{kj} | C_i)$$

- Y el prior se va actualizando considerando una distribución Beta( $\alpha, \beta$ ), con  $\alpha$ : veces que el termino aparece, y  $\beta$ : numero de veces que el termino está ausente.

$$p(\theta | \alpha, \beta) = c \theta^\alpha (1 - \theta)^\beta p(\theta)$$

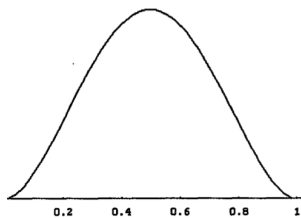


Figure 3: Beta distribution ( $\alpha = 2, \beta = 2$ )

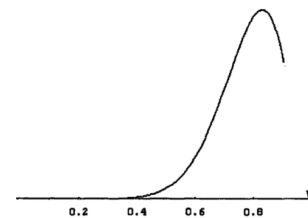
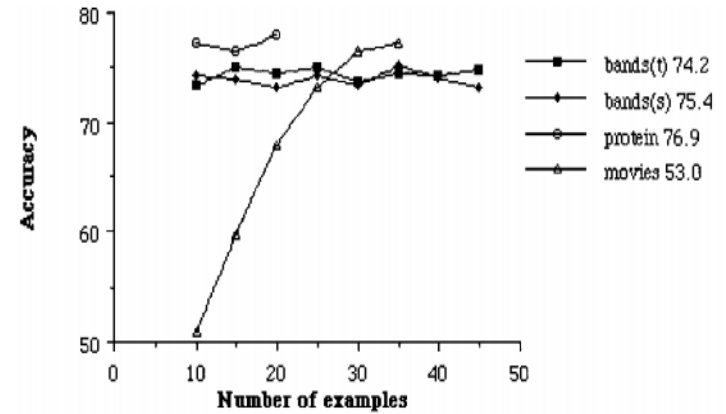
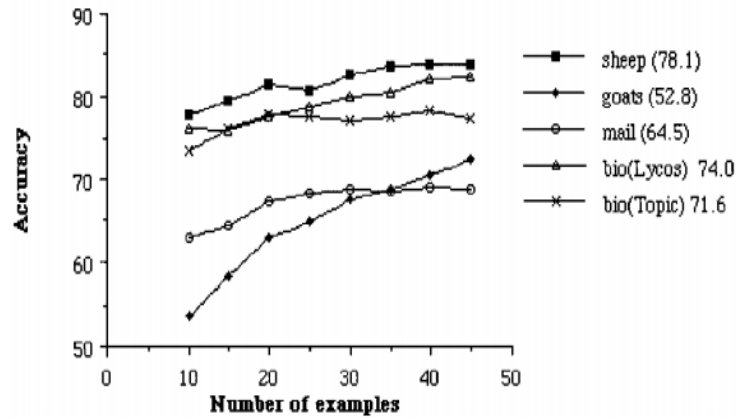


Figure 4: Beta distribution ( $\alpha = 10, \beta = 2$ )

