



# Modelos Clásicos de Recuperación

M. Andrea Rodríguez Tastets

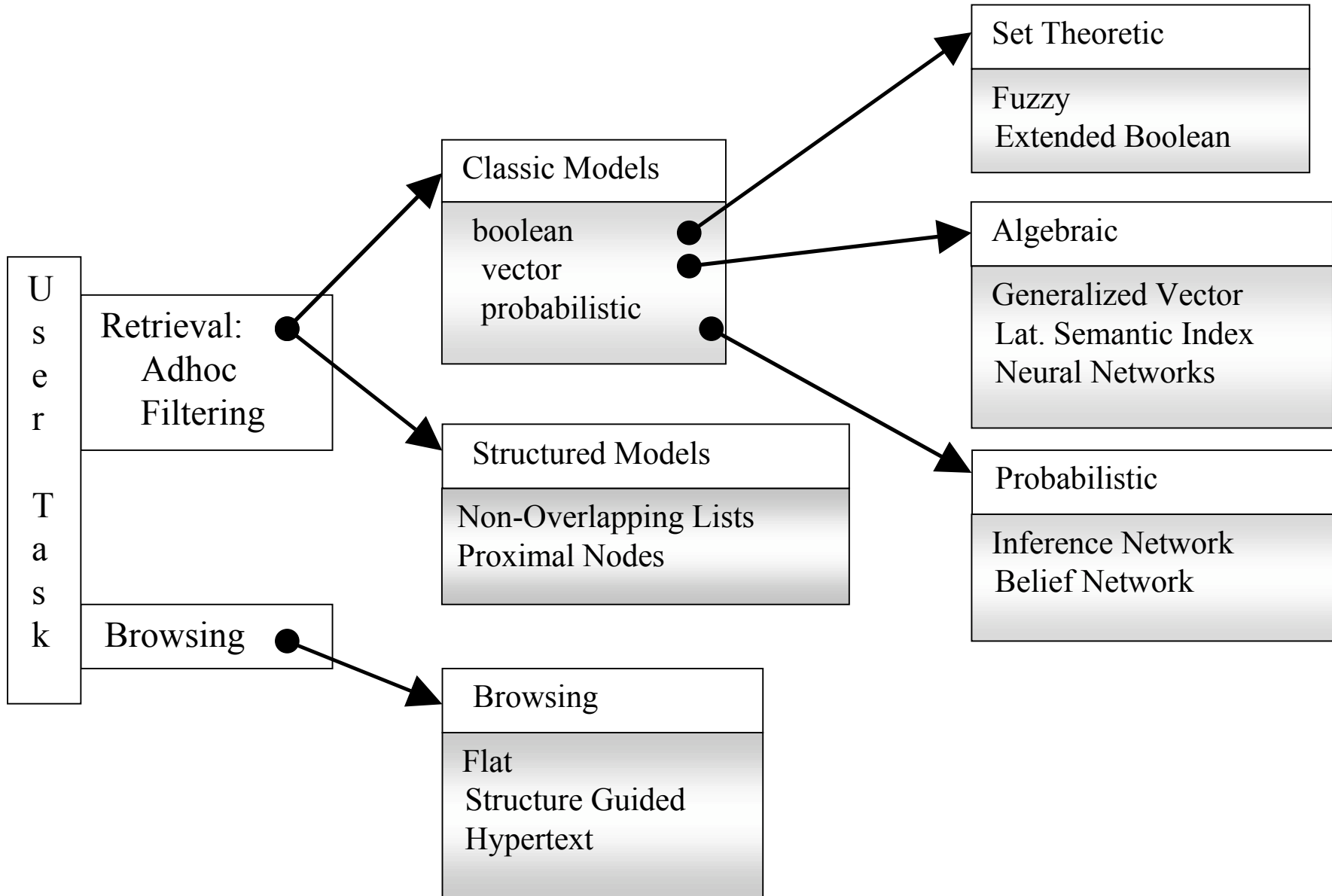
DIIC - Universidad de Concepción

<http://www.inf.udec.cl/~andrea>

# Introducción

- Un *ranking* es un orden de documentos recuperados que “ojalá” refleje la relevancia respecto a una consulta.
- Un ranking está basado en premisas fundamentales respecto a la noción de relevancia, tales como:
  - conjunto común de términos
  - compartir términos con pesos
  - similitud de relevancia
- Cada premisa en particular lleva a un modelo de Recuperación de Información

# Modelos



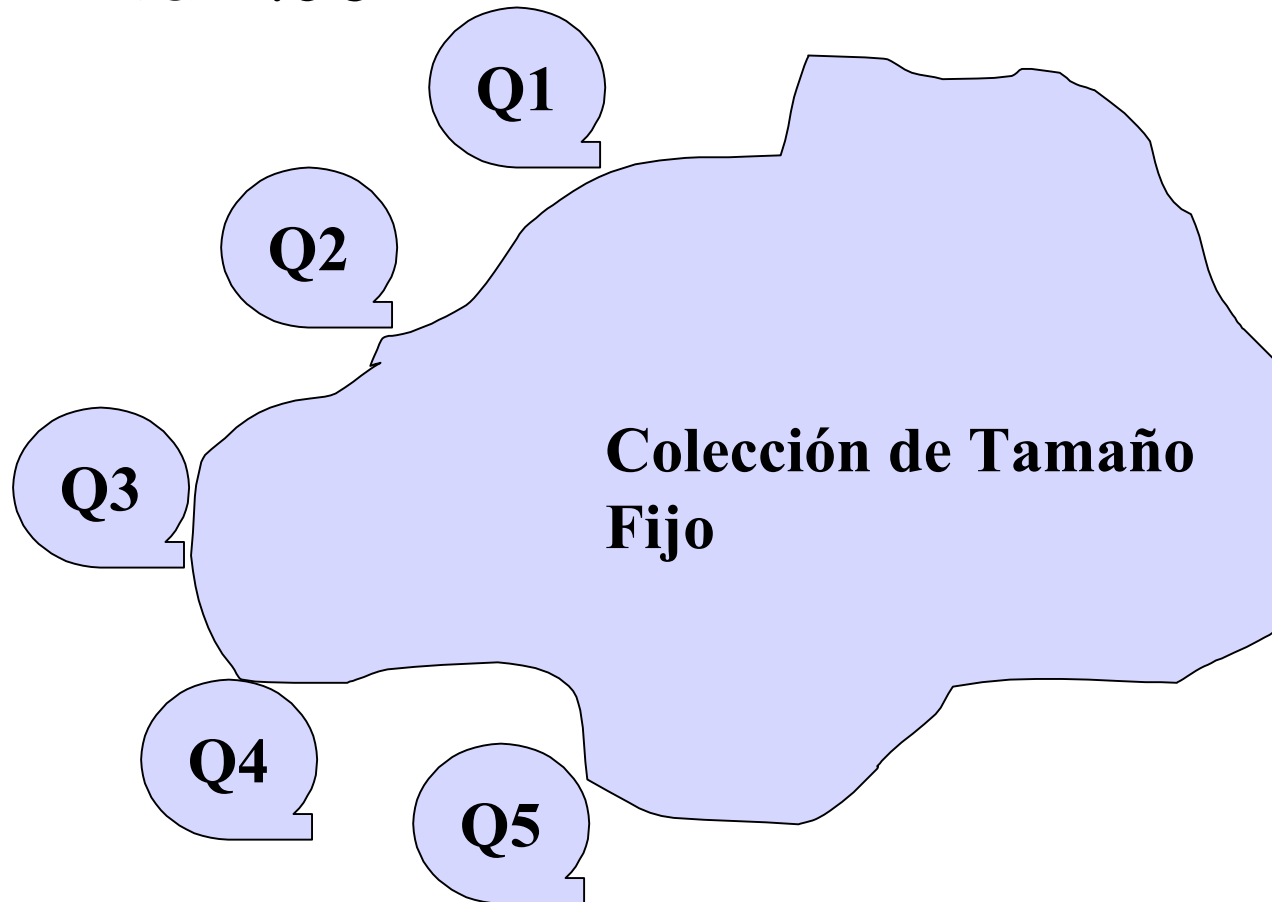
# Modelos IR

- Vista lógica versus tarea de recuperación

	<b>Index Terms</b>	<b>Full Text</b>	<b>Full Text + Structure</b>
<b>Retrieval</b>	Classic Set Theoretic Algebraic Probabilistic	Classic Set Theoretic Algebraic Probabilistic	Structured
<b>Browsing</b>	Flat	Flat Hypertext	Structure Guided Hypertext

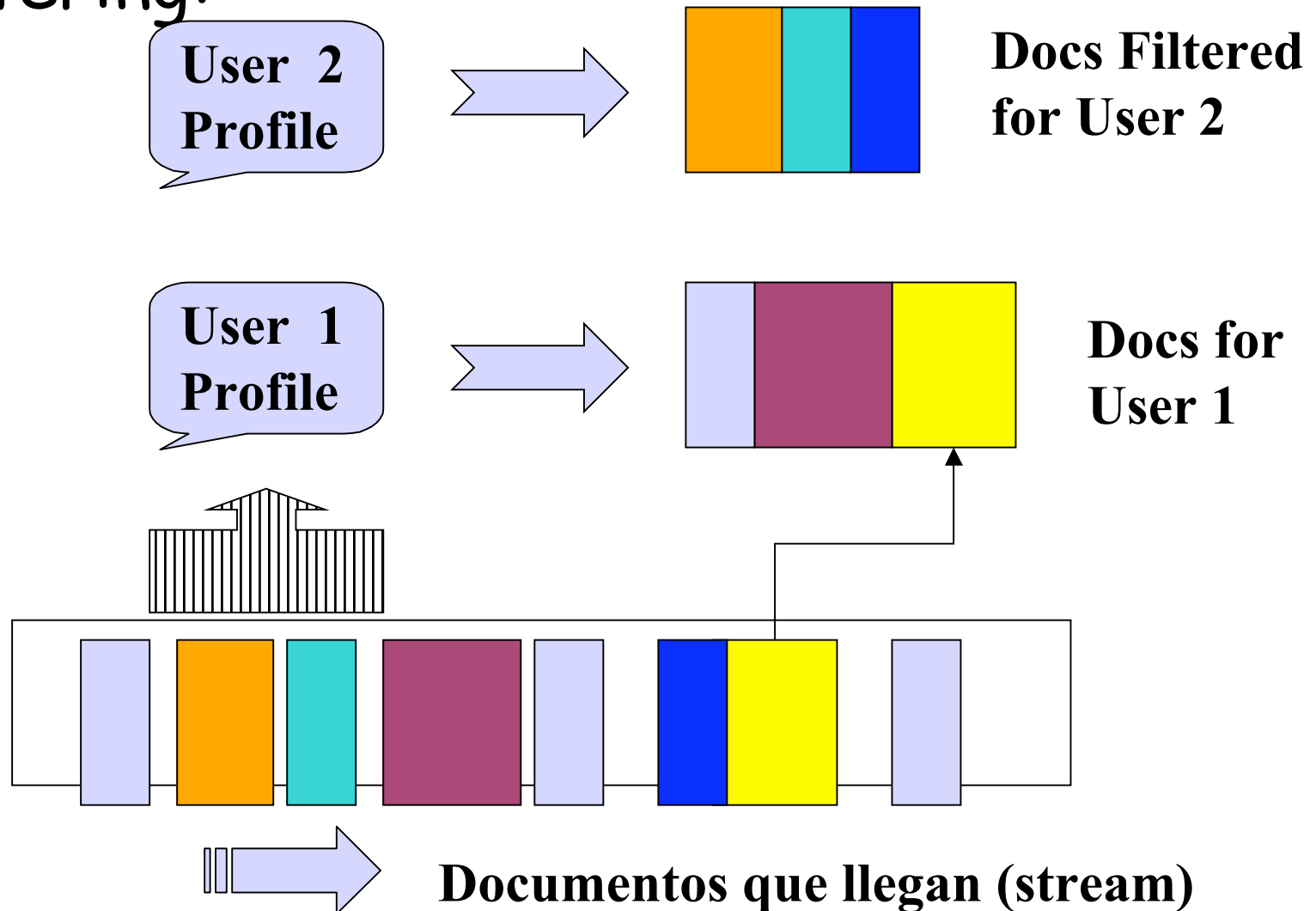
# Recuperación: Ad Hoc x Filtering

- Ad-hoc



# Recuperación: Ad Hoc x Filtering

- Filtering:



# Modelos Clásicos: Conceptos

- Cada documento es representado por un conjunto de palabras representativas o keywords.
- Un índice es una palabra en un documento que es útil para identificar el contenido del documento.
- Usualmente estos índices son sustantivos porque ellos tienen significado por sí solos.
- Sin embargo, máquinas de búsqueda asumen que todas las palabras son términos índices (representación de full text)

# Modelos Clásicos: Conceptos

- No todos los términos son igualmente útiles para representar un documento: términos menos frecuentes permiten identificar un conjunto más selecto de documentos.
- La *importancia* de los términos índices es representada por pesos asociados a ellos. Sean
  - $k_i$  un término índice
  - $d_j$  un documento
  - $w_{ij}$  es un peso asociado con  $(k_i, d_j)$
- El peso  $w_{ij}$  cuantifica la importancia del término índice para describir el contenido de documentos



# Modelos Clásicos: Conceptos

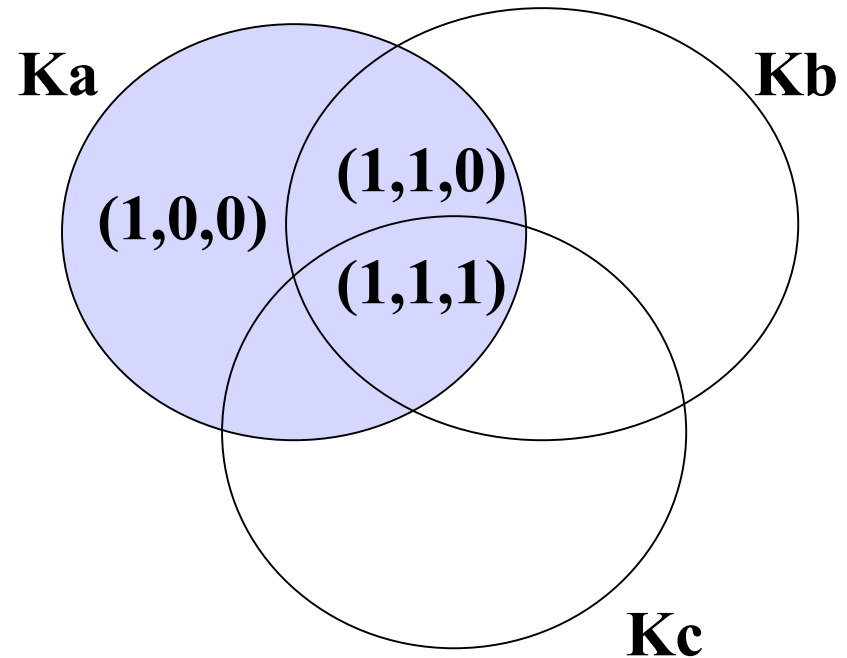
- $K_i$  es un término índice
- $d_j$  es un documento
- $t$  es el total de términos índices
- $K = (k_1, k_2, \dots, k_t)$  es el conjunto de términos índices
- $w_{ij} \geq 0$  es un peso asociado con  $(k_i, d_j)$
- $w_{ij} = 0$  indica que el término  $i$  no pertenece a un documento  $j$
- $\text{vec}(d_j) = (w_{1j}, w_{2j}, \dots, w_{tj})$  es un vector de pesos asociado a documento  $d_j$
- $g_i(\text{vec}(d_j)) = w_{ij}$  es una función la cual retorna el peso asociado al par  $(k_i, d_j)$

# Modelo Boolean

- Modelo Simple basado en teoría de conjunto
- Consultas expresadas como expresiones booleanas
  - precisa especificación
  - formalismo claro
  - $q = ka \quad (kb \vee kc)$
- Términos están o no presentes. Así,  $w_{ij} \in \{0,1\}$
- Considere
  - $q = ka \quad (kb \vee kc)$
  - $vec(q_{dnf}) = (1,1,1) \vee (1,1,0) \vee (1,0,0)$
  - $vec(q_{cc}) = (1,1,0)$  es un componente conjuntivo

# Modelo Booleano

- $q = ka \quad (kb \vee kc)$



- $sim(q,dj) = 1$  if  $\exists vec(q_{cc}) \mid (vec(q_{cc}) \quad vec(q_{dnf}))$   
(  $ki, gi(vec(dj)) = gi(vec(q_{cc}))$  )  
 $0$  otherwise

# Desventajas Modelo Boolean

- Recuperación basada en un criterio binario con ninguna noción de correspondencia parcial
- No ranking es producido
- La consulta debe ser traducidas a una expresión Boolean la cual no es cómoda para usuarios
- Las consultas son muchas veces demasiado simples
- Consecuentemente, el modelo puede retornan muchos o pocos documentos en respuesta a una consulta del usuario.

# Desventajas Modelo Boolean

- Recuperación basada en un criterio binario con ninguna noción de correspondencia parcial
- No ranking es producido
- La consulta debe ser traducidas a una expresión Boolean la cual no es cómoda para usuarios
- Las consultas son muchas veces demasiado simples
- Consecuentemente, el modelo puede retornan muchos o pocos documentos en respuesta a una consulta del usuario.

# Ejemplo

	k1	k2	k3	k4	k5	k6
d1	5	0	2	1	2	1
d2	2	5	2	0	0	3
d3	0	2	1	5	0	0
d4	3	0	1	5	1	0
d5	0	1	1	0	2	0

q: k1 ( k3  $\vee$   $\neg$ k6)

# Modelo Vectorial

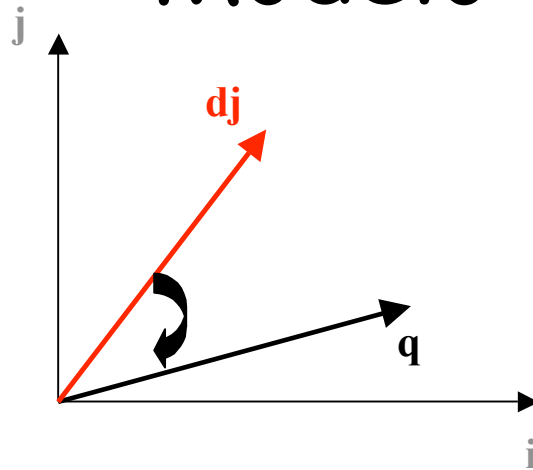
- El uso de pesos binarios es muy limitado
- Considera correspondencia parcial
- El peso de los términos es usado para calcular *el grado de similitud* entre consulta y cada documento.
- Conjunto ordenado (ranked) de documentos en base a mejor correspondencia con la consulta

# Modelo Vectorial

- Define:
  - $w_{ij} > 0$  cuando  $k_i$   $d_j$
  - $w_{iq} \geq 0$  asociado con el par  $(k_i, q)$
  - $vec(d_j) = (w_{1j}, w_{2j}, \dots, w_{tj})$   
 $vec(q) = (w_{1q}, w_{2q}, \dots, w_{tq})$
  - Por cada término  $k_i$  se asocia un vector unitario  $vec(i)$
  - Los vectores unitarios  $vec(i)$  y  $vec(j)$  son asumidos ortogonales (i.e., los índices son considerados independientes en un documento)
- El vector t-unario  $vec(i)$  forma una base ortogonal para el espacio t-dimensional
- En este espacio, consultas y documentos son representados como vectores de pesos



# Modelo Vectorial



- $Sim(q, dj) = \cos( )$   
 $= [vec(dj) \quad vec(q)] / |vec(dj)| \quad |vec(q)|$   
 $= [ \quad w_{ij} \quad w_{iq} ] / w_{ij}^2 \quad w_{iq}^2$
- Ya que  $w_{ij} \geq 0$  y  $w_{iq} \geq 0$ ,  $0 \leq sim(q, dj) \leq 1$
- Un documento es recuperado aunque corresponda parcialmente a una consulta.

# Modelo Vectorial

$$\text{Sim}(q,d_j) = [\text{vec}(d_j) \cdot \text{vec}(q)] / |\text{vec}(d_j)| \cdot |\text{vec}(q)|$$

- Como calcular los pesos  $w_{ij}$  y  $w_{iq}$ ?
- Un buen peso debe tomar en cuenta dos efectos:
  - cuantificación de intra-documentos (similitud)  
*factor tf*, la frecuencia del término en un documento
  - cuantificación inter-documentos (disimilitud)  
*factor idf*, la frecuencia inversa en los documentos

$$w_{ij} = \text{tf}(i,j) * \text{idf}(i)$$

# Modelo Vectorial

- Sean,
  - $N$  el número total de documentos
  - $n_i$  el número de documentos que contienen  $k_i$
  - $freq(i,j)$  frecuencia de  $k_i$  en  $d_j$
- Un factor normalizado de  $tf$  está dado por
  - $f(i,j) = freq(i,j) / \max(freq(l,j))$
  - donde el máximo es computado sobre todos los términos que ocurren en el documento  $d_j$
- El factor  $idf$  es calculado por
  - $idf(i) = \log(N/n_i)$
  - el  $\log$  es usado para hacer los valores de  $tf$  y  $idf$  comparables. Puede ser interpretado como la cantidad de información asociada con el término  $k_i$ .

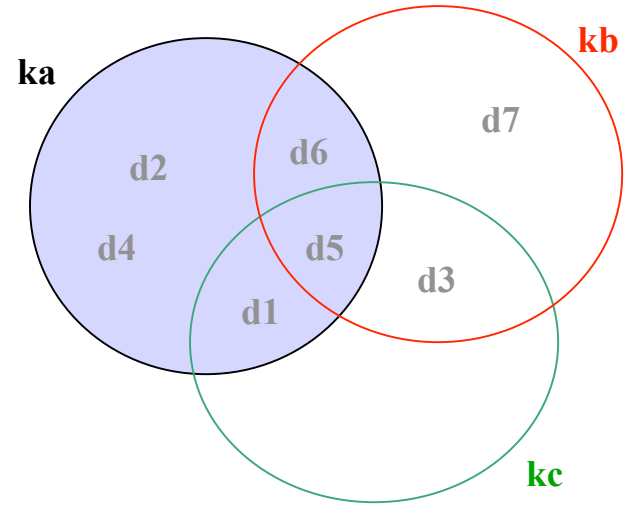
# Modelo Vectorial

- El esquema best term-weighting usa los pesos definidos por
  - $w_{ij} = f(i,j) * \log(N/n_i)$
  - esta estrategia es conocida como esquema de pesos *tf-idf*
- Para los términos de la consulta, se sugiere:
  - $w_{iq} = (0.5 + [0.5 * \text{freq}(i,q) / \max(\text{freq}(l,q))]) * \log(N/n_i)$
- El modelo vectorial con pesos *tf-idf* es una buena estrategia de ranking para una colección general
- El modelo vectorial es usualmente tan bueno como las alternativas de ranking conocidas. Es también simple y rápido de calcular.

# Modelo Vectorial

- Vectajas:
  - Peso de términos mejora la calidad del conjunto de respuesta
  - Correspondencia parcial permite la recuperación de documentos que aproximan las condiciones de la consulta.
  - Fórmula del ranking por coseno ordena los documentos de acuerdo al grado de similitud con la consulta
- Desventajas:
  - Asume independencia de los índices (algo aún incierto)

# Modelo Vectorial: Ejemplo I



$$q = ka \quad (kb \nabla kc)$$

# Modelo Probabilístico

- Objetivo: capturar el problema de IR usando un enfoque probabilístico
- Dada una consulta de un usuario, hay un conjunto respuesta ideal
- Consulta como especificación de las propiedades de este conjunto ideal de respuesta (clustering)
- Pero, cuáles son estas propiedades?
- Estima al comienzo cuáles pueden ser (e.i., una estimación inicial del conjunto respuesta ideal)
- Mejora por iteración

# Modelo Probabilístico

- Un conjunto inicial de documentos es recuperado de alguna manera
- El usuario inspecciona estos documentos buscando por los relevantes (en la práctica, sólo los top top 10-20 necesitan ser inspeccionados)
- Un sistema IR usa esta información para refinar la descripción del conjunto ideal
- Por repetición, se espera que el conjunto respuesta ideal vaya mejorando
- Siempre se debe pensar en que al comienzo se hace una estimación del conjunto ideal de respuesta
- La descripción del conjunto ideal es modelada en términos probabilísticos



# Principio del Ranking Probabilístico

- Dada una consulta de usuario  $q$  y un documento  $d_j$ , el modelo probabilístico trata de estimar la probabilidad que el usuario encontrará el documento  $d_j$  interesante (i.e., relevante). Este modelo asume que esta probabilidad de relevancia depende sólo de la consulta y del documento. El conjunto de respuesta ideal es denotado por  $R$  y debe maximizar la probabilidad de relevancia. Documentos en el conjunto  $R$  se dicen ser relevantes.
- Pero,
  - como determinar las probabilidades?
  - cuál es el espacio de muestreo?

# El Ranking

- El ranking probabilístico es calculado como:
  - $sim(q,dj) = P(R | vec(dj)) / P(\bar{R} | vec(dj))$
- Definición
  - $w_{ij} \in \{0,1\}$
  - $P(R | vec(dj))$ : probabilidad de que un documento dado es relevante
  - $P(\bar{R} | vec(dj))$ : probabilidad de un documento no es relevante

# El Ranking

- $\text{sim}(d_j, q) = P(R \mid \text{vec}(d_j)) / P(R \mid \text{vec}(d_j))$

$$= \frac{[P(\text{vec}(d_j) \mid R) * P(R)]}{[P(\text{vec}(d_j) \mid \bar{R}) * P(\bar{R})]}, \quad \text{Por Bayes's}$$

$$\sim \frac{P(\text{vec}(d_j) \mid R)}{P(\text{vec}(d_j) \mid \bar{R})}, [P(R), P(\bar{R})] \text{ para todos los docs.}$$

- $P(\text{vec}(d_j) \mid R)$  : probabilidad de aleatoriamente seleccionar el documento  $d_j$  desde el conjunto  $R$  de documentos relevantes

# El Ranking

- $\text{sim}(d_j, q) \sim \frac{P(\text{vec}(d_j) | R)}{P(\text{vec}(d_j) | R)}$   
 $\sim \frac{[ P(k_i | R) \quad P(k_i | R) ]}{[ P(k_i | R) ] [ P(k_i | R) ]}$

- $P(k_i | R)$  : probabilidad que el término índice  $k_i$  este presente en un documento aleatoriamente seleccionado desde el conjunto  $R$  de documentos relevantes

# El Ranking

- $\text{sim}(d_j, q) \sim \log \frac{[P(k_i | R)] [P(k_j | R)]}{[P(k_i | R)] [P(k_j | R)]}$

$$\sim w_{iq} w_{ij} \left( \log \frac{P(k_i | R)}{P(k_i | R)} + \log \frac{P(k_j | R)}{P(k_j | R)} \right)$$

donde  $P(k_i | R) = 1 - P(k_i | R)$

$$P(k_i | R) = 1 - P(k_i | R)$$

# El Ranking Inicial

- $\text{sim}(d_j, q) \sim$   
 $\sim w_{iq} * w_{ij} * \left( \log \frac{P(k_i | R)}{P(k_i | R)} + \log \frac{P(k_i | R)}{P(k_i | R)} \right)$

- Probabilidades  $P(k_i | R)$  y  $P(k_i | R)$ ?

- Estimaciones basadas en presunciones:

- $P(k_i | R) = 0.5$

- $P(k_i | R) = \frac{n_i}{N}$

donde  $n_i$  es el número de documentos que contienen  $k_i$

- Use esta estimación inicial para recuperar un ranking inicial

- Mejoras sobre el ranking inicial

# Mejorando el Ranking Inicial

- $\text{sim}(d_j, q) \sim w_{iq} * w_{ij} * (\log \frac{P(k_i | R)}{P(k_i | \bar{R})} + \log \frac{P(k_i | \bar{R})}{P(k_i | R)})$
- Sean
  - $V$  : conjunto de documentos inicialmente recuperados
  - $V_i$  : subconjunto de documentos recuperados que contienen  $k_i$
- Re-evaluar estimaciones :
  - $P(k_i | R) = \frac{V_i}{V}$
  - $P(k_i | \bar{R}) = \frac{n_i - V_i}{N - V}$
- Repetir recursivamente