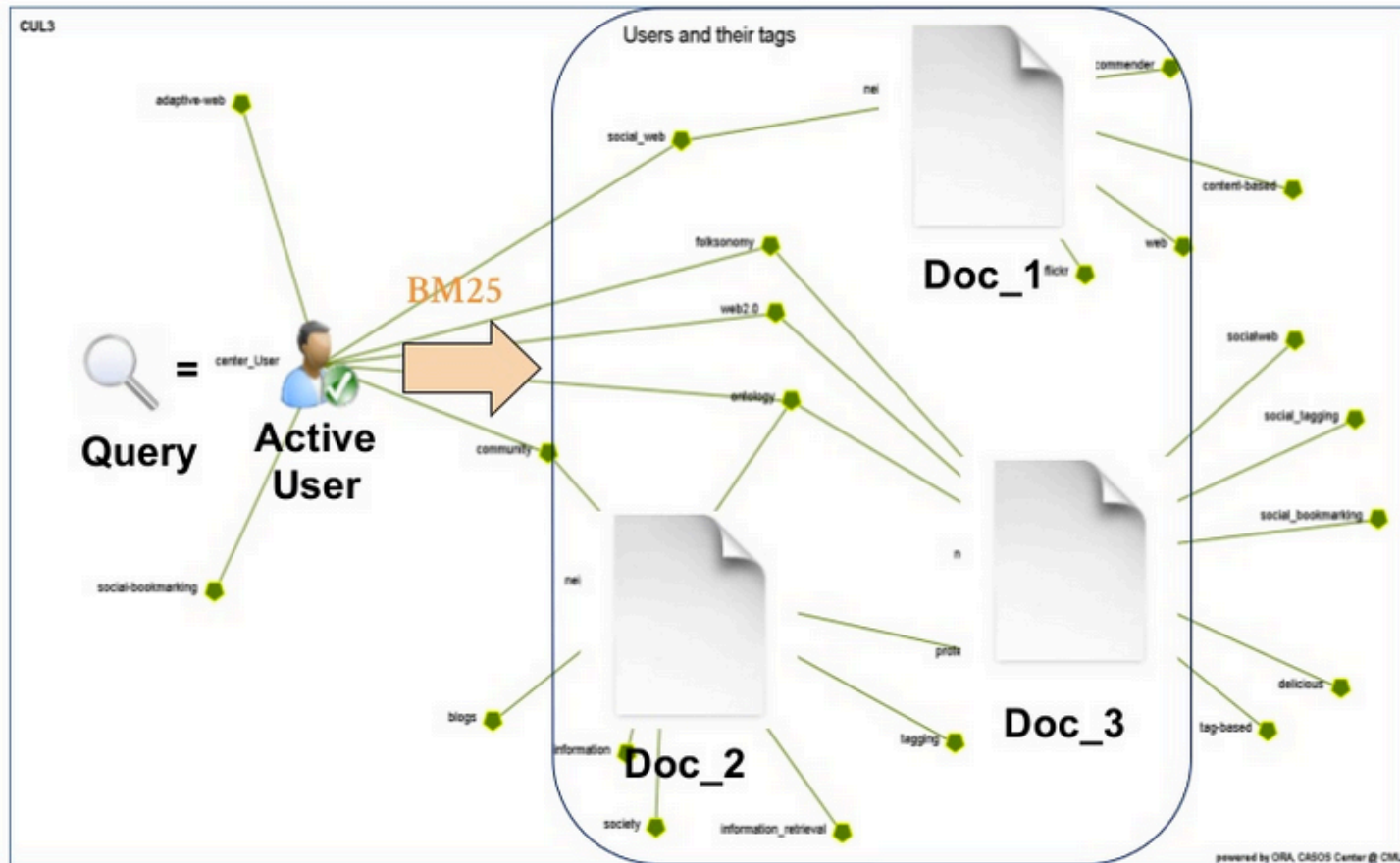
# Representación del Modelo de Usuario III

Resultados al observar user feedback



# Extensiones

## Tag-Based Recommendation



# MM-LDA (Ramage, 2009)

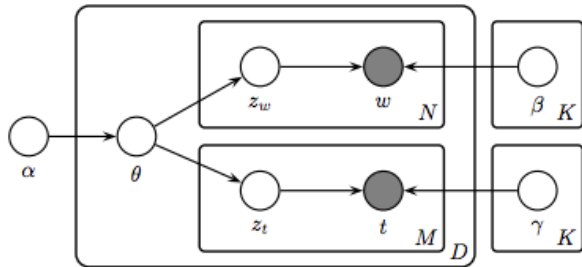


Figure 5: Graphical representation of MM-LDA.

## GENERAL CORPUS

	(MM-)LDA	K-means
Words	0.260	.139
Tags	0.270	.219
Words+Tags	<b>0.307</b>	.225

Figure 7: F-scores for (MM-)LDA and K-means on 13,320 documents. Including tags improves both models significantly versus words alone. MM-LDA (bold) significantly outperforms all other conditions.

## SPECIFIC CORPUS

		(MM-)LDA	K-means
Programming Languages	Words	.288	.189
	Tags	.463	.567
	Words+Tags	.297	.556
Social Sciences	Words	.300	.196
	Tags	.310	.307
	Words+Tags	.302	.308

Figure 10: F-scores for (MM-)LDA and K-means on two representative ODP subtrees. For these tasks, clustering on tags alone can outperform alternatives that use word information.

Ref: Ramage, D., Heymann, P., Manning, C. D., & Garcia-Molina, H. (2009, February). Clustering the tagged web. In Proceedings of the Second ACM International Conference on Web Search and Data Mining (pp. 54-63). ACM.

## Combinando Varios Features en un mismo Modelo

Context-Aware Event Recommendation in Event-based Social Networks by Augusto Q. Macedo, Leandro B. Marinho and Rodrygo L. T. Santos

- Recommendations on event-based social networks (EBSNs) often undermines the users' ability to choose the events that best fit their interests.
- The event recommendation problem is intrinsically cold-start. Indeed, events published in EBSNs are typically short-lived and, by definition, are always in the future, having little or no trace of historical attendance.
- THEN: exploit several contextual signals available from EBSNs. In particular, besides content-based signals based on the events' description and collaborative signals derived from users' RSVPs, we exploit social signals based on group memberships, location signals based on the users' geographical preferences, and temporal signals derived from the users' time preferences.

# Referencias

- Manning, C. D., Raghavan, P., & Schütze, H. (2008). Introduction to information retrieval (Vol. 1, p. 6). Cambridge: Cambridge university press